

# Introduction to Queueing Theory and Stochastic Teletraffic Models

by Moshe Zukerman

Copyright M. Zukerman © 2000–2007

## Preface

The aim of this textbook is to provide students with basic knowledge of stochastic models that may be applied to telecommunications research areas such as traffic modelling, resource provisioning and traffic management. These study areas are often collectively called *teletraffic*. This book assumes prior knowledge of a programming language, mathematics, probability and stochastic processes normally taught in an electrical engineering course. For students who have some but not sufficiently strong background in probability and stochastic processes, we provide, in the first few chapters, a revision of the relevant concepts in these areas.

The book aims to enhance intuition and physical meaning understanding of the theoretical concepts it introduces. The famous mathematician Pierre-Simon Laplace is quoted to say that “Probability is common sense reduced to calculation” [10]; as the content of this book falls under the field of applied probability, Laplace’s quote very much applies. Accordingly, the book maintains clear linkage between common sense and the mathematical models and techniques it uses.

A unique feature of this book is the considerable attention given to guided projects involved computer simulations and analyzes. By successfully completing the programming assignments, students learn to simulate and analyze stochastic models such as queueing systems and by interpreting the results, they gain insight into the queueing performance effects and principles of telecommunications systems modelling. Although the book, at times, provides intuitive explanations, it still presents the important concepts and ideas required for the understanding teletraffic, queueing theory fundamentals and related queueing behavior of telecommunications networks and systems. These concepts and ideas form a strong base for the more mathematically inclined students who can follow up with the extensive literature on probability models and queueing theory. A small sample of it is listed at the end of this book.

As mentioned above, the first two chapters provide a revision of probability and stochastic processes topics relevant to the queueing and teletraffic models of this book. The content of these chapters is mainly based on [10, 20, 51, 53, 54, 55]. These chapters are intended for students who have some background in these topics. Students with no background in probability and stochastic processes are encouraged to study the original textbooks that include far more explanations, illustrations, discussions, examples and homework assignments. For students with background, we provide here a comprehensive summary of the key topics with relevant homework assignments that are especially tailored for understanding the queueing and

teletraffic models discussed in later chapters. Chapter 3 discusses general queueing notation and concepts and it should be studied well. Chapter 4 aims to assist the student to perform simulations of queueing systems. Simulations are useful and important in the many cases where exact analytical results are not available. An important learning objective of this book is to train students to perform queueing simulations. Chapter 5 provides analyses of deterministic queues. Many queueing theory books tend to exclude deterministic queues; however, the study of such queues is useful for beginners in that it helps them better understand non-deterministic queueing models. Chapters 6 – 12 provide analyses of a wide range of queueing and teletraffic models that fall under the category of continuous time Markov chain processes. Chapter 13 provides an example of a discrete time queue that is modelled as a discrete time Markov chain. In Chapter 14, various aspects of a single server queue with Poisson arrivals and general service times are studied, mainly focussing on mean value results as in [9]. Then, in Chapter 15, some selected results of a single server queue with a general arrival process and general service times are provided. Next, queueing networks are discussed in Chapter 16. Finally, in Chapter 17, stochastic processes that have been used as traffic models are discussed with special focus on their characteristics that affect queueing performance.

Throughout the book there is an emphasis on linking the theory with telecommunications applications as demonstrated by the following examples. Section 1.17 describes how properties of Gaussian distribution can be applied to link dimensioning. Section 6.4 shows, in the context of an M/M/1 queueing model, how optimally to set a link service rate such that delay requirements are met and how the level of multiplexing affects the spare capacity required to meet such delay requirement. An application of M/M/∞ queueing model to a multiple access performance problem [9] is discussed in Section 7.4. In Sections 8.4 and 9.3, discussions on dimensioning and related utilization issues of a multi-channel system are presented. Section 16.3 guides the reader to simulate a mobile cellular network. Section 17.6 describes a traffic model applicable to the Internet.

Last but not least, the author wish thank all the students and colleagues that provided comments and questions that helped developing and editing the manuscript over the years.

# Contents

|          |  |          |
|----------|--|----------|
| <b>1</b> | <b>Revision of Relevant Probability Topics</b>                         | <b>8</b> |
| 1.1      | Events, Sample Space, and Random Variables . . . . .                   | 8        |
| 1.2      | Probability, Conditional Probability and Independence . . . . .        | 8        |
| 1.3      | Probability and Distribution Functions . . . . .                       | 10       |
| 1.4      | Joint Distribution Functions . . . . .                                 | 10       |
| 1.5      | Conditional Probability for Random Variables . . . . .                 | 11       |
| 1.6      | Independence between Random Variables . . . . .                        | 11       |
| 1.7      | Convolution . . . . .  | 12       |
| 1.8      | Examples of Discrete Random Variables . . . . .                        | 12       |
| 1.8.1    | Bernoulli . . . . .  | 12       |
| 1.8.2    | Geometric . . . . .  | 13       |
| 1.8.3    | Binomial . . . . .   | 13       |
| 1.8.4    | Poisson . . . . .  | 14       |
| 1.8.5    | Pascal . . . . .   | 14       |
| 1.9      | Continuous Random Variables and their Probability Functions . . . . .  | 15       |
| 1.10     | Selected Continuous Random Variables . . . . .                         | 17       |
| 1.10.1   | Uniform . . . . .  | 17       |
| 1.10.2   | Exponential . . . . .  | 17       |
| 1.10.3   | Hyper-Exponential . . . . .  | 19       |
| 1.10.4   | Erlang . . . . .   | 19       |
| 1.10.5   | Hypo-Exponential . . . . .   | 19       |
| 1.10.6   | Gaussian . . . . .   | 19       |
| 1.10.7   | Pareto . . . . .   | 20       |
| 1.11     | Moments . . . . .  | 20       |
| 1.12     | Covariance and Correlation . . . . .                                   | 25       |
| 1.13     | Transforms . . . . .   | 27       |
| 1.13.1   | Z-transform . . . . .  | 30       |
| 1.13.2   | Laplace Transform . . . . .  | 32       |
| 1.14     | Multivariate Random Variables and Transform . . . . .                  | 34       |
| 1.15     | Probability Inequalities and Their Dimensioning Applications . . . . . | 34       |
| 1.16     | Limit Theorems . . . . .   | 36       |
| 1.17     | Link Dimensioning . . . . .  | 37       |

|          |   |           |
|----------|---|-----------|
| <b>2</b> | <b>Selected Stochastic Processes Topics</b>               | <b>39</b> |
| 2.1      | General Concepts . . . . .                                | 39        |
| 2.2      | Two Orderly and Memoryless Point Processes . . . . .      | 41        |
| 2.2.1    | Bernoulli Process . . . . .                               | 42        |
| 2.2.2    | Poisson Process . . . . .                                 | 44        |
| 2.3      | Markov Modulated Poisson Process . . . . .                | 48        |
| 2.4      | Discrete Time Markov Chains . . . . .                     | 48        |
| 2.4.1    | Definitions and Preliminaries . . . . .                   | 48        |
| 2.4.2    | Transition Probability Matrix . . . . .                   | 49        |
| 2.4.3    | Chapman-Kolmogorov Equation . . . . .                     | 49        |
| 2.4.4    | Marginal Probabilities . . . . .                          | 50        |
| 2.4.5    | Properties and Classification of States . . . . .         | 50        |
| 2.4.6    | Steady State Probabilities . . . . .                      | 52        |
| 2.4.7    | Reversibility . . . . .                                   | 54        |
| 2.4.8    | Multi-Dimensional Markov Chains . . . . .                 | 55        |
| 2.5      | Continuous Time Markov Chains . . . . .                   | 56        |
| 2.5.1    | Definitions and Preliminaries . . . . .                   | 56        |
| 2.5.2    | Examples . . . . .  | 57        |
| 2.5.3    | Birth and Death Process . . . . .                         | 57        |
| 2.5.4    | First Passage Time . . . . .                              | 58        |
| 2.5.5    | Transition Probability Function . . . . .                 | 58        |
| 2.5.6    | Steady State Probabilities . . . . .                      | 59        |
| 2.5.7    | Simulations . . . . .                                     | 61        |
| 2.5.8    | Reversibility . . . . .                                   | 61        |
| 2.5.9    | Multi-Dimensional Continuous Time Markov Chains . . . . . | 62        |
| <b>3</b> | <b>General Queueing Concepts</b>                          | <b>63</b> |
| 3.1      | Notation . . . . .  | 63        |
| 3.2      | Utilization . . . . .                                     | 63        |
| 3.3      | Little's Formula . . . . .                                | 64        |
| 3.4      | Work Conservation . . . . .                               | 66        |
| 3.5      | PASTA . . . . .   | 66        |
| 3.6      | Queueing Models . . . . .                                 | 67        |
| <b>4</b> | <b>Simulations</b>  | <b>68</b> |

|          |   |           |
|----------|---|-----------|
| 4.1      | Confidence Intervals . . . . .                        | 68        |
| 4.2      | Simulation of a G/G/1 Queue . . . . .                 | 68        |
| <b>5</b> | <b>Deterministic Queues</b>                           | <b>71</b> |
| 5.1      | D/D/1 . . . . .                                       | 71        |
| 5.2      | D/D/ $k$ . . . . .                                    | 72        |
| 5.3      | D/D/ $k/k$ . . . . .                                  | 73        |
| 5.4      | Summary of Results . . . . .                          | 74        |
| <b>6</b> | <b>M/M/1</b>  | <b>75</b> |
| 6.1      | Steady State Queue Size Probabilities . . . . .       | 75        |
| 6.2      | Delay Statistics . . . . .                            | 76        |
| 6.3      | Using Z-Transform . . . . .                           | 78        |
| 6.4      | Multiplexing . . . . .                                | 78        |
| 6.5      | The Departure Process . . . . .                       | 80        |
| 6.6      | Mean Busy Period and First Passage Time . . . . .     | 82        |
| 6.7      | A Markov Chain Simulation of M/M/1 . . . . .          | 83        |
| <b>7</b> | <b>M/M/<math>\infty</math></b>                        | <b>85</b> |
| 7.1      | Steady State Equations . . . . .                      | 85        |
| 7.2      | Solving the Steady State Equations . . . . .          | 85        |
| 7.3      | Insensitivity . . . . .                               | 86        |
| 7.4      | Applications . . . . .                                | 86        |
| <b>8</b> | <b>Erlang B Formula</b>                               | <b>88</b> |
| 8.1      | Solving the Steady State Equations . . . . .          | 88        |
| 8.2      | Recursion and Jagerman Formula . . . . .              | 88        |
| 8.3      | The Special Case: M/M/1/1 . . . . .                   | 89        |
| 8.4      | Dimensioning and Utilization . . . . .                | 90        |
| 8.5      | Insensitivity and Many Classes of Customers . . . . . | 91        |
| 8.6      | A Markov Chain Simulation of M/M/ $k/k$ . . . . .     | 93        |
| <b>9</b> | <b>M/M/<math>k</math></b>                             | <b>95</b> |
| 9.1      | Steady State Equations and Their Solution . . . . .   | 95        |
| 9.2      | Erlang C Formula . . . . .                            | 95        |
| 9.3      | Dimensioning and Utilization . . . . .                | 97        |

|  |            |
|--|------------|
| <b>10 Engset Loss Formula</b>  | <b>98</b>  |
| 10.1 Steady State Equations and Their Solution . . . . .                   | 98         |
| 10.2 Blocking Probability . . . . .  | 99         |
| 10.3 Obtaining the Blocking Probability by a Recursion . . . . .           | 100        |
| 10.4 Insensitivity . . . . .   | 100        |
| 10.5 Load Classifications and Definitions . . . . .                        | 100        |
| 10.6 The Many Sources Limit . . . . .                                      | 101        |
| 10.7 Obtaining the Blocking Probability by Successive Iterations . . . . . | 102        |
| <b>11 A Queue with State Dependent Arrivals and Service</b>                | <b>103</b> |
| 11.1 Steady State Queue Size Probabilities . . . . .                       | 103        |
| 11.2 Solving the Steady State Equations . . . . .                          | 104        |
| <b>12 Three Queueing Models with Finite Buffers</b>                        | <b>105</b> |
| 12.1 M/M/1/ $k$ . . . . .  | 105        |
| 12.2 MMPP(2)/M/1/ $k$ . . . . .  | 106        |
| 12.3 M/E <sub>n</sub> /1/ $k$ . . . . .                                    | 110        |
| <b>13 Discrete Time Queue</b>  | <b>112</b> |
| <b>14 M/G/1</b>  | <b>115</b> |
| 14.1 Pollaczek Khinchin Formula: Residual Service Approach [9] . . . . .   | 115        |
| 14.2 Pollaczek Khinchin Formula: by Kendall's Recursion [33] . . . . .     | 117        |
| 14.3 Special Cases: M/M/1 and M/D/1 . . . . .                              | 118        |
| 14.4 Busy Period . . . . .   | 118        |
| 14.5 M/G/1 with Priorities . . . . .                                       | 119        |
| 14.6 Nonpreemptive . . . . .   | 119        |
| 14.7 Preemptive Resume . . . . .   | 120        |
| <b>15 G/G/1</b>  | <b>122</b> |
| 15.1 Reich's Formula . . . . .   | 122        |
| 15.2 The G/GI/1 Queue and Its G/GI/1/ $k$ Equivalent . . . . .             | 123        |
| <b>16 Queueing Networks</b>  | <b>126</b> |
| 16.1 Jackson Networks . . . . .  | 126        |
| 16.2 Erlang Fixed-Point Approximation . . . . .                            | 129        |

|   |            |
|---|------------|
| 16.3 A Markov Chain Simulation of a Mobile Cellular Network . . . . . | 130        |
| <b>17 Stochastic Processes as Traffic Models</b>                      | <b>133</b> |
| 17.1 Parameter Fitting . . . . .                                      | 133        |
| 17.2 Poisson Process . . . . .  | 134        |
| 17.3 Markov Modulated Poisson Process (MMPP) . . . . .                | 134        |
| 17.4 Autoregressive Gaussian Process . . . . .                        | 134        |
| 17.5 Exponential Autoregressive (1) Process . . . . .                 | 135        |
| 17.6 Poisson Pareto Burst Process . . . . .                           | 136        |

# 1 Revision of Relevant Probability Topics

Probability theory provides the foundation for queueing theory and stochastic teletraffic models, therefore it is important that the student masters the probability concepts required for the material that follows. We aim to provide in this chapter sufficient coverage for readers that have some probability background. Although the cover here is comprehensive in the sense that it discusses all the probability concepts and techniques used in later chapters, it does not include the many examples and exercises that are normally included in a probability textbook to help readers grasp the material better. Therefore, readers without prior probability background may be aided by additional probability texts such as [10] and [54].

## 1.1 Events, Sample Space, and Random Variables

Consider an experiment with an uncertain outcome. The term “experiment” refers to any uncertain scenario such as tomorrow’s weather, tomorrow’s share price of a certain company, or the result of flipping a coin. The *sample space* is a set of all possible outcomes of an experiment. An *event* is a subset of the sample space. Consider, for example, an experiment which consists of tossing a die. The sample space is  $\{1, 2, 3, 4, 5, 6\}$ , and an event could be the set  $\{2, 3\}$ , or  $\{6\}$ , or the empty set  $\{\}$  or even the entire sample space  $\{1, 2, 3, 4, 5, 6\}$ . Events are called *mutually exclusive* if their intersection is the empty set. A set of events is said to be *exhaustive* if its union is equal to the sample space.

A *random variable* is a real valued function defined on the sample space. This definition appears somewhat contradictory to the wording “random variable” as a random variable is not at all random, because it is actually a deterministic function which assigns a real valued number to each possible outcome of an experiment. It is the outcome of the experiment that is random and therefore the name: random variable. If we consider the flipping a coin experiment, the possible outcomes are Head (H) and Tail (T), hence the sample space is  $\{H, T\}$ , and a random variable  $X$  could assign  $X = 1$  for the outcome H, and  $X = 0$  for the outcome T.

If  $X$  is a random variable than  $Y = g(X)$  for some function  $g(\cdot)$  is also a random variable. In particular, some functions of interest are  $Y = cX$  for some constant  $c$  and  $Y = X^n$  for some integer  $n$ .

If  $X_1, X_2, X_3, \dots, X_n$  is a sequence of random variables, than  $Y = \sum_{i=1}^n X_i$  is also a random variable.

## 1.2 Probability, Conditional Probability and Independence

Consider a sample space  $S$ . Let  $A$  be a set in  $S$ , the probability of  $A$  is the function on  $S$ , denoted  $P(A)$ , that satisfies the following three axioms:

1.  $0 \leq P(A) \leq 1$
2.  $P(S) = 1$
3. The probability of the union of mutually exclusive events is equal to the sum of the probabilities of these events.



Normally, higher probability signifies higher likelihood of occurrence. In particular, if we conduct a very large number of experiments, and we generate the *histogram* by measuring how many times each of the possible occurrences actually occurred. Then we normalized the histogram by dividing all its values by the total number of experiments to obtain the relative frequencies. These measurable relative frequencies are represented by the theoretical concept of probability.

We use the notation  $P(A \mid B)$  for the *conditional probability* of  $A$  given  $B$ , which is the probability of the event  $A$  given that we know that event  $B$  has occurred. If we know that  $B$  has occurred, it is our new sample space, and for  $A$  to occur, the relevant experiments outcomes must be in  $A \cap B$ , hence the new probability of  $A$ , namely the probability  $P(A \mid B)$ , is the ratio between the probability of  $A \cap B$  and the probability of  $B$ . Accordingly,

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}. \quad (1)$$

**Remark:**

The intersection of  $A$  and  $B$  is also denoted by  $A, B$  or  $AB$  in addition to  $A \cap B$ .

If events  $A$  and  $B$  are *independent*, which means that if one of them occurs, the probability of the other to occur is not affected, then

$$P(A \mid B) = P(A) \quad (2)$$

and hence, by Eq. (1), if  $A$  and  $B$  are independent then,

$$P(A \cap B) = P(A)P(B). \quad (3)$$

Let  $B_1, B_2, B_3, \dots, B_n$  be a sequence of mutually exclusive and exhaustive events in  $S$ , and let  $A$  be another event in  $S$ . Then,

$$A = \bigcup_{i=1}^n (A \cap B_i) \quad (4)$$

and since the  $B_i$ s are mutually exclusive, the events  $A \cap B_i$ s are also mutually exclusive. Hence,

$$P(A) = \sum_{i=1}^n P(A \cap B_i). \quad (5)$$

Thus, by Eq. (1),

$$P(A) = \sum_{i=1}^n P(A \mid B_i) \times P(B_i). \quad (6)$$

The latter is a very useful formula for deriving probability of a given event by conditioning and unconditioning on a set of mutually exclusive and exhaustive events. It is called *the law of total probability*. Therefore, by Eqs. (6) and (1) (again), we obtain the following formula for conditional probability between two events:

$$P(B_1 \mid A) = \frac{P(A \mid B_1)P(B_1)}{\sum_{i=1}^n P(A \mid B_i) \times P(B_i)}. \quad (7)$$

The latter is known as Bayes' formula.

### 1.3 Probability and Distribution Functions

Random variables are related to events. When we say that random variable  $X$  takes value  $x$ , this means that  $x$  represents a certain outcome of an experiment which is an event, so  $\{X = x\}$  is an event. Therefore, we may assign probabilities to all possible values of the random variable. This function denoted  $P_X(x) = P(X = x)$  will henceforth be called *probability function*. The *distribution function* of random variable  $X$  is defined for all  $x \in R$  ( $R$  being the set of all real numbers), is defined as

$$F_X(x) = P(X \leq x). \quad (8)$$

Accordingly, the *complementary distribution function*  $\bar{F}_X(x)$  is defined by

$$\bar{F}_X(x) = P(X > x). \quad (9)$$

Consequently, for any random variable, for every  $x \in R$ ,  $F(x) + \bar{F}(x) = 1$ .

### 1.4 Joint Distribution Functions

In some cases, we are interested in the probability that two or more random variables are within a certain range. For this purpose, we define, *the joint distribution function* for  $n$  random variables  $X_1, X_2, \dots, X_n$ , as follows:

$$F_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n). \quad (10)$$

Having the joint distribution function, we can obtain the distribution function of a single random variable, say,  $X_1$ , as

$$F_{X_1}(x_1) = F_{X_1, X_2, \dots, X_n}(x_1, \infty, \dots, \infty). \quad (11)$$

When the random variables  $X_1, X_2, \dots, X_n$  are discrete, we can use their *joint probability function* which is defined by

$$P_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n). \quad (12)$$

The probability function of a single random variable can then be obtained by

$$P_{X_1}(x_1) = \sum_{x_2} \cdots \sum_{x_n} P_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n). \quad (13)$$

A random variable is called *discrete* if it takes at most a countable number of possible values. On the other hand, a *continuous* random variable takes an uncountable number of possible values. In this section and in sections 1.5, 1.6, 1.7, when we mention random variables or their probability and distribution function, we consider them all to be discrete. Then in Section 1.9, we will introduce the analogous definitions and notation relevant to their continuous counterparts.

## 1.5 Conditional Probability for Random Variables

The conditional probability concept, which we defined for events, can also apply to random variables. Because  $\{X = x\}$ , namely, the random variable  $X$  takes value  $x$ , is an event, by the definition of conditional probability (1) we have

$$P(X = x \mid Y = y) = \frac{P(X = x, Y = y)}{P(Y = y)}. \quad (14)$$

Let  $P_{X|Y}(x \mid y) = P(X = x \mid Y = y)$  be the conditional probability of a discrete random variable  $X$  given  $Y$ , we obtain by (14)

$$P_{X|Y}(x \mid y) = \frac{P_{X,Y}(x, y)}{P_Y(y)}. \quad (15)$$

Noticing that

$$P_Y(y) = \sum_x P_{X,Y}(x, y), \quad (16)$$

we obtain by (15)

$$\sum_x P_{X|Y}(x \mid y) = 1. \quad (17)$$

This means that if we condition on the event  $\{Y = y\}$  for a specific  $y$ , the probability function of  $X$  given  $\{Y = y\}$  is a legitimate probability function. This is consistent with our discussion above. The event  $\{Y = y\}$  is the new sample space and  $X$  has a legitimate probability and distribution functions there. By (15)

$$P_{X,Y}(x, y) = P_{X|Y}(x \mid y)P_Y(y) \quad (18)$$

and by symmetry

$$P_{X,Y}(x, y) = P_{Y|X}(y \mid x)P_X(x) \quad (19)$$

so the latter and (16) gives

$$P_Y(y) = \sum_x P_{X,Y}(x, y) = \sum_x P_{Y|X}(y \mid x)P_X(x) \quad (20)$$

which is another version of the law of total probability (6).

## 1.6 Independence between Random Variables

The definition of independence between random variables is very much related to the definition of independence between events because when we say that random variables  $U$  and  $V$  are independent, it is equivalent to say that the events  $\{U = u\}$  and  $\{V = v\}$  are independent for every  $u$  and  $v$ . Accordingly, random variables  $U$  and  $V$  are said to be independent if

$$P_{U,V}(u, v) = P_U(u)P_V(v) \quad \text{for all } u, v. \quad (21)$$

Notice that by (19) and (21), we obtain an equivalent definition of independent random variables  $U$  and  $V$  which is

$$P_{U|V}(u \mid v) = P_U(u) \quad (22)$$

which is equivalent to  $P(A \mid B) = P(A)$  which we used to define independent events  $A$  and  $B$ .

## 1.7 Convolution

Consider independent random variables  $U$  and  $V$  that have probability functions  $P_U(u)$  and  $P_V(v)$ , respectively, and their sum which is another random variable  $X = U + V$ . Let us now derive the probability function  $P_X(x)$  of  $X$ .

$$\begin{aligned} P_X(x) &= P(U + V = x) \\ &= \sum_u P(U = u, V = x - u) \\ &= \sum_u P_U(u)P_V(x - u). \end{aligned}$$

The latter is called the *convolution* of the probability functions  $P_U(u)$  and  $P_V(v)$ .

Let us now extend the result from two to  $k$  random variables. Consider  $k$  independent random variables  $X_i$ ,  $i = 1, 2, 3, \dots, k$ . Let  $P_{X_i}(x_i)$  be the probability function of  $X_i$ , for  $i = 1, 2, 3, \dots, k$ , and let  $Y = \sum_{i=1}^k X_i$ . If  $k = 3$ , we first compute the convolution of  $X_1$  and  $X_2$  to obtain the probability function of  $U = X_1 + X_2$  using the above convolution formula and then we use the formula again to obtain the probability function of  $Y = U + X_3 = X_1 + X_2 + X_3$ . Therefore, for an arbitrary  $k$ , we obtain

$$P_Y(y) = \sum_{x_2, x_3, \dots, x_k: x_2, x_3, \dots, x_k \leq y} \left( P_{X_1}(y - \sum_{i=2}^k x_i) \prod_{i=2}^k P_{X_i}(x_i) \right). \quad (23)$$

If all the random variable  $X_i$ ,  $i = 1, 2, 3, \dots, k$ , are independent and identically distributed (IID) random variables, with probability function  $P_{X_1}(x)$  then the probability function  $P_Y(y)$  is called the  $k$ -fold convolution of  $P_{X_1}(x)$ .

## 1.8 Examples of Discrete Random Variables

We now consider several discrete random variables and their corresponding probability distributions. As mentioned above, our cover here is not exhaustive. There are many important discrete random variables which are not discussed here.

### 1.8.1 Bernoulli

We begin with the Bernoulli random variable. It represents an outcome of an experiment which has only two possible outcomes. Let us call them “success” and “failure”. These two outcomes are mutually exclusive and exhaustive events. The Bernoulli random variable assigns the value  $X = 1$  to the “success” outcome and the value  $X = 0$  to the “failure” outcome. Let  $p$  be the probability of the “success” outcome, and because “success” and “failure” are mutually exclusive and exhaustive, the probability of the “failure” outcome is  $1 - p$ . The probability function in terms of the Bernoulli random variable is:

$$\begin{aligned} P(X = 1) &= p \\ P(X = 0) &= 1 - p. \end{aligned} \quad (24)$$

### 1.8.2 Geometric

The geometric random variable  $X$  represents the number of independent Bernoulli trials, each of which with  $p$  being the probability of success, required until the first success. For  $X$  to be equal to  $i$  we must have  $i - 1$  failures and then one success in  $k$  independent Bernoulli trials. Therefore, we obtain

$$P(X = i) = (1 - p)^{i-1}p \quad \text{for } i = 1, 2, 3, \dots \quad (25)$$

### 1.8.3 Binomial

Assume that  $n$  independent Bernoulli trials are performed. Let  $X$  be a random variable representing the number of successes in these  $n$  trials. Such random variable is called a binomial random variable with parameters  $n$  and  $p$ . Its probability function is:

$$P(X = i) = \binom{n}{i} p^i (1 - p)^{n-i} \quad i = 0, 1, 2, \dots, n.$$

Notice that a Binomial random variable with parameters 1 and  $p$  is a Bernoulli random variable. The Bernoulli and binomial random variables have many applications. In particular, it is used as a model for voice as well as data sources. Such sources alternates between two states “on” and “off”. During the “on” state the source is active and transmits at the rate equal to the transmission rate of its equipment (e.g. a modem), and during the “off” state, the source is idle. If  $p$  is the proportion of time that the source is active, and if we consider a superposition of  $n$  independent identical sources, then the binomial distribution gives us the probability of the number of sources which are simultaneously active which is important for resource provisioning.

## Homework 1.1

Consider a state with voter population  $N$ . There are two candidates in the state election for governor and the winner is chosen based on a simple majority. Let  $N_1$  and  $N_2$  be the total number of votes obtained by candidates 1 and 2, respectively, from voters other than Johnny. Johnny just voted for candidate 1, and he would like to know the probability that his vote affects the election results, namely,  $0 \geq N_1 - N_2 \geq -1$ . Assume that each other voter (excluding Johnny) votes independently for candidates 1 and 2 with probabilities  $p_1$  and  $p_2$ , respectively, and also that  $p_1 + p_2 < 1$  to allow for the case that a voter chooses not to vote for either candidate. Derive a formula for the probability that Johnny’s vote affects the election results and provide an algorithm and a computer program to compute it for the case  $N = 2,000,000$  and  $p_1 = p_2 = 0.4$ .

## Guide

By the definition of conditional probability,

$$P(N_1 = n_1, N_2 = n_2) = P(N_1 = n_1)P(N_2 = n_2 \mid N_1 = n_1)$$

so

$$P(N_1 = n_1, N_2 = n_2) = \binom{N-1}{n_1} p_1^{n_1} (1-p_1)^{N-n_1-1} \binom{N-n_1-1}{n_2} p_2^{n_2} (1-p_2)^{N-n_1-1-n_2}.$$

Then the required probability is

$$\sum_{k=0}^{\lfloor (N-1)/2 \rfloor} P(N_1 = k, N_2 = k) + \sum_{k=0}^{\lfloor (N-1)/2 \rfloor - 1} P(N_1 = k, N_2 = k+1).$$

□

### 1.8.4 Poisson

A Poisson random variable with parameter  $\lambda$  has the following probability function:

$$P(X = i) = e^{-\lambda} \frac{\lambda^i}{i!} \quad i = 0, 1, 2, 3, \dots \quad (26)$$

The importance of the Poisson random variable lies in its property to approximate the binomial random variable in case when  $n$  is very large and  $p$  is very small so that  $np$  is not too large and not too small. It will be shown rigorously using Z-transform in Subsection 1.13.1 that if  $n$  increases and  $np$  stays constant, a binomial random variable with parameters  $n$  and  $p$  approaches Poisson with parameter  $\lambda = np$ . The Poisson random variable accurately models the number of calls arriving at a telephone exchange or Internet service provider in a short period of time, a few seconds or a minute, say. In this case, the population of customers (or packets)  $n$  is large. The probability  $p$  of a customer (or packet) making a call within a given short period of time is small, and the calls are typically independent. Therefore, models based on Poisson random variables have been used successfully for design and dimensioning of telecommunications networks and systems for many years. When we refer to items in a queueing system in this book, they will be called customers, jobs or packets interchangeably.

### Homework 1.2

Consider a Poisson random variable  $X$  with parameter  $\lambda = 500$ . Write a program that computes the probabilities  $P(X = i)$  for  $0 \leq i \leq 800$  and plot the function  $P_X(x)$ . □

### 1.8.5 Pascal

The Pascal random variable  $X$  with parameters  $k$  (integer and  $\geq 1$ ) and  $p$  (real within  $(0,1]$ ), represents a sum of  $k$  geometric random variables each with parameter  $p$ . For  $X$  to be equal to  $i$ , we must have a successful Bernoulli trial at the  $i$ th trial because this is the successful trial associated with the  $k$ th geometric random variable. Then there must also be exactly  $k-1$  successes among the first  $i-1$  trials. The probability to have a success at the  $i$ th trial is equal to the probability of a Bernoulli random variable with parameter  $p$  equal to 1, which is  $p$ , and the probability of having  $k-1$  successes among the first  $i-1$  is equal to the probability of

having a binomial random variable with parameters  $p$  and  $i - 1$  equal to  $k - 1$ , for  $i \geq k \geq 1$ , which is equal to

$$\binom{i-1}{k-1} p^{k-1} (1-p)^{i-k} \quad k = 1, 2, \dots, i$$

and since the two random variables here, namely, the Bernoulli and the Binomial are independent (because the underlying Bernoulli trials are independent), we can multiply their probabilities to obtain

$$P(X = 1) = \binom{i-1}{k-1} p^k (1-p)^{i-k} \quad i = k, k+1, k+2, \dots \quad (27)$$

## 1.9 Continuous Random Variables and their Probability Functions

Continuous random variable are related to cases whereby the set of possible outcomes is uncountable. A continuous random variable  $X$  is a function that assigns a real number to outcome of an experiment, and is characterized by the existence of a function  $f(\cdot)$  defined for all  $x \in R$ , which has the property that for any set  $A \subset R$ ,

$$P(X \in A) = \int_A f(x) dx. \quad (28)$$

Such function is the *probability density function* (or simply the *density*) of  $X$ . Since the continuous random variable  $X$  must take a value in  $R$  with probability 1,  $f$  must satisfy,

$$\int_{-\infty}^{+\infty} f(x) dx = 1. \quad (29)$$

If we consider Eq. (28), letting  $A = [a, b]$ , we obtain,

$$P(a \leq x \leq b) = \int_a^b f(x) dx. \quad (30)$$

An interesting point to notice is that the probability of a continuous random variable taking a particular value is equal to zero. If we set  $a = b$  in Eq. (30), we obtain

$$P(x = a) = \int_a^a f(x) dx = 0. \quad (31)$$

As a result, the distribution function  $F(x)$  is equal to both  $P(X \leq x)$  and to  $P(X < x)$ . Similarly, the complementary distribution function is equal to both  $P(X \geq x)$  and to  $P(X > x)$ .

By Eq. (30), we obtain

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(s) ds. \quad (32)$$

Hence, the probability density function is the derivative of the distribution function.

An important concept which gives rise to a continuous version of the law of total probability is the continuous equivalence of Eq. (12), namely, the joint distribution of continuous random variables. Let  $X$  and  $Y$  be two continuous random variables. The joint density of  $X$  and  $Y$  denoted  $f_{X,Y}(x, y)$  is a nonnegative function that satisfies

$$P(\{X, Y\} \in A) = \iint_{\{X, Y\} \in A} f_{X,Y}(x, y) dx dy. \quad (33)$$

Another important concept is the conditional density of one continuous random variable on another. Let  $X$  and  $Y$  be two continuous random variables with joint density  $f_{X,Y}(x, y)$ . For any  $y$ , such that the density of  $Y$  takes a positive value at  $Y = y$  (i.e. such that  $f_Y(y) > 0$ ), the conditional density of  $X$  given  $Y$  is defined as

$$f_{X|Y}(x | y) = \frac{f_{X,Y}(x, y)}{f_Y(y)}. \quad (34)$$

For every given fixed  $y$ , it is a legitimate density because

$$\int_{-\infty}^{\infty} f_{X|Y}(x | y) dx = \int_{-\infty}^{\infty} \frac{f_{X,Y}(x, y)}{f_Y(y)} dx = \frac{f_Y(y)}{f_Y(y)} = 1. \quad (35)$$

Notice the equivalence between the conditional probability (1) and the conditional density (34). By (34)

$$f_{X,Y}(x, y) = f_Y(y) f_{X|Y}(x | y) \quad (36)$$

so

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy = \int_{-\infty}^{\infty} f_Y(y) f_{X|Y}(x | y) dy. \quad (37)$$

Recall again that  $f_{X,Y}(x, y)$  is defined only for  $y$  values such that  $f_Y(y) > 0$ .

Let define event  $A$  as  $X \in A$ . Thus,

$$P(A) = P(X \in A) = \int_A f_X(x) dx = \int_A \int_{-\infty}^{\infty} f_Y(y) f_{X|Y}(x | y) dy dx. \quad (38)$$

Hence,

$$P(A) = \int_{-\infty}^{\infty} f_Y(y) \int_A f_{X|Y}(x | y) dx dy \quad (39)$$

and therefore

$$P(A) = \int_{-\infty}^{\infty} f_Y(y) P(A | Y = y) dy \quad (40)$$

which is the continuous equivalence of the Law of Total Probability (6).

We will now discuss the concept of convolution as applied to continuous random variables. Consider independent random variables  $U$  and  $V$  that have densities  $f_U(u)$  and  $f_V(v)$ , respectively, and their sum which is another random variable  $X = U + V$ . Let us now derive the density  $f_X(x)$  of  $X$ .

$$\begin{aligned} f_X(x) &= P(U + V = x) \\ &= \int_u f(U = u, V = x - u) \\ &= \int_u f_U(u) f_V(x - u). \end{aligned}$$

The latter is the *convolution* of the densities  $f_U(u)$  and  $f_V(v)$ .

As in the discrete case the convolution  $f_Y(y)$ , of  $k$  densities  $f_{X_i}(x_i)$ ,  $i = 1, 2, 3, \dots, k$ , of random variables  $X_i$ ,  $i = 1, 2, 3, \dots, k$ , respectively, is given by

$$f_Y(y) = \iint_{x_2, x_3, \dots, x_k: x_2 + x_3 + \dots + x_k \leq y} \left( f_{X_1}(y - \sum_{i=2}^k x_i) \prod_{i=2}^k f_{X_i}(x_i) \right). \quad (41)$$

And again, in the special case where all the random variable  $X_i$ ,  $i = 1, 2, 3, \dots, k$ , are IID, the density  $f_Y$  is the  $k$ -fold convolution of  $f_{X_1}$ .



## 1.10 Selected Continuous Random Variables

We will now discuss several continuous random variables and their corresponding probability distributions: uniform, exponential, hyper-exponential, Erlang, hypo-exponential Gaussian, multivariate Gaussian and Pareto. These are selected because of their applicability in teletraffic and related queueing models and consequently their relevance to the material in this book.

### 1.10.1 Uniform

The probability density function of the uniform random variable takes nonnegative values over the interval  $(a, b)$  and is given by

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise.} \end{cases} \quad (42)$$

Of particular interest is the special case - the uniform  $(0,1)$  random variable. Its probability density function is given by

$$f(x) = \begin{cases} 1 & \text{if } 0 < x < 1 \\ 0 & \text{otherwise.} \end{cases} \quad (43)$$

The uniform  $(0,1)$  random variable is very important in simulations. Almost all computers programs have a function which generates uniform  $(0,1)$  random deviates. By a simple transformation such uniform  $(0,1)$  random deviates can be translated to sequence of random deviates of any distribution as follows. Let  $U_1(0,1)$  be the first uniform  $(0,1)$  random deviate, and let  $F(x)$  be a distribution function of an arbitrary random variable. Set,

$$U_1(0,1) = F(x_1) \quad (44)$$

so  $x_1 = F^{-1}(U_1(0,1))$  is the first  $F(\cdot)$  random deviate. Then generating the second uniform  $(0,1)$  random deviate, the second  $F(\cdot)$  random number is obtained in the same way, etc.

To see why this method works, let  $U$  be a uniform  $(0,1)$  random variable. Let  $F(x)$  be an arbitrary cumulative distribution function. Let the random variable  $Y$  be defined by:  $Y = F^{-1}(U)$ . That is,  $U = F(Y)$ . We will now show that the distribution of  $Y$ , namely  $P(Y \leq x)$ , is equal to  $F(x)$ . Notice that  $P(Y \leq x) = P[F^{-1}(U) \leq x] = P[U \leq F(x)]$ . Because  $U$  is a uniform  $(0,1)$  random variable, then  $= P[U \leq F(x)] = F(x)$ . Thus,  $P(Y \leq x) = F(x)$ .  $\square$

### 1.10.2 Exponential

The exponential random variable has one parameter  $\mu$  and its probability density function is given by,

$$f(x) = \begin{cases} \mu e^{-\mu x} & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases} \quad (45)$$

Its distribution function is given by

$$F(x) = \int_0^x \mu e^{-\mu s} ds = 1 - e^{-\mu x} \quad x \geq 0. \quad (46)$$

The most useful way to describe the exponential random variable is by its complementary distribution function. It is given by,

$$\bar{F}(x) = e^{-\mu x} \quad x \geq 0. \quad (47)$$

An important application of the exponential random variable is the time until the next call (or connection request) arrives at a switch. Interestingly, such time does not depend on how long ago was the last call that arrived. This property is called the *memoryless* property of a random variable. In particular, a random variable is called memoryless if

$$P(X > s + t \mid X > t) = P(X > s). \quad (48)$$

If our lifetime were memoryless, then the probability we survive at least 80 years given that we have survived 70 years is equal to the probability that a newborn baby lives to be 10 years. Of course human lifetime is not memoryless, but, as mentioned above, inter-arrivals of phone calls at a telephone exchange are. To show that exponential random variable is memoryless we show that Eq. (48) holds using the conditional probability definition together with the complementary distribution function of an exponential random variable as follows.

$$\begin{aligned} P(X > s + t \mid X > t) &= \frac{P(X > s + t \cap X > t)}{P(X > t)} \\ &= \frac{P(X > s + t)}{P(X > t)} \\ &= \frac{e^{-\mu(s+t)}}{e^{-\mu t}} \\ &= e^{-\mu s} = P(X > s). \end{aligned}$$

Not only is exponential random variable memoryless, it is actually, the only memoryless continuous random variable.

### Homework 1.3

Write a computer program that generates a sequence of 100 random deviates from an exponential distribution with  $\mu = 1$ .  $\square$

Let  $X_1$  and  $X_2$  be exponentially distributed random variables with parameters  $\lambda_1$  and  $\lambda_2$ . We are interested to know the distribution of  $X = \min[X_1, X_2]$ . In other words, we are interested in the distribution of the time that passes until the first one of the two random variables  $X_1$  and  $X_2$  occurs. This is as if we have a competition between the two and we are interested in the time of the winner whichever it is. Then

$$P(X > t) = P(\min[X_1, X_2] > t) = P(X_1 > t, X_2 > t) = e^{-\lambda_1 t} e^{-\lambda_2 t} = e^{-(\lambda_1 + \lambda_2)t}. \quad (49)$$

Thus, the distribution of  $T$  is exponential with parameter  $\lambda_1 + \lambda_2$ .

Another interesting question related to the competition between two exponential random variables is what is the probability that one of them, say  $X_1$  wins. That is, we are interested in the probability of  $X_1 < X_2$ . This is obtained using the continuous version of the law of total probability (40) as follows:

$$P(X_1 < X_2) = \int_0^\infty (1 - e^{-\lambda_1 t}) \lambda_2 e^{-\lambda_2 t} dt = \frac{\lambda_1}{\lambda_1 + \lambda_2}. \quad (50)$$

Equivalently,

$$P(X_1 > X_2) = \frac{\lambda_2}{\lambda_1 + \lambda_2}. \quad (51)$$

As expected,  $P(X_1 < X_2) + P(X_1 > X_2) = 1$ . Notice that as  $X_1$  and  $X_2$  are continuous time random variables, the probability that they are equal to each other is equal to zero.

### 1.10.3 Hyper-Exponential

Let  $X_i$  for  $i = 1, 2, 3, \dots, k$  be  $k$  independent exponential random variables with parameters  $\lambda_i$ ,  $i = 1, 2, 3, \dots, k$ , respectively. Let  $p_i$  for  $i = 1, 2, 3, \dots, k$  be  $k$  nonnegative real numbers such that  $\sum_{i=1}^k p_i = 1$ . A random variable  $X$  that is equal to  $X_i$  with probability  $p_i$  is called Hyper-exponential. By the Law of total probability, its density is

$$f_X(x) = \sum_{i=1}^k p_i f_{X_i}(x). \quad (52)$$

### 1.10.4 Erlang

A random variable  $X$  has Erlang distribution with parameters  $\lambda$  (positive real) and  $k$  (positive integer) if its density is given by

$$f_X(x) = \frac{\lambda^k x^{k-1} e^{-\lambda x}}{(k-1)!}. \quad (53)$$

### Homework 1.4

Let  $X_i$ ,  $i = 1, 2, \dots, k$  be  $k$  independent exponentially distributed random variables each with parameter  $\lambda$ , prove by induction that the random variable  $X$  defined by the sum  $X = \sum_{i=1}^k X_i$  has Erlang distribution with parameter  $k$  and  $\lambda$ . In other words,  $f_X(x)$  of (53) is a  $k$ -fold convolution of  $\lambda e^{-\lambda x}$ .  $\square$

### 1.10.5 Hypo-Exponential

Let  $X_i$ ,  $i = 1, 2, \dots, k$  be  $k$  independent exponentially distributed random variables each with parameters  $\lambda_i$ , respectively. The random variable  $X$  defined by the sum  $X = \sum_{i=1}^k X_i$  is called hypo-exponential. In other words, the density of  $X$  is a convolution of the  $k$  densities  $\lambda_i e^{-\lambda_i x}$ ,  $i = 1, 2, \dots, k$ . The Erlang distribution is a special case of hypo-exponential when all the  $k$  random variables are identically distributed.

### 1.10.6 Gaussian

A continuous random variable, which commonly used in many applications, is the Gaussian (also called Normal) random variable. We say that the random variable  $X$  has Gaussian distribution with parameters  $m$  and  $\sigma^2$  if its density is given by

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m)^2/2\sigma^2} \quad -\infty < x < \infty. \quad (54)$$

This density is symmetric and bell shaped.

The wide use of the Gaussian random variable is rooted in the so-called **The Central Limit Theorem**. This theorem is the most important result in probability theory. Loosely speaking, it says that the sum of a large number of independent random variables (not necessarily of the same distribution, but each has a finite variance) has Gaussian (normal) distribution. This is also true if the distribution of these random variables are very different from Gaussian. This theorem explains why so many populations in nature and society have bell shaped Gaussian histograms, and justifies the use of the Gaussian distribution as their model.

### 1.10.7 Pareto

Another continuous random variable often used in telecommunication modelling is the **Pareto** random variable. This random variable, for a certain parameter range, it can be useful in modelling lengths of data bursts in data and multimedia networks [1]. We choose to define the Pareto random variable with parameters  $\gamma$  and  $\delta$  by its complementary distribution function which is given by

$$P(X > x) = \begin{cases} \left(\frac{x}{\delta}\right)^{-\gamma}, & x \geq \delta \\ 1, & \text{otherwise.} \end{cases}$$

### Homework 1.5

Write a computer program that generates a sequence of 100 random deviates from a Pareto distribution with  $\gamma = 1.2$  and  $\delta = 4$ .  $\square$

## 1.11 Moments

The *mean* (or the expectation) of a discrete random variable is defined by

$$E[X] = \sum_{\{n: P(n) > 0\}} n P_X(n). \quad (55)$$

Equivalently, the mean of a continuous random variable is defined by

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx. \quad (56)$$

A very useful expression for the mean of a continuous nonnegative random variable  $Z$  (i.e. a random variable  $Z$  with the property that its density  $f(z) = 0$  for  $z < 0$ ) is:

$$E[Z] = \int_0^{\infty} P(Z > z) dz = \int_0^{\infty} [1 - F_Z(z)] dz. \quad (57)$$

The discrete equivalence of the latter is:

$$E[Z] = \sum_{n=0}^{\infty} P(Z > n) = \sum_{n=0}^{\infty} [1 - F_Z(n)]. \quad (58)$$

### Homework 1.6

Use geometrical arguments to show (57) and (58).  $\square$

As mentioned above, function of a random variable is also a random variable. The mean of a function of random variables denoted  $g(\cdot)$  by

$$E[g(X)] = \sum_{\{k:P_X(k)>0\}} g(k)P_X(k) \quad (59)$$

for a discrete random variable and

$$E[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x)dx \quad (60)$$

for a continuous random variable. If  $a$  and  $b$  are constants then for a random variable  $X$  (either discrete or continuous) we have:

$$E[aX] = aE[X], \quad (61)$$

$$E[X - b] = E[X] - b, \quad (62)$$

and

$$E[aX - b] = aE[X] - b. \quad (63)$$

The  **$n$ th moment** of the random variable  $X$  is defined by  $E[X^n]$ . Substituting  $g(X) = X^n$  in (59) and in (60), the  $n$ th moment of  $X$  is given by:

$$E[X^n] = \sum_{\{k:P_X(k)>0\}} k^n P_X(k) \quad (64)$$

for a discrete random variable and

$$E[X^n] = \int_{-\infty}^{\infty} x^n f_X(x)dx \quad (65)$$

for a continuous random variable. Similarly, the  $n$ th central moment of random variable  $X$  is defined by  $E[(X - E[X])^n]$ . Substituting  $g(X) = (X - E[X])^n$  in (59) and in (60), the  **$n$ th central moment** of  $X$  is given by:

$$E[(X - E[X])^n] = \sum_{\{k:P(k)>0\}} (k - E[X])^n P_X(k) \quad (66)$$

for a discrete random variable and

$$E[(X - E[X])^n] = \int_{-\infty}^{\infty} (x - E[X])^n f_X(x)dx \quad (67)$$

for a continuous random variable. By definition the first moment is the mean. The second central moment is called the **variance**. It is defined as

$$var[X] = E[(X - E[X])^2]. \quad (68)$$

The variance of a random variable  $X$  is given by

$$var[X] = \sum_{\{k:P(k)>0\}} (k - E[X])^2 P_X(k) \quad (69)$$

if  $X$  is discrete, and by

$$\text{var}[X] = \int_{-\infty}^{\infty} (x - E[X])^2 f_X(x) dx \quad (70)$$

if it is continuous.

By (68) we obtain

$$\text{var}[X] = E[(X - E[X])^2] = E[X^2 - 2XE[X] + (E[X])^2] = E[X^2] - (E[X])^2. \quad (71)$$

In the following table we provide the mean and the variance of some of the above described random variables.

| random variable | parameters                        | mean                              | variance       |
|-----------------|-----------------------------------|-----------------------------------|----------------|
| Bernoulli       | $0 \leq p \leq 1$                 | $p$                               | $p(1 - p)$     |
| binomial        | $n$ and $0 \leq p \leq 1$         | $np$                              | $np(1 - p)$    |
| Poisson         | $\lambda > 0$                     | $\lambda$                         | $\lambda$      |
| uniform         | $a$ and $b$                       | $(a + b)/2$                       | $(b - a)^2/12$ |
| exponential     | $\mu > 0$                         | $1/\mu$                           | $1/\mu^2$      |
| Gaussian        | $m$ and $\sigma$                  | $m$                               | $\sigma^2$     |
| Pareto          | $\delta > 0$ and $1 < \gamma < 2$ | $\frac{\delta\gamma}{(\gamma-1)}$ | $\infty$       |

Notice that since the binomial random variable is a sum of  $n$  independent Bernoulli random variables, its mean and its variance are  $n$  times the mean and variance, respectively, of the Bernoulli random variable. Notice also that by letting  $p \rightarrow 0$ , and  $np \rightarrow \lambda$ , both the mean and the variance of the binomial random variable approach  $\lambda$ , which is the value of both the mean and variance of the Poisson random variable.

While the mean provides the average, or the average of possible values a random variable can take weighted according to its probability function or density, the variance is a measure of the level of variation of the possible values of the random variable. Another measure of such variation is the **standard deviation** denoted  $\sigma_X$ , or simply  $\sigma$ , and defined by

$$\sigma_X = \sqrt{\text{var}[X]}. \quad (72)$$

Hence the variance is often denoted by  $\sigma^2$ .

Notice that the first central moment  $E[x - E[X]]$  is not very useful because it is always equal to zero, the second central moment  $E[(x - E[X])^2]$ , which is the variance, and its square root, the standard deviation, are used for measuring the level of variation of a random variable.

The mean of sum of random variables is always the sum of their means, namely,

$$E\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n E[X_i] \quad (73)$$

but the variance of sum of random variables is not always equal to the sum of their variances. It is true for independent random variables. That is, if the random variables  $X_1, X_2, X_3, \dots, X_n$  are independent, then

$$\text{var}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \text{var}[X_i]. \quad (74)$$

Also, if  $X_1, X_2, X_3, \dots, X_n$  are independent, then

$$E[\Pi_{i=1}^n X_i] = \Pi_{i=1}^n E[X_i]. \quad (75)$$

In many application it is useful to use the concept of **Conditional Expectation (or Mean)** to derive moments of unknown distributions. It is defined by:

$$E[X | Y] = E_X[P(X | Y)], \quad (76)$$

where the subscript  $X$  indicates that the mean is over  $X$ . For example, the conditional expectation of two discrete random variables is defined by

$$E[X | Y = j] = \sum_i iP(X = i | Y = j). \quad (77)$$

If  $X$  and  $Y$  are continuous, their conditional expectation is defined as

$$E[X | Y = y] = \int_{x=-\infty}^{\infty} xf_{X|Y}(x | y)dx. \quad (78)$$

It is important to realize that  $E[X | Y]$  is a random variable which is a function of the random variable  $Y$ . Therefore, if we consider its mean (in the case that  $Y$  is discrete) we obtain

$$\begin{aligned} E_Y[E[X | Y]] &= \sum_j E[X | Y = j]P(Y = j) \\ &= \sum_j \sum_i iP(X = i | Y = j)P(Y = j) \\ &= \sum_i i \sum_j P(X = i | Y = j)P(Y = j) \\ &= \sum_i iP(X = i) = E[X]. \end{aligned}$$

Thus, we have obtained the following formula for the mean  $E[X]$

$$E[X] = E_Y[E[X | Y]]. \quad (79)$$

The latter also applies to continuous random variables. In this case we have:

$$\begin{aligned} E_Y[E[X | Y]] &= \int_{y=-\infty}^{\infty} E[X | Y = y]f_Y(y)dy \\ &= \int_{y=-\infty}^{\infty} \int_{x=-\infty}^{\infty} xf_{X|Y}(x | y)dx f_Y(y)dy \\ &= \int_{x=-\infty}^{\infty} x \int_{y=-\infty}^{\infty} f_{X|Y}(x | y)f_Y(y)dy \\ &= \int_{x=-\infty}^{\infty} xf_X(x)dx = E[X]. \end{aligned}$$

## Homework 1.7

Show that  $E[X] = E_Y[E[X | Y]]$  holds also for the case where  $X$  is discrete and  $Y$  is continuous and vice versa.  $\square$

Note that  $P(X = x | Y = y)$  is itself a random variable that is a function of the values  $y$  taken by random variable  $Y$ . Therefore, by definition  $E_Y[P(X = x | Y = y)] = \sum_y P(X = x | Y = y)P(Y = y)$  which lead to another way to express the Law of Total Probability:

$$P_X(x) = E_Y[P(X = x | Y = y)]. \quad (80)$$

Define the **Conditional Variance** as

$$\text{var}[X | Y] = E[(X - E[X | Y])^2 | Y]. \quad (81)$$

This gives rise to the following useful formula for the variance of a random variable known as EVVE:

$$\text{var}[X] = E[\text{var}[X | Y]] + \text{var}[E[X | Y]]. \quad (82)$$

To show EVVE, we recall (71):  $\text{var}[X] = E[X^2] - (E[X])^2$ , and (79):  $E[X] = E_Y[E[X | Y]]$ , we obtain

$$\text{var}[X] = E[E[X^2 | Y]] - (E[E[X | Y]])^2. \quad (83)$$

Then using  $E[X^2] = \text{var}[X] + (E[X])^2$  gives

$$\text{var}[X] = E[\text{var}[X | Y] + (E[X | Y])^2] - (E[E[X | Y]])^2 \quad (84)$$

or

$$\text{var}[X] = E[\text{var}[X | Y]] + E[E[X | Y]]^2 - (E[E[X | Y]])^2. \quad (85)$$

Now considering again the formula  $\text{var}[X] = E[X^2] - (E[X])^2$ , but instead of the random variable  $X$  we put the random variable  $E[X | Y]$ , we obtain

$$\text{var}[E[X | Y]] = E[E[X | Y]]^2 - (E[E[X | Y]])^2, \quad (86)$$

observing that the right-hand side of (86) equals to the last two terms in the right-hand side of (85), we obtain EVVE.

To illustrate the use of conditional mean and variance, consider the following example. Every second the number of Internet flows that arrive at a router, denoted  $\phi$ , has mean  $\phi_e$  and variance  $\phi_v$ . The number of packets in each flow, denoted  $\varsigma$ , has mean  $\varsigma_e$  and variance  $\varsigma_v$ . Assume that the number of packets in each flow and the number of flows arriving per second are independent. Let  $W$  the total number of packets arriving at the router per second which has mean  $W_e$  and variance  $W_v$ . Assume  $W = \varsigma\phi$ . The network designer, aiming to meet certain quality of service (QoS) requirements, makes sure that the router serves the arriving packets at the rate of  $s_r$  per second, such that  $s_r = W_e + 4\sqrt{W_v}$ . To compute  $s_r$  one needs to have the values of  $W_e$  and  $W_v$ . Because  $\phi$  and  $\varsigma$  are independent  $E[W|\phi] = \phi\varsigma_e$  and by (79)

$$W_e = E[W] = E[E[W|\phi]] = E[\phi]E[\varsigma] = \phi_e\varsigma_e.$$

Note that the relationship

$$W_e = \phi_e\varsigma_e \quad (87)$$

is also obtained directly by (75). In fact, the above proves (75) for the case of two random variables.

Also by EVVE,

$$\text{Var}[W] = E[\text{Var}[W|\phi]] + \text{Var}[E[W|\phi]] = \varsigma_v E[\phi^2] + (\varsigma_e)^2 \text{Var}[\phi].$$

Therefore

$$W_v = \phi_v\varsigma_v + \varsigma_v\phi_e^2 + \phi_v\varsigma_e^2. \quad (88)$$



**Homework 1.8**

1. Provide detailed derivations of Equations (87) and (88) using (79) and (82).
2. Derive Equations (87) and (88) in a different way, considering the independence of the number of packets in each flow and the number of flows arriving per second.  $\square$

**1.12 Covariance and Correlation**

When random variables are positively dependent, namely, if when one of them obtains high values, the others are likely to obtain high value also, then the variance of their sum may be much higher than the sum of the individual variances. This is very significant for bursty multimedia traffic modeling and resource provisioning. For example, let time be divided into consecutive small time intervals, if  $X_i$  is the amount of traffic arrives during the  $i$ th interval, and assume that we use a buffer that can store traffic that arrives in many intervals, the probability of buffer overflow will be significantly affected by the variance of the amount of traffic arrives in a time period of many intervals, which in turn is strongly affected by the dependence between the  $X_i$ s. Therefore, there is a need to define a quantitative measure for dependence between random variables. Such measure is called the **covariance**. The covariance of two random variables  $X_1$  and  $X_2$ , denoted by  $cov(X_1, X_2)$ , is defined by

$$cov(X_1, X_2) = E[(X_1 - E[X_1])(X_2 - E[X_2])]. \quad (89)$$

Intuitively, by Eq. (89), if high value of  $X_1$  implies high value of  $X_2$ , and low value of  $X_1$  implies low value of  $X_2$ , the covariance is high. By Eq. (89),

$$cov(X_1, X_2) = E[X_1 X_2] - E[X_1]E[X_2]. \quad (90)$$

Hence, by (75), if  $X_1$  and  $X_2$  are independent then  $cov(X_1, X_2) = 0$ . The variance of the sum of two random variables  $X_1$  and  $X_2$  is given by

$$var[X_1 + X_2] = var[X_1] + var[X_2] + 2cov(X_1, X_2). \quad (91)$$

This is consistent with our comments above. The higher the dependence between the two random variables, as measured by their covariance, the higher the variance of their sum, and if they are independence, hence  $cov(X_1, X_2) = 0$ , the variance of their sum is equal to the sum of their variances. Notice that the reverse is not always true:  $cov(X_1, X_2) = 0$  does not necessarily imply that  $X_1$  and  $X_2$  are independent.

Notice also that negative covariance results in lower value for the variance of their sum than the sum of the individual variances.

**Homework 1.9**

Prove that  $cov(X_1, X_2) = 0$  does not necessarily imply that  $X_1$  and  $X_2$  are independent.

**Guide**

The proof is by a counter example. Consider two random variables  $X$  and  $Y$  and assume that both have Bernoulli distribution with parameter  $p$ . Consider random variable  $X_1$  defined by  $X_1 = X+Y$  and another random variable  $X_2$  defined by  $X_2 = X-Y$ ; show that  $cov(X_1, X_2) = 0$  and that  $X_1$  and  $X_2$  are not independent.  $\square$

Let the sum of the random variables  $X_1, X_2, X_3, \dots, X_k$  be denoted by

$$S_k = X_1 + X_2 + X_3 + \dots + X_k.$$

Then

$$var(S_k) = \sum_{i=1}^k var[X_i] + 2 \sum_{i < j} cov[X_i, X_j] \quad (92)$$

where  $\sum_{i < j} cov[X_i, X_j]$  is a sum over all  $cov[X_i, X_j]$  such that  $i$  and  $j$  is a pair selected without repetitions out of  $1, 2, 3, \dots, k$  so that  $i < j$ .

**Homework 1.10**

Prove Eq. (92).

**Guide**

First show that  $S_k - E[S_k] = \sum_{i=1}^k (X_i - E[X_i])$  and that

$$(S_k - E[S_k])^2 = \sum_{i=1}^k (X_i - E[X_i])^2 + 2 \sum_{i < j} (X_i - E[X_i])(X_j - E[X_j]).$$

Then take expectations of both sides of the latter.  $\square$

If we consider  $k$  independent random variables denoted  $X_1, X_2, X_3, \dots, X_k$ , then by substituting  $cov[X_i, X_j] = 0$  for all relevant  $i$  and  $j$  in (92), we obtain

$$var(S_k) = \sum_{i=1}^k var[X_i]. \quad (93)$$

**Homework 1.11**

Use Eq. (92) to explain the relationship between the variance of a Bernoulli random variable and a binomial random variable.

**Guide**

Notice that a binomial random variable with parameter  $k$  and  $p$  is a sum of  $k$  independent Bernoulli random variables with parameter  $p$ .  $\square$

The covariance can take any value between  $-\infty$  and  $+\infty$ , and in some cases, it is convenient to have a normalized dependence measure - a measure that takes values between -1 and 1. Such measure is the **correlation**. Noticing that the covariance is bounded by

$$\text{cov}(X_1, X_2) \leq \sqrt{\text{var}[X_1]\text{var}[X_2]}, \quad (94)$$

the correlation of two random variables  $X$  and  $Y$  denoted by  $\text{corr}(X, Y)$  is defined by

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}, \quad (95)$$

assuming  $\text{var}[X] \neq 0$  and  $\text{var}[Y] \neq 0$ .

### Homework 1.12

Prove that  $|\text{corr}(X, Y)| \leq 1$ .

### Guide

Let  $C = \text{corr}(X, Y)$ , and show  $C^2 - \sigma_X^2 \sigma_Y^2 \leq 0$ , by noticing that  $C^2 - \sigma_X^2 \sigma_Y^2$  is a discriminant of the quadratic  $a^2 \sigma_X^2 + 2aC + \sigma_Y^2$  which must be nonnegative because  $E[a(X - E[X]) + (Y - E[Y])]^2$  is nonnegative.  $\square$

## 1.13 Transforms

Transforms are very useful in analysis of probability models and queueing systems. We will first consider the following general definition [10] for a transform function  $\Gamma$  of a random variable  $X$ :

$$\Gamma_X(\omega) = E[e^{\omega X}] \quad (96)$$

where  $\omega$  is a complex scalar. Transforms have two important properties:

1. There is a one-to-one correspondence between transforms and probability distributions. This is why they are sometimes called *characteristics* functions. This means that for any distribution function there is a unique transform function that characterizes it and for each transform function there is a unique probability distribution it characterizes. Unfortunately it is not always easy to convert a transform to its probability distribution, and therefore we in some cases that we are able to obtain the transform but not its probability function, we use it as means to characterize the random variable statistics instead of the probability distribution.
2. Having a transform function of a random variable we can generate its moments. This is why transforms are sometimes called *moment generating* functions. In many cases, it is easier to obtain the moments having the transform than having the actual probability distribution.

We will now show how to obtain the moments of a continuous random variable  $X$  with density function  $f_X(x)$  from its transform function  $\Gamma_X(\omega)$ . By definition,

$$\Gamma_X(\omega) = \int_{-\infty}^{\infty} e^{\omega x} f_X(x) dx. \quad (97)$$

Taking derivative with respect to  $\omega$  leads to

$$\Gamma'_X(\omega) = \int_{-\infty}^{\infty} x e^{\omega x} f_X(x) dx. \quad (98)$$

Letting  $\omega \rightarrow 0$ , we obtain

$$\lim_{\omega \rightarrow 0} \Gamma'_X(\omega) = E[X], \quad (99)$$

and in general, taking the  $n$ th derivative and letting  $\omega \rightarrow 0$ , we obtain

$$\lim_{\omega \rightarrow 0} \Gamma_X^{(n)}(\omega) = E[X^n]. \quad (100)$$

### Homework 1.13

Derive Eq. (100) using (97) – (100) completing all the missing steps.  $\square$

Consider for example the exponential random variable  $X$  with parameter  $\lambda$  having density function  $f_X(x) = \lambda e^{-\lambda x}$  and derive its transform function. By definition,

$$\Gamma_X(\omega) = E[e^{\omega X}] = \lambda \int_{x=0}^{\infty} e^{\omega x} e^{-\lambda x} dx, \quad (101)$$

which gives after some derivations

$$\Gamma_X(\omega) = \frac{\lambda}{\lambda - \omega}. \quad (102)$$

### Homework 1.14

Derive Eq. (102) from (101)  $\square$

Let  $X$  and  $Y$  be random variables and assume that  $Y = aX + b$ . The transform of  $Y$  is given by

$$\Gamma_Y(\omega) = E[e^{\omega Y}] = E[e^{\omega(aX+b)}] = e^{\omega b} E[e^{\omega a X}] = e^{\omega b} \Gamma_X(\omega a). \quad (103)$$

Let random variable  $Y$  be the sum of independent random variables  $X_1$  and  $X_2$ , i.e.,  $Y = X_1 + X_2$ . The transform of  $Y$  is given by

$$\Gamma_Y(\omega) = E[e^{\omega Y}] = E[e^{\omega(X_1+X_2)}] = E[e^{\omega X_1}] E[e^{\omega X_2}] = \Gamma_{X_1}(\omega) \Gamma_{X_2}(\omega). \quad (104)$$

This result applies to a sum of  $n$  independent random variables, so the transform of a sum of independent random variable equals to the product of their transform. If  $Y = \sum_{i=1}^n X_i$  and all the  $X_i$ s are  $n$  independent and identically distributed (IID) random variables, then

$$\Gamma_Y(\omega) = E[e^{\omega Y}] = [\Gamma_{X_1}(\omega)]^n. \quad (105)$$

Let us now consider a Gaussian random variable  $X$  with parameters  $m$  and  $\sigma$  and density

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m)^2/2\sigma^2} \quad -\infty < x < \infty. \quad (106)$$

Its transform is derived as follows

$$\begin{aligned} \Gamma_X(\omega) &= E[e^{\omega X}] \\ &= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m)^2/2\sigma^2} e^{\omega x} dx \\ &= e^{(\sigma^2 \omega^2/2) + m\omega} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m)^2/2\sigma^2} e^{\omega x} e^{-(\sigma^2 \omega^2/2) - m\omega} dx \\ &= e^{(\sigma^2 \omega^2/2) + m\omega} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m-\sigma^2 \omega)^2/2\sigma^2} dx \\ &= e^{(\sigma^2 \omega^2/2) + m\omega}. \end{aligned}$$

Let us use the transform just derived to obtain the mean and variance of a Gaussian random variable with parameters  $m$  and  $\sigma$ . Taking the first derivative and putting  $\omega = 0$ , we obtain

$$E[X] = \Gamma'_X(0) = m. \quad (107)$$

Taking the second derivative and setting  $\omega = 0$ , we obtain

$$E[X^2] = \Gamma_X^{(2)}(0) = \sigma^2 + m^2. \quad (108)$$

Thus,

$$\text{var}[X] = E[X^2] - E[X]^2 = \sigma^2 + m^2 - m^2 = \sigma^2. \quad (109)$$

A Gaussian random variable with mean equal to zero and variance equal to one is called *standard Gaussian*. It is well known that if  $Y$  is Gaussian with mean  $m$  and standard deviation  $\sigma$ , then the random variable  $X$  defined as

$$X = \frac{Y - m}{\sigma} \quad (110)$$

is standard Gaussian.

Substituting  $\sigma = 1$  and  $m = 0$  in the above transform of a Gaussian random variable, we obtain that

$$\Gamma_X(\omega) = e^{(\omega^2/2)} \quad (111)$$

is the transform of a standard Gaussian random variable.

### Homework 1.15

Show the consistency between the results obtained for transform of a Gaussian random variable, (103), (110) and (111).  $\square$

Let  $X_i$ ,  $i = 1, 2, 3, \dots, n$  be  $n$  independent random variables and let  $Y$  be a random variable that equals  $X_i$  with probability  $p_i$  for  $i = 1, 2, 3, \dots, N$ . Therefore, by the Law of Total Probability,

$$P(Y = y) = \sum_{i=1}^N p_i P(X_i = y) \quad (112)$$

or for continuous densities

$$f_Y(y) = \sum_{i=1}^n p_i f_{X_i}(y). \quad (113)$$

Its transform is given by

$$\begin{aligned} \Gamma_Y(\omega) &= E[e^{\omega Y}] \\ &= \int_{-\infty}^{\infty} f_Y(y) e^{\omega y} \\ &= \int_{-\infty}^{\infty} \left[ \sum_{i=1}^n p_i f_{X_i}(y) \right] e^{\omega y} \\ &= \int_{-\infty}^{\infty} \sum_{i=1}^n p_i f_{X_i}(y) e^{\omega y} \\ &= \sum_{i=1}^n p_i \Gamma_{X_i}(\omega). \end{aligned}$$

Notice that if the  $X_i$  are exponential random variables then, by definition,  $Y$  is hyper-exponential. Particular transforms include the Z, the Laplace, and the Fourier transforms.

The **Z-transform**  $\Pi_X(z)$  applies to integer valued random variable  $X$  and is defined by

$$\Pi_X(z) = E[z^X].$$

This is a special case of (96) by setting  $z = e^\omega$ .

The **Laplace transform** applies to nonnegative valued random variable  $X$  and is defined by

$$L_X(s) = E[e^{-sX}] \quad \text{for } s \geq 0.$$

This is a special case of (96) by setting  $\omega = -s$ .

The **Fourier transform** applies to both nonnegative and negative valued random variable  $X$  and is defined by

$$\Upsilon_X(s) = E[e^{i\theta X}],$$

where  $i = \sqrt{-1}$  and  $\theta$  is real. This is a special case of (96) by setting  $\omega = i\theta$ .

We will only use the Z and Laplace transforms in this book.

### 1.13.1 Z-transform

Consider a discrete and nonnegative random variable  $X$ , and let  $p_i = P(X = i)$ ,  $i = 0, 1, 2, \dots$  with  $\sum_{i=0}^{\infty} p_i = 1$ . The Z-transform of  $X$  is defined by

$$\Pi_X(z) = E[z^X] = \sum_{i=0}^{\infty} p_i z^i, \quad (114)$$

where  $z$  is a real number that satisfies  $0 \leq z \leq 1$ . Note that in many applications the Z-transforms is defined for complex  $z$ . However, for the purpose of this book, we will only consider real  $z$  within  $0 \leq z \leq 1$ .

**Homework 1.16**

Prove the following properties of the Z-transform  $\Pi_X(z)$ :

1.  $\lim_{z \rightarrow 1^-} \Pi_X(z) = 1$  ( $z \rightarrow 1^-$  is defined as  $z$  approaches 1 from below).
2.  $p_i = \Pi_X^{(i)}(0)/i!$  where  $\Pi_X^{(i)}(z)$  is the  $i$ th derivative of  $\Pi_X(z)$ .
3.  $E[X] = \lim_{z \rightarrow 1^-} \Pi_X^{(1)}(z)$ .  $\square$

For simplification of notation, in the following, we will use  $\Pi_X^{(i)}(1) = \lim_{z \rightarrow 1^-} \Pi_X^{(i)}(z)$ , but the reader must keep in mind that a straightforward substitution of  $z = 1$  in  $\Pi_X^{(i)}(z)$  is not always possible and the limit needs to be derived. An elegant way to show the 3rd property is to consider  $\Pi_X(z) = E[z^X]$ , and exchanging the operation of derivative and expectation, we obtain  $\Pi_X^{(1)}(z) = E[Xz^{X-1}]$ , so  $\Pi_X^{(1)}(1) = E[X]$ . Similarly,

$$\Pi_X^{(i)}(1) = E[X(X-1)\dots(X-i+1)]. \quad (115)$$

**Homework 1.17**

Show that the variance  $\text{var}[X]$  is given by

$$\text{var}[X] = \Pi_X^{(2)}(1) + \Pi_X^{(1)}(1) - (\Pi_X^{(1)}(1))^2. \quad \square \quad (116)$$

**Homework 1.18**

Derive a formula for  $E[X^i]$  using the Z-transform.  $\square$

As a Z-transform is a special case of the transform  $\Gamma_Y(\omega) = E[e^{\omega Y}]$ , the following results hold. If random variables  $X$  and  $Y$  are related by  $Y = aX + b$  for real numbers  $a$  and  $b$  then

$$\Pi_Y(z) = z^b \Pi_X(za). \quad (117)$$

Let random variable  $Y$  be the sum of independent random variables  $X_1, X_2, \dots, X_n$  ( $Y = \sum_{i=1}^n X_i$ ), The Z-transform of  $Y$  is given by

$$\Pi_Y(z) = \Pi_{X_1}(z) \Pi_{X_2}(z) \Pi_{X_3}(z) \dots \Pi_{X_n}(z). \quad (118)$$

If  $X_1, X_2, \dots, X_n$  are also identically distributed, then

$$\Pi_Y(z) = [\Pi_{X_1}(z)]^n. \quad (119)$$

Let us now consider several examples of Z-transforms of nonnegative discrete random variables. If  $X$  is a Bernoulli random variable with parameter  $p$ , then its Z-transform is given by

$$\Pi_X(z) = (1-p)z^0 + pz^1 = 1-p+pz. \quad (120)$$

Its mean is  $E[X] = \Pi_X^{(1)}(1) = p$  and by (116) its variance is  $p(1-p)$ .

If  $X$  is a Geometric random variable with parameter  $p$ , then its Z-transform is given by

$$\Pi_X(z) = p \sum_{i=1}^{\infty} (1-p)^{i-1} z^i = \frac{pz}{1 - (1-p)z}. \quad (121)$$

Its mean is  $E[X] = \Pi_X^{(1)}(1) = 1/p$  and by (116) its variance is  $(1-p)/p^2$ .

If  $X$  is a Binomial random variable with parameter  $p$ , then we can obtain its Z-transform either by definition or by realizing that a Binomial random variable is a sum of  $n$  IID Bernoulli random variables. Therefore its Z-transform is given by

$$\Pi_X(z) = (1 - p + pz)^n = [1 + (z - 1)p]^n. \quad (122)$$

### Homework 1.19

Verify that the latter is consistent with the Z-transform obtained using  $\Pi_X(z) = \sum_{i=0}^{\infty} p_i z^i$ .  $\square$

The mean of  $X$  is  $E[X] = \Pi_X^{(1)}(1) = np$  and by (116) its variance is  $np(1-p)$ .

If  $X$  is a Poisson random variable with parameter  $\lambda$ , then its Z-transform is given by

$$\Pi_X(z) = \sum_{i=0}^{\infty} p_i z^i = e^{-\lambda} \sum_{i=0}^{\infty} \frac{\lambda^i z^i}{i!} = e^{(z-1)\lambda}. \quad (123)$$

Its mean is  $E[X] = \Pi_X^{(1)}(1) = \lambda$  and by (116) its variance is also equal to  $\lambda$ .

We can now see the relationship between the Binomial and the Poisson random variables. If we consider the Z-transform of the Binomial random variable  $\Pi_X(z) = (1 - p + pz)^n$ , and set  $\lambda = np$  as a constant so that  $\Pi_X(z) = (1 + (z - 1)\lambda/n)^n$  and let  $n \rightarrow \infty$ , we obtain

$$\lim_{n \rightarrow \infty} (1 - p + pz)^n = \lim_{n \rightarrow \infty} [1 + (z - 1)\lambda/n]^n = e^{(z-1)\lambda} \quad (124)$$

which is exactly the Z-transform of the Poisson random variable. This proves the convergence of the binomial to the Poisson random variable if we keep  $np$  constant and let  $n$  go to infinity.

### 1.13.2 Laplace Transform

The Laplace transform of a non-negative random variable  $X$  with density  $f_X(x)$  is defined as

$$\mathcal{L}_X(s) = E[e^{-sX}] = \int_0^{\infty} e^{-sx} f_X(x) dx. \quad (125)$$

As it is related to the transform  $\Gamma_X(\omega) = E[e^{j\omega X}]$  by setting  $\omega = -js$ , similar derivations to those made for  $\Gamma_X(\omega)$  above give the following.

If  $X_1, X_2, \dots, X_n$  are  $n$  independent random variables then

$$\mathcal{L}_{X_1+X_2+\dots+X_n}(s) = \mathcal{L}_{X_1}(s) \mathcal{L}_{X_2}(s) \dots \mathcal{L}_{X_n}(s). \quad (126)$$

Let  $X$  and  $Y$  be random variables and  $Y = aX + b$ . The Laplace transform of  $Y$  is given by

$$\mathcal{L}_Y(s) = e^{-sb} \mathcal{L}_X(sa). \quad (127)$$



The  $n$ th moment of random variable  $X$  is given by

$$E[X^n] = (-1)^n \mathcal{L}_X^{(n)}(0) \quad (128)$$

where  $\mathcal{L}_X^{(n)}(0)$  is the  $n$ th derivative of  $\mathcal{L}_X(s)$  at  $s = 0$  (or at the limit  $s \rightarrow 0$ ). Therefore,

$$\text{var}[X] = E[X^2] - (E[X])^2 = (-1)^2 \mathcal{L}_X^{(2)}(0) - ((-1) \mathcal{L}_X^{(1)}(0))^2 = \mathcal{L}_X^{(2)}(0) - (\mathcal{L}_X^{(1)}(0))^2. \quad (129)$$

Let  $X$  be an exponential random variable with parameter  $\lambda$ . Its Laplace transform is given by

$$\mathcal{L}_X(s) = \frac{\lambda}{\lambda + s}. \quad (130)$$

### Homework 1.20

Derive (126)–(130) using the derivations made for  $\Gamma_X(\omega)$  as a guide.  $\square$

Now consider  $N$  to be a nonnegative discrete (integer) random variable of a probability distribution that has the Z-transform  $\Pi_N(z)$ , and let  $Y = X_1 + X_2 + \dots + X_N$ , where  $X_1, X_2, \dots, X_N$  are nonnegative IID random variables with a common distribution that has the Laplace transform  $\mathcal{L}_X(s)$ . Let us derive the Laplace transform of  $Y$ . Conditioning and unconditioning on  $N$ , we obtain

$$\mathcal{L}_Y(s) = E[e^{-sY}] = E_N[E[e^{-s(X_1+X_2+\dots+X_N|N)}]]. \quad (131)$$

Therefore, by independence of the  $X_i$ ,

$$\mathcal{L}_Y(s) = E_N[E[e^{-sX_1} + E[e^{-sX_2} + \dots + E[e^{-sX_N}]]] = E_N[(\mathcal{L}_X(s))^N]. \quad (132)$$

Therefore

$$\mathcal{L}_Y(s) = \Pi_N[(\mathcal{L}_X(s))]. \quad (133)$$

An interesting example of (133) is the case where the  $X_i$  are IID exponentially distributed each with parameter  $\lambda$ , and  $N$  is geometrically distributed with parameter  $p$ . In this case, we already know that since  $X$  is an exponential random variable, we have  $\mathcal{L}_X(s) = \lambda/(\lambda + s)$ , so

$$\mathcal{L}_Y(s) = \Pi_N\left(\frac{\lambda}{\lambda + s}\right). \quad (134)$$

We also know that  $N$  is geometrically distributed, so  $\Pi_N(z) = pz/[1 - (1 - p)z]$ . Therefore, from (134), we obtain,

$$\mathcal{L}_Y(s) = \frac{\frac{p\lambda}{\lambda+s}}{1 - \frac{(1-p)\lambda}{\lambda+s}} \quad (135)$$

and after some algebra we obtain

$$\mathcal{L}_Y(s) = \frac{p\lambda}{s + p\lambda}. \quad (136)$$

This result is interesting. We have shown that  $Y$  is exponentially distributed with parameter  $p\lambda$ .

**Homework 1.21**

Derive the density, the Laplace transform, the mean and the variance of  $Y$  in the following three cases.

1. Let  $X_1$  and  $X_2$  be exponential random variables with parameters  $\mu_1$  and  $\mu_2$ , respectively. In this case,  $Y$  is a hyperexponential random variable with density  $f_Y(y) = pf_{X_1}(y) + (1 - p)f_{X_2}(y)$ .
2. Let  $X_1$  and  $X_2$  be exponential random variables with parameters  $\mu_1$  and  $\mu_2$ , respectively. The hypoexponential random variable  $Y$  is defined by  $Y = X_1 + X_2$ .
3. Let  $Y$  be an Erlang random variable, namely,  $Y = \sum_{i=1}^k X_i$  where the  $X_i$ s are IID exponentially distributed random variables with parameter  $\mu$ .
4. Now plot the standard deviation to mean ratio for the cases of hyperexponential and Erlang random variables over a wide range of parameter values and discuss implications. For example, show that for Erlang( $k$ ) the standard deviation to mean ratio approaches zero as  $k$  approaches infinity.  $\square$

**1.14 Multivariate Random Variables and Transform**

A *multivariate random variable* is a vector  $X = (X_1, X_2, \dots, X_k)$  where each of the  $k$  components is a random variable. A multivariate random variable is also known as *random vector*. These  $k$  components of a random vector are related to events (outcomes of experiments) on the same sample space and they can be continuous or discrete. They also have a legitimate well defined joint distribution (or density) function. The distribution of each individual component  $X_i$  of the random vector is its marginal distribution. A transform of a random vector  $X = (X_1, X_2, \dots, X_k)$  is called *multivariate transform* and is defined by

$$\Gamma_X(\omega_1, \omega_2, \dots, \omega_k) = E[s^{\omega_1 X_1, \omega_2 X_2, \dots, \omega_k X_k}]. \quad (137)$$

**1.15 Probability Inequalities and Their Dimensioning Applications**

In the course of design of telecommunications networks, a fundamental important problem is how much capacity a link should have. If we consider the demand as a non-negative random variable  $X$  and the link capacity as fix scalar  $C > 0$ , we will be interested in the probability that the demand exceeds the capacity  $P(X > C)$ . The more we know about the distribution the more accurate out estimation of  $P(X > C)$ .

If we know only the mean, we use the so-called **Markov inequality**:

$$P(X > C) \leq \frac{E[X]}{C}. \quad (138)$$

**Homework 1.22**

Prove Eq. (138).

**Guide**

Define a new random variable  $U(C)$  a function of  $X$  and  $C$  defined by:  $U(C) = 0$  if  $X < C$ , and  $U(C) = C$  if  $X \geq C$ . Notice  $U(C) \leq X$ , so  $E[U(C)] \leq E[X]$ . Also,  $E[U(C)] = CP(U(C) = C) = CP(X \geq C)$ , and Eq. (138) follows.  $\square$

If we know the mean and the variance of  $X$ , then we can use the so-called **Chebyshev inequality**:

$$P(|X - E[X]| > C) \leq \frac{\text{var}[X]}{C^2}. \quad (139)$$

**Homework 1.23**

Prove Eq. (139).

**Guide**

Define a new random variable  $(X - E[X])^2$  and apply the Markov inequality putting  $C^2$  instead of  $C$  obtaining:

$$P((X - E[X])^2 \geq C^2) \leq \frac{E[(X - E[X])^2]}{C^2} = \frac{\text{var}[X]}{C^2}.$$

Notice that the two events  $(X - E[X])^2 \geq C^2$  and  $|X - E[X]| \geq C$  are identical.  $\square$

Another version of Chebyshev inequality is

$$P(|X - E[X]| > C^* \sigma) \leq \frac{1}{(C^*)^2} \quad (140)$$

for  $C^* > 0$ .

**Homework 1.24**

Prove and provide interpretation to Eq. (140).

**Guide**

Observe that the right hand side of (140) is equal to  $\frac{\text{var}[X]}{\text{var}[X](C^*)^2}$ .  $\square$

**Homework 1.25**

For a wide range of parameter value, study numerically how tight the bounds provided by Markov versus Chebyshev inequalities are. Discuss the differences and provide interpretations.  $\square$

A further refinement of the Chebyshev inequality is the following **Kolmogorov inequality**. Let  $X_1, X_2, X_3, \dots, X_k$  be a sequence of mutually independent random variables (not necessarily identically distributed) and let  $S_k = X_1 + X_2 + X_3 + \dots + X_k$  and  $\sigma(S_k)$  be the standard

deviation of  $S_k$ . Then for every  $\epsilon > 0$ ,

$$P(|S_k - E[S_k]| < \theta \sigma(S_k) \text{ for all } k = 1, 2, \dots, n) \geq 1 - \frac{1}{\theta^2}. \quad (141)$$

The interested reader may consult Feller [20] for the proof of the Kolmogorov inequality. We are however more interested in its teletraffic implication. If we let time be divided into consecutive intervals and we assume that  $X_i$  is the number of packets arrive during the  $i$ th interval, and if the number of packets arrive during the different intervals are mutually independent, then it is rare that we will have within a period of  $n$  consecutive intervals any period of  $k$  consecutive intervals ( $k \leq n$ ) during which the number of packets arriving is significantly more than the average.

## 1.16 Limit Theorems

Let  $X_1, X_2, X_3, \dots, X_k$  be a sequence of IID random variables with mean  $\lambda$  and variance  $\sigma^2$ . Let  $\bar{S}_k$  be the *sample mean* of these  $k$  random variables defined by

$$\bar{S}_k = \frac{X_1 + X_2 + X_3 + \dots + X_k}{k}.$$

This gives

$$E[\bar{S}_k] = \frac{E[X_1] + E[X_2] + E[X_3] + \dots + E[X_k]}{k} = \frac{k\lambda}{k} = \lambda.$$

Recalling that the  $X_i$ s are independent, we obtain

$$\text{var}[\bar{S}_k] = \frac{\sigma^2}{k}. \quad (142)$$

### Homework 1.26

Prove Eq. (142).  $\square$

Applying Chebyshev's inequality, we obtain

$$P(|\bar{S}_k - \lambda| \geq \epsilon) \leq \frac{\sigma^2}{k\epsilon^2} \text{ for all } \epsilon > 0. \quad (143)$$

Noticing that as  $k$  approaches infinity, the right-hand side of (143) approaches zero which implies that the left hand-side approaches zero as well. This leads to the so-called **the weak law of large numbers** that states the following. Let  $X_1, X_2, X_3, \dots, X_k$  be  $k$  IID random variables with common mean  $\lambda$ . Then

$$P\left(\left|\frac{X_1 + X_2 + X_3 + \dots + X_k}{k} - \lambda\right| \geq \epsilon\right) \rightarrow 0 \text{ as } k \rightarrow \infty \text{ for all } \epsilon > 0. \quad (144)$$

What the weak law or large number essentially says is that the sample mean approaches the mean as the sample size increases.

Next we state the Central Limit Theorem that we have mentioned in Section 1.10.6. Let  $X_1, X_2, X_3, \dots, X_k$  be  $k$  IID random variables with common mean  $\lambda$  and variance  $\sigma^2$ . Let random variable  $Y_k$  be defined as

$$Y_k = \frac{X_1 + X_2 + X_3 + \dots + X_k - k\lambda}{\sigma\sqrt{k}}. \quad (145)$$

Then,

$$\lim_{k \rightarrow \infty} P(Y_k \leq y) = \Phi(y) \quad (146)$$

where  $\Phi(\cdot)$  is the distribution function of a standard Gaussian random variable given by

$$\Phi(y) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^y e^{-t^2/2} dt.$$

### Homework 1.27

Prove that  $E[Y_k] = 0$  and that  $var[Y_k] = 1$  from first principles without using the Central Limit Theorem.  $\square$

As we mentioned in Section 1.10.6, the Central Limit Theorem is considered the most important result in probability. Notice that it implies that the sum of  $k$  IID random variable with common mean  $\lambda$  and variance  $\sigma^2$  is approximately Gaussian with mean  $k\lambda$  and variance  $k\sigma^2$  *regardless* of the distribution of these variables.

## 1.17 Link Dimensioning

Before we end this chapter on probability, we demonstrate how the probability concepts discussed so far can be used to provide a simple method for link dimensioning. Consider  $N$  independent sources (end-terminals), sharing a transmission link of capacity  $C$  [Mb/s]. Each of the sources transmits data in accordance with an on-off process, meaning, it alternates between two states: 1) the on state during which the source transmits at a rate  $R$  [Mb/s], and 2) the off state during which the source is idle. Assume that the proportion of time a source is in the on-state is  $p$ , so it is in the off-state  $1 - p$  of the time. The question is how much capacity should the link have so it can serve all  $N$  sources such that the probability that the demand exceeds the total link capacity is no higher than  $\alpha$ .

We first derive the distribution of the total traffic demanded by the  $N$  sources. Without loss of generality, let us normalize the traffic generated by a source during on period by setting  $R = 1$ . Realizing that the demand generated by a single source is Bernoulli distributed with parameter  $p$ , we obtain that the demand generated by all  $N$  sources has Binomial distribution with parameters  $p$  and  $N$ . Accordingly, finding the desired capacity is reduced to finding the smallest  $C$  such that

$$\sum_{i=C+1}^N \binom{N}{i} p^i (1-p)^{N-i} \leq \alpha. \quad (147)$$

Since the left hand-side of (147) increases as  $C$  decreases, And since its value is zero if  $C = N$ , all we need to do to find the optimal  $C$  is to compute the value of the left hand-side of (147) for  $C$  values of  $N - 1, N - 2, \dots$  until we find the first  $C$  value for which the inequality (147) is violated. Increasing that  $C$  value by one will give us the desired optimal  $C$  value.

If  $N$  is large we can use the Central Limit Theorem and approximate the Binomial distribution by a Gaussian distribution. Accordingly, the demand can be approximated by a Gaussian random variable with mean  $Np$  and variance  $Np(1-p)$  and simply find  $C$  such that the probability of our Gaussian random variable to exceed  $C$  is  $\alpha$ .

It is well known that Gaussian random variables obey the so-called 68-95-99.7% Rule which means that the following apply to a random variable  $X$  with mean  $m$  and standard deviation  $\sigma$ .

$$\begin{aligned} P(m - \sigma \leq X \leq m + \sigma) &= 0.68 \\ P(m - 2\sigma \leq X \leq m + 2\sigma) &= 0.95 \\ P(m - 3\sigma \leq X \leq m + 3\sigma) &= 0.997. \end{aligned}$$

Therefore, if  $\alpha = 0.0015$  then  $C$  should be three standard deviations from the mean, namely,

$$C = Np + 3\sqrt{Np(1-p)}. \quad (148)$$

## 2 Selected Stochastic Processes Topics

Aiming to understand behaviors of various natural and artificial processes, researchers often model them as collections of random variables where the mathematically defined statistical characteristics and dependencies of such random variables are fitted to those of the real processes. The research in the field of stochastic processes has therefore three facets:

**Theory:** mathematical explorations of stochastic processes models that aim to better understand their properties.

**Measurements:** taken on the real process in order to identify its statistical characteristics.

**Modelling:** fitting the measured statistical characteristics of the real process with those of a model and development of new models of stochastic processes that well match the real process.

This chapter provides background on basic theoretical aspects of stochastic processes which form a basis for queueing theory and teletraffic models discussed in the later chapters.

### 2.1 General Concepts

For a given *index set*  $T$ , a *stochastic process*  $\{X_t, t \in T\}$  is an indexed collection of random variables. They may or may not be identically distributed. In many applications the index  $t$  is used to model time. Accordingly, the random variable  $X_t$  for a given  $t$  can represent, for example, the number of telephone calls that have arrived at an exchange by time  $t$ .

If the index set  $T$  is countable, the stochastic process is called a *discrete time* process, or a *time series* [7, 12, 44]. Otherwise, the stochastic process is called a *continuous time* process. Considering our previous example, where the number of phone calls arriving at an exchange by time  $t$  is modelled as a continuous time process  $\{X_t, t \in T\}$ , we can alternatively, use a discrete time process to model, essentially, the same thing. This can be done by defining the discrete time process  $\{X_n, n = 1, 2, 3, \dots\}$ , where  $X_n$  is a random variable representing, for example, the number of calls arriving within the  $n$ th minute.

A stochastic process  $\{X_t, t \in T\}$  is called *discrete space* stochastic process if the random variables  $X_t$  are discrete, and it is called *continuous space* stochastic process if it is continuous. We therefore have four types of stochastic processes:

1. Discrete Time Discrete Space
2. Discrete Time Continuous Space
3. Continuous Time Discrete Space
4. Continuous Time Continuous Space.

A discrete time stochastic process  $\{X_n, n = 1, 2, 3, \dots\}$  is *strictly stationary* if for any subset of  $\{X_n\}$ , say,  $\{X_{n(1)}, X_{n(2)}, X_{n(3)}, \dots, X_{n(k)}\}$ , for any integer  $m$  the joint probability function  $P(X_{n(1)}, X_{n(2)}, X_{n(3)}, \dots, X_{n(k)})$ , is equal to the joint probability function

$P(X_{n(1)+m}, X_{n(2)+m}, X_{n(3)+m}, \dots, X_{n(k)+m})$ . In other words,  $P(X_{n(1)+m}, X_{n(2)+m}, X_{n(3)+m}, \dots, X_{n(k)+m})$  is independent of  $m$ . In this case, the probability structure of the process does not change with time. An equivalent definition for strict stationarity is applied also for a continuous time process accept that in that case  $m$  is non-integer. Notice that for the process to be strictly stationary, the value of  $k$  is unlimited as the joint probability should be independent of  $m$  for any subset of  $\{X_n, n = 1, 2, 3, \dots\}$ . If  $k$  is limited to some value  $k^*$ , we say that the process is *stationary of order  $k^*$* .

A equivalent definition applies to a continuous time stochastic process. A continuous time stochastic process  $X_t$  is said to be strictly stationary if its statistical properties do not change with a shift of the origin. In other words the process  $X_t$  statistically the same as the process  $X_{t-d}$  for any value of  $d$ .

An important stochastic process is the *Gaussian Process* defined as a process that has the property that the joint probability function (density) associated with any set of times is multivariate Gaussian. The importance of the Gaussian process lies in its property to be an accurate model for superposition of many independent processes. This makes the Gaussian process a useful model for heavily multiplexed traffic which arrive at switches or routers deep in a major telecommunications network. Fortunately, the Gaussian process is not only useful, but it is also simple and amenable to analysis. Notice that for a multivariate Gaussian distribution, all moments of order higher than two are zero, and therefore, for a Gaussian process, stationarity of order two also called *weak stationarity* implies strict stationarity. For a time series  $\{X_n, n = 1, 2, 3, \dots\}$ , weak stationarity implies that, for all  $n$ ,  $E[X_n]$  is constant, denoted  $E[X]$ , independent of  $n$ . Namely, for all  $n$ ,

$$E[X] = E[X_n]. \quad (149)$$

Weak stationarity (because it is stationarity of order two) also implies that the covariance between  $X_n$  and  $X_{n+k}$ , for any  $k$ , is independent of  $n$ , and is only a function of  $k$ , denoted  $U(k)$ . Namely, for all  $n$ ,

$$U(k) = \text{cov}(X_n, X_{n+k}). \quad (150)$$

Notice that, the case of  $k = 0$  in Eq. (150), namely,

$$U(0) = \text{cov}(X_n, X_n) = \text{var}[X_n] \quad (151)$$

implies that the variance of  $X_n$  is also independent of  $n$ . Also for all integer  $k$ ,

$$U(-k) = U(k) \quad (152)$$

because  $\text{cov}(X_n, X_{n+k}) = \text{cov}(X_{n+k}, X_n) = \text{cov}(X_n, X_{n-k})$ . The function  $U(k)$ ,  $k = 0, 1, 2, \dots$ , is called the *autocovariance function*. The value of the autocovariance function at  $k$ ,  $U(k)$ , is also called the autocovariance of lag  $k$ .

Important parameters are the so-called Autocovariance Sum, denoted  $S$ , and Asymptotic Variance Rate (AVR) denoted  $v$  [4, 5]. They are defined by:

$$S = \sum_{i=1}^{\infty} U(i) \quad (153)$$

and

$$v = \sum_{i=-\infty}^{\infty} U(i). \quad (154)$$



Notice that

$$v = 2S + \text{var}[X_n]. \quad (155)$$

Another important definition of the AVR which justifies its name is

$$v = \lim_{n \rightarrow \infty} \frac{S_n}{n}. \quad (156)$$

We will further discuss these concepts in Section 17.1.

## Homework 2.1

Prove that the above two definitions are equivalent; namely, prove that

$$\lim_{n \rightarrow \infty} \frac{S_n}{n} = 2S + \text{var}[X_n]. \quad (157)$$

## Guide

Use Eq. (92)  $\square$

The *autocorrelation function* at lag  $k$ , denoted  $C(k)$ , is the normalized version of the autocovariance function, and since by weak stationarity, for all  $i$  and  $j$ ,  $\text{var}[X_j] = \text{var}[X_i]$ , it is given by:

$$C(k) = \frac{U(k)}{\text{var}[X_n]}. \quad (158)$$

A stochastic process is called *ergodic* if every realization contains sufficient information on the probabilistic structure of the process. For example, let us consider a process which can be in either one of two realization: either  $X_n = 1$  for all  $n$ , or  $X_n = 0$  for all  $n$ . Assume that each one of these two realizations occur with probability 0.5. If we observe any one of these realizations, regardless of the duration of the observations, we shall never conclude that  $E[A] = 0.5$ . We shall only have the estimations of either  $E[A] = 0$  or  $E[A] = 1$ , depends on which realization we happen to observe. Such a process is not ergodic.

Assuming  $\{X_n, n = 1, 2, 3, \dots\}$  is ergodic and stationary, and we observe  $m$  observations of this  $\{X_n\}$  process, denoted by  $\{\hat{A}_n, n = 1, 2, 3, \dots, m\}$ , then the mean of the process  $E[A]$  can be estimated by

$$\hat{E}[A] = \frac{1}{m} \sum_{n=1}^m \hat{A}_n, \quad (159)$$

and the autocovariance function  $U(k)$  of the process can be estimated by

$$\hat{U}(k) = \frac{1}{m-k} \sum_{n=k+1}^m (\hat{A}_{n-k} - E[A])(\hat{A}_n - E[A]). \quad (160)$$

## 2.2 Two Orderly and Memoryless Point Processes

In this section we consider a very special class of stochastic processes called *point* processes that also possess two properties: *orderliness* and *memorylessness*. After providing, somewhat

intuitive, definitions of these concepts, we will discuss two processes that belong to this special class: one is discrete-time - called the *Bernoulli process* and the other is continuous-time - called the *Poisson process*.

We consider here a physical interpretation, where a *point process* is a sequence of events which we call *arrivals* occurring at random in points of time  $t_i$ ,  $i = 1, 2, \dots$ ,  $t_{i+1} > t_i$ , or  $i = \dots, -2, -1, 0, 1, 2, \dots$ ,  $t_{i+1} > t_i$ . The index set, namely, the time, or the set where the  $t_i$  get their values from, can be continuous or discrete, although in most books the index set is considered to be the real line, or its non-negative part. We call our events arrivals to relate is to the context of queueing theory, where a point process typically corresponds to points of arrivals, i.e.,  $t_i$  is the time of the  $i$ th arrival that joins a queue. A point process can be defined by its *counting process*  $\{N(t), t \geq 0\}$ , where  $N(t)$  is the number of arrivals occurred within  $[0, t)$ . A counting process  $\{N(t)\}$  has the following properties:

1.  $N(t) \geq 0$ ,
2.  $N(t)$  is integer,
3. if  $s > t$ , then  $N(s) \geq N(t)$  and  $N(s) - N(t)$  is the number of occurrences within  $(t, s]$ .

Note that  $N(t)$  is not an independent process because for example, if  $t_2 > t_1$  then  $N(t_2)$  is dependent on the number of arrivals in  $[0, t_1)$ , namely,  $N(t_1)$ .

Another way to define a point process is by the stochastic process of the interarrival times  $\Delta_i$  where  $\Delta_i = t_{i+1} - t_i$ .

One important property of a counting process is the so-called *Orderliness* which means that the probability that two or more arrivals happen at once is negligible. Mathematically, for a continuous-time counting process to be *orderly*, it should satisfy:

$$\lim_{\Delta t \rightarrow 0} P(X(t + \Delta t) - X(t) > 1 \mid X(t + \Delta t) - X(t) \geq 1) = 0. \quad (161)$$

Another very important property is the *memorylessness*. A stochastic process is *memoryless* if at any point in time, the future evolution of the process is statistically independent of its past.

### 2.2.1 Bernoulli Process

The Bernoulli process is a discrete time stochastic process made up of a sequence of IID Bernoulli distributed random variables  $\{X_i, i = 0, 1, 2, 3, \dots\}$  where for all  $i$ ,  $P(X_i = 1) = p$  and  $P(X_i = 0) = 1 - p$ . In other words, we divide time into consecutive equal time slots. At each time-slot we conduct a Bernoulli experiment. Then for each time-slot  $i$ , we conduct a bernoulli experiment. If  $X_i = 1$ , we say that there was an *arrival* at time-slot  $i$ . Otherwise, if  $X_i = 0$ , we say that there was no arrival at time slot  $i$ .

The Bernoulli process is both orderly and memoryless. It is orderly because, by definition, no more than one arrival can occur at any time slot as the Bernoulli random variable takes values of more than one with probability zero. It is also memoryless because the Bernoulli trials are independent, so at any discrete point in time  $n$ , the future evolution of the process is independent of its past.

The counting process for the Bernoulli process is another discrete-time stochastic process  $\{N(n), n \geq 0\}$  which is a sequence of Binomial random variables  $N(n)$  representing the total number of arrivals occurring within the first  $n$  time slots. Notice that since we start from slot 0,  $N(n)$  does not include slot  $n$  in the counting. That is, we have

$$P[N(n) = i] = \binom{n}{i} p^i (1-p)^{n-i} \quad i = 0, 1, 2, \dots, n. \quad (162)$$

The concept of an interarrival time for the Bernoulli process can be explained as follows. Let us assume without loss of generality that there was an arrival at time-slot  $k$ , the interarrival time will be the number of slots between  $k$  and the first time-slot to have an arrival following  $k$ . We do not count time-slot  $k$  but we do count the time-slot of the next arrival. Because the Bernoulli process is memoryless, the interarrival times are IID, so we can drop the index  $i$  of  $\Delta_i$ , designating the  $i$  interarrival time, and consider the probability function of the random variable  $\Delta$  representing any interarrival time. Because  $\Delta$  represents a number of Bernoulli trials until a success, it is geometrically distributed, and its probability function is given by

$$P(\Delta = i) = p(1-p)^{i-1} \quad i = 1, 2, \dots \quad (163)$$

Another important statistical measure is the time it takes  $n$  until the  $i$ th arrival. This time is a sum of  $i$  interarrival times which is a sum of  $i$  geometric random variables which we already know has a Pascal distribution with parameters  $p$  and  $i$ , so we have

$$P[\text{the } i\text{th arrival occurs in time slot } n] = \binom{n-1}{i-1} p^i (1-p)^{n-i} \quad i = i, i+1, i+2, \dots \quad (164)$$

The reader may notice that the on-off sources discussed in Section 1.17 could be modeled as Bernoulli processes where the on periods are represented by consecutive successes of Bernoulli trials and the off periods by failures. In this case, for each on-off process, the length of the on and the off periods are both geometrically distributed. Accordingly, the **multiplexing** of  $N$  Bernoulli processes with parameter  $p$  is another discrete-time stochastic process where the number of arrivals during the different slots are IID and binomial distributed with parameters  $N$  and  $p$ .

## Homework 2.2

Prove the last statement.  $\square$

Another important concept is **merging** of processes which is different from multiplexing. Let us use a sensor network example to illustrate it. Consider  $N$  sensors that are spread around a country to detect certain events. Time is divided into consecutive time-slots and a sensor is silent if it does not detect an event in a given time-slot and active (transmitting an alarm signal) if it does. Assume that time-slots during which the  $i$ th sensor is active follow a Bernoulli process with parameter  $p_i$ , namely, the probability that sensor  $i$  detects an event in a given time-slot is equal to  $p_i$ , and that the probability of such detection is independent from time-slot to time-slot. We also assume that the  $N$  Bernoulli processes are independent. Assume that an alarm is sound during a time-slot when at least one of the sensors is active. We are interested in the discrete-time process representing alarm sounds. The probability that an alarm is sound in a given time-slot is the probability that at least one of the sensors is active which is one

minus the probability that they are all silent. Therefore the probability that the alarm is sound is given by

$$P_a = 1 - \prod_{i=1}^N (1 - p_i). \quad (165)$$

Now, considering the independence of the processes, we can realize that the alarms follow a Bernoulli process with parameter  $P_a$ .

In general, an arrival in the process that results from merging of  $N$  Bernoulli processes is the process of time-slots during which at least one of the  $N$  processes records an arrival. Unlike multiplexing in which we are interested in the total number of arrivals, in merging we are only interested to know if there was at least one arrival within a time-slot without any interest of how many arrivals there were in total.

Let us now consider **splitting**. Consider a Bernoulli process with parameter  $p$  and then color each arrival, independently of all other arrivals, in red with probability  $q$  and in blue with probability  $1 - q$ . Then in each time slot we have a red arrival with probability  $pq$  and a blue one with probability  $p(1 - q)$ . Therefore, the red arrivals follow a Bernoulli process with parameter  $pq$  and the blue arrivals follow a Bernoulli process with parameter  $p(1 - q)$ .

### 2.2.2 Poisson Process

The Poisson process is a continuous-time point process which is also memoryless and orderly. It applies to many cases where a certain event occurs at different points in time. Such occurrences of the events could be, for example, arrivals of phone call requests at a telephone exchange. As mentioned above such a process can be described by its *counting process*  $\{N(t), t \geq 0\}$  representing the total number of occurrences by time  $t$ .

A counting process  $\{N(t)\}$  is defined as a *Poisson process* with rate  $\lambda > 0$  if it satisfies the following three conditions.

1.  $N(0) = 0$ .
2. The number of occurrences in two non-overlapping intervals are independent. That is, for any  $s > t > u > v > 0$ , the random variable  $N(s) - N(t)$ , and the random variable  $N(u) - N(v)$  are independent. This means that the Poisson process has what is called *independent increments*.
3. The number of occurrences in an interval of length  $t$  has a Poisson distribution with mean  $\lambda t$ .

These three conditions will be henceforth called the *Three Poisson process conditions*.

By definition, the Poisson process has what is called *stationary increments* [47, 53], that is, for any  $t_2 > t_1$ , the random variable  $X(t_2) - X(t_1)$ , and the random variable  $X(t_2 + u) - X(t_1 + u)$  have the same distribution for any  $u > 0$ . In both cases, the distribution is Poisson with parameter  $\lambda(t_2 - t_1)$ . Intuitively, if we choose the time interval  $\Delta = t_2 - t_1$  to be arbitrarily small (almost a “point” in time), then probability of having an occurrence there is the same regardless of where the “point” is. Loosely speaking, every point in time has the same chance

of having a occurrence. Therefore, occurrences are equally likely to happen at all times. This property is also called *time-homogeneity* [10].

Another important property of the Poisson process is that the inter-arrival times of occurrences is exponentially distributed with parameter  $\lambda$ . This is shown by considering  $s$  to be an occurrence and  $T$  the time until the next occurrence, noticing that  $P(T > t) = P(N(t) = 0) = e^{-\lambda t}$ , and recalling the properties of independent and stationary increments. As a result, the mean interarrival time is given by

$$E[T] = \frac{1}{\lambda}. \quad (166)$$

By the memoryless property of the exponential distribution, the time until the next occurrence is always exponentially distributed and therefore, at any point in time, not necessarily at points of occurrences, the future evolution of the Poisson process is independent of the past, and is always probabilistically the same. The Poisson process is therefore memoryless. Actually, the independence of the past can be explained also by the Poisson process property of *independent increments* [53], and the fact that the future evolution is probabilistically the same can also be explained by the stationary increments property.

An interesting paradox emerges when one considers the Poisson process. If we consider a random point in time, the time until the next occurrence event has exponential distribution with parameter  $\lambda$ . Because the Poisson process in reverse is also a Poisson process, then at any point in time, the time passed from the last Poisson occurrence event also has exponential distribution with parameter  $\lambda$ . Therefore, if we pick a random point in time the mean length of the interval between two consecutive Poisson occurrences must be  $1/\lambda + 1/\lambda = 2/\lambda$ . How can we explain this phenomenon, if we know that the time between consecutive Poisson occurrences must be exponentially distributed with mean  $1/\lambda$ ? The explanation is that if we pick a point of time at random we are likely to pick an interval that is longer than the average.

### Homework 2.3

Demonstrate the above paradox as follows. Generate a Poisson process with rate  $\lambda = 1$  for a period of time of length  $T = 10,000$ . Pick a point in time from a uniform distribution within the interval  $[1, 10000]$ . Record the length of the interval (between two consecutive Poisson occurrences) that includes the chosen point in time Repeat the experiment 1000 times.  $\square$

A **superposition** of a number of Poisson processes is another point process that comprises all the points of the different processes. Another important property of the Poisson process is that superposition of two Poisson processes with parameters  $\lambda_1$  and  $\lambda_2$  is a Poisson process with parameter  $\lambda_1 + \lambda_2$ . Notice that in such a case, at any point in time, the time until the next occurrence is a competition between two exponential random variables one with parameter  $\lambda_1$  and the other with parameter  $\lambda_2$ . Let  $T$  be the time until the winner of the two occurs, and let  $T_1$  and  $T_2$  be the time until the next occurrence of the first process and the second process, respectively. Then by (49)

$$P(T > t) = e^{-(\lambda_1 + \lambda_2)t}. \quad (167)$$

Thus, the interarrival time of the superposition is exponentially distributed with parameter  $\lambda_1 + \lambda_2$ . This is consistent with the fact that the superposition of the two processes is a Poisson process with parameter  $\lambda_1 + \lambda_2$ .

## Homework 2.4

Prove that a superposition of  $N$  Poisson processes with parameters  $\lambda_1, \lambda_2, \dots, \lambda_N$ , is a Poisson process with parameter  $\lambda_1 + \lambda_2 + \dots + \lambda_N$ .  $\square$

Another interesting question related to superposition of Poisson processes is the question of what is the probability that the next event that occurs will be of a particular process. This is equivalent to the question of having say two exponential random variables  $T_1$  and  $T_2$  with parameters  $\lambda_1$  and  $\lambda_2$ , respectively, and we are interested in the probability of  $T_1 < T_2$ . By (50),

$$P(T_1 < T_2) = \frac{\lambda_1}{\lambda_1 + \lambda_2}. \quad (168)$$

Before we introduce further properties of the Poisson process, we shall introduce the following definition: a function  $g(\cdot)$  is  $o(\Delta t)$  if

$$\lim_{\Delta t \rightarrow 0} \frac{g(\Delta t)}{\Delta t} = 0. \quad (169)$$

Examples of functions which are  $o(\Delta t)$  are  $g(x) = x^v$  for  $v > 1$ . Sum or product of two functions which are  $o(\Delta t)$  is also  $o(\Delta t)$ , and a constant times a function which is  $o(\Delta t)$  is  $o(\Delta t)$ .

If a counting process  $\{N(t)\}$  is a *Poisson process* then, for a small interval  $\Delta t$ , we have:

1.  $P(N(\Delta t) = 0) = 1 - \lambda\Delta t + o(\Delta t)$
2.  $P(N(\Delta t) = 1) = \lambda\Delta t + o(\Delta t)$
3.  $P(N(\Delta t) \geq 2) = o(\Delta t)$ .

The above three conditions will henceforth be called *small interval conditions*. To show the first, we know that  $X(\Delta t)$  has a Poisson distribution, therefore

$$P(N(\Delta t) = 0) = e^{-\lambda\Delta t} \quad (170)$$

and developing it into a series gives,

$$P(X(\Delta t) = 0) = 1 - \lambda\Delta t + o(\Delta t). \quad (171)$$

The second is shown by noticing that  $P(N(\Delta t) = 1) = \lambda\Delta t P(N(\Delta t) = 0)$  and using the previous result. The third is obtained by  $P(N(\Delta t) \geq 2) = 1 - P(N(\Delta t) = 1) - P(N(\Delta t) = 0)$ . In fact, these three small interval conditions plus the stationarity and independence properties together with  $N(0) = 0$ , can serve as an alternative definition of the Poisson process. These properties imply that the number of occurrences per interval has a Poisson distribution.

## Homework 2.5

Prove the last statement. Namely, show that the three small-interval conditions plus the stationarity and independence properties together with  $N(0) = 0$  are equivalent to the Three Poisson Conditions.  $\square$

In many networking applications, it is of interest to study the effect of **splitting** of packet arrival processes. In particular, we will consider two types of splitting: *random splitting* and *regular splitting*. To explain the difference between the two, consider an arrival process of packets to a certain switch called Switch X. This packet arrival process is assumed to follow a Poisson process with parameter  $\lambda$ . Some of these packets are then forwarded to Switch A and the others to Switch B. We are interested in the process of packets arriving from Switch X to Switch A, designated *X-A Process*.

Under random splitting, every packet that arrives at Switch X is forwarded to A with probability  $p$  and to B with probability  $1-p$  independently of any other event associated with other packets. In this case, the packets stream from X to A follows a Poisson process with parameter  $p\lambda$ .

### Homework 2.6

Prove that the X-A Process is a Poisson process with parameter  $p\lambda$ .  $\square$

It may be interesting to notice that the interarrival times in the X-A Process are exponentially distributed because they are geometric sums of exponential random variables.

Under regular splitting, the first packet that arrives at Switch X is forwarded to A the second to B, the third to A, the fourth to B, etc. In this case, the packets stream from X to A (the X-A Process) will follow a stochastic process which is a point process where the interarrival times are Erlang distributed with parameter  $\lambda$  and 2.

### Homework 2.7

1. Prove the last statement.
2. Derive the mean and the variance of the interarrival times of the X-A process in the two cases above: random splitting and regular splitting.
3. Consider now 3-way splitting. Derive and compare the mean and the variance of the interarrival times for the regular and random splitting cases.
4. Repeat the above for  $n$ -way splitting and let  $n$  increase arbitrarily. What can you say about the burstiness/variability of regular versus random splitting.  $\square$

The properties of the Poisson process, namely, independence and time-homogeneity, make it unique in its ability to randomly inspect other continuous-time stochastic processes in a way that the sample it provides gives us enough information on what is called *time-averages*. Examples of time-averages are the proportion of time a process  $X(t)$  is in state  $i$ , i.e., the proportion of time during which  $X(t) = i$ . Or the overall mean of the process defined by

$$E[X(t)] = \frac{\int_0^T X(t)dt}{T} \quad (172)$$

for an arbitrarily large  $T$ . These properties that an occurrence can occur at any time with equal probability, regardless of times of past occurrences, gave rise to the expression a *pure chance process* for the Poisson process.

## Homework 2.8

Consider a Poisson process with parameter  $\lambda$ . You know that there was exactly one occurrence during the interval  $[0,1]$ . Prove that the time of the occurrence is uniformly distributed within  $[0,1]$ .

### Guide

For  $0 \leq t \leq 1$ , consider

$$P(\text{occurrence within } [0, t] \mid \text{exactly one occurrence within } [0, 1])$$

and use the definition of conditional probability. Notice that the latter is equal to:

$$\frac{P(\text{one occurrence within } [0, t] \text{ and no occurrence within } [t, 1])}{P(\text{exactly one occurrence within } [0, 1])}$$

or

$$\frac{P(\text{one occurrence within } [0, t])P(\text{no occurrence within } [t, 1])}{P(\text{exactly one occurrence within } [0, 1])}.$$

Then recall that the number of occurrences in any interval of size  $T$  has Poisson distribution with parameter  $\lambda T$ .  $\square$

## 2.3 Markov Modulated Poisson Process

The stochastic process called Markov modulated Poisson process (MMPP) is a point process that behaves as a Poisson process with parameter  $\lambda_i$  for a period of time that is exponentially distributed with parameter  $\delta_i$ . Then it moves to mode (state)  $j$  where it behaves like a Poisson process with parameter  $\lambda_j$  for a period of time that is exponentially distributed with parameter  $\mu_j$ . In general, the MMPP can have an arbitrary number of modes, so it requires a transition probability matrix as an additional set of parameters to specify the probability that it moves to mode  $j$  given that it is in mode  $i$ . However, we are mostly interested in the simplest case of MMPP – the two mode MMPP denoted MMPP(2) and defined by only four parameters:  $\lambda_0$ ,  $\lambda_1$ ,  $\delta_0$ , and  $\delta_1$ . The MMPP(2) behaves as a Poisson process with parameter  $\lambda_0$  for a period of time that is exponentially distributed with parameter  $\mu_0$ . Then moves to mode 1 where it behaves like a Poisson process with parameter  $\lambda_1$  for a period of time that is exponentially distributed with parameter  $\mu_1$ . Then it switches back to mode 0, etc. alternating between the two modes 0 and 1.

## 2.4 Discrete Time Markov Chains

### 2.4.1 Definitions and Preliminaries

Markov chains are certain discrete space stochastic processes which are amenable for analysis and hence are very popular for analysis, traffic characterization and modeling of queueing and telecommunications networks and systems. They can be classified into two groups: discrete



time Markov chains discussed here and continues time Markov chains discussed in the next section.

A discrete-time Markov chain is a discrete time stochastic process  $\{X_n, n = 0, 1, 2, \dots\}$  with the Markov property; namely, that at any point in time  $n$ , the future evolution of the process is dependent only on the state of the process  $X_n$ , and is independent of the past evolution of the process. The state of the process can be a scalar or a vector. In this section, for simplicity we will mainly discuss the case where the state of the process is a scalar, but we will also demonstrate how to extend the discussion to a multiple dimension case.

The discrete time Markov chain  $\{X_n, n = 0, 1, 2, \dots\}$  at any point in time may take many possible values. The set of these possible values is finite or countable and it is called the state space of the Markov chain, denoted by  $\Theta$ . A *time-homogeneous Markov chain* is a process in which

$$P(X_{n+1} = i \mid X_n = j) = P(X_n = i \mid X_{n-1} = j) \quad \text{for all } n.$$

We will only consider, in this section, Markov chains which are time-homogeneous.

A discrete-time time-homogeneous Markov chain is characterized by the property that, for any  $n$ , given  $X_n$ , the distribution of  $X_{n+1}$  is fully defined regardless of states that occur before time  $n$ . That is,

$$P(X_{n+1} = j \mid X_n = i) = P(X_{n+1} = j \mid X_n = i, X_{n-1} = i_{n-1}, X_{n-2} = i_{n-2}, \dots). \quad (173)$$

### 2.4.2 Transition Probability Matrix

A Markov chain is characterized by the so-called *Transition Probability Matrix*  $\mathbf{P}$  which is a matrix of one step transition probabilities  $P_{ij}$  defined by

$$P_{ij} = P(X_{n+1} = j \mid X_n = i) \quad \text{for all } n. \quad (174)$$

We can observe in the latter that the event  $\{X_{n+1} = j\}$  depends only on the state of the process at  $X_n$  and the transition probability matrix  $\mathbf{P}$ .

Since the  $P_{ij}$ s are probabilities and since when you transit out of state  $i$ , you must enter some state, all the entries in  $\mathbf{P}$  are non-negatives, less or equal to 1, and the sum of entries in each row of  $\mathbf{P}$  must add up to 1.

### 2.4.3 Chapman-Kolmogorov Equation

Having defined the one-step transition probabilities  $P_{ij}$  in (174), let us define the  $n$ -step transition probability from state  $i$  to state  $j$  as

$$P_{ij}^{(n)} = P(X_n = j \mid X_0 = i). \quad (175)$$

The following is known as the Chapman-Kolmogorov equation:

$$P_{ij}^{(n)} = \sum_{k \in \Theta} P_{ik}^{(m)} P_{kj}^{(n-m)}, \quad (176)$$

for any  $m$ , such that  $0 < m < n$ .

Let  $\mathbf{P}^{(n)}$  be the matrix that its entries are the  $P_{ij}^{(n)}$  values.

## Homework 2.9

First prove the Chapman-Kolmogorov equation and then use it to prove:

1.  $\mathbf{P}^{(k+n)} = \mathbf{P}^{(k)} \times \mathbf{P}^{(n)}$
2.  $\mathbf{P}^{(n)} = \mathbf{P}^n$ .     $\square$

### 2.4.4 Marginal Probabilities

Consider the marginal distribution  $\pi_n(i) = P(X_n = i)$  of the Markov chain at time  $n$ , over the different states  $i \in \Theta$ . Assuming that the process started at time 0, the initial distribution of the Markov chain is  $\pi_0(i) = P(X_0 = i)$ ,  $i \in \Theta$ . Then  $\pi_n(i)$ ,  $i \in \Theta$ , can be obtained based on the marginal probability  $\pi_{n-1}(i)$  as follows

$$\pi_n(j) = \sum_{k \in \Theta} P_{kj} \pi_{n-1}(k), \quad (177)$$

or based on the initial distribution by

$$\pi_n(j) = \sum_{k \in \Theta} P_{kj}^{(n)} \pi_0(k), \quad j \in \Theta \quad (178)$$

or, in matrix notation

$$\pi_n(j) = \sum_{k \in \Theta} P_{kj}^{(n)} \pi_0(k), \quad (179)$$

Let the vector  $\Pi_n$  be defined by  $\Pi_n = \{\pi_n(j), j = 0, 1, 2, 3, \dots\}$ . The vector  $\Pi_n$  can be obtained by

$$\Pi_n = \Pi_{n-1} \mathbf{P} = \Pi_{n-2} \mathbf{P}^2 = \dots = \Pi_0 \mathbf{P}^n. \quad (180)$$

### 2.4.5 Properties and Classification of States

One state  $i$  is said to be *accessible* from a second state  $j$  if there exists  $n$ ,  $n = 0, 1, 2, \dots$ , such that

$$P_{ji}^{(n)} > 0. \quad (181)$$

This means that there is a positive probability for the Markov chain to reach state  $i$  at some time in the future if it is now in state  $j$ .

A state  $i$  is said to *communicate* with state  $j$  if  $i$  is accessible from  $j$  and  $j$  is accessible from  $i$ .

## Homework 2.10

Prove the following:

1. A state communicates with itself.
2. If state  $a$  communicates with  $b$ , then  $b$  communicates with  $a$ .

3. If state  $a$  communicates with  $b$ , and  $b$  communicates with  $c$ , then  $a$  communicates with  $c$ .  $\square$

A *communicating class* is a set of states that every pair of states in it communicates with each other.

### Homework 2.11

Prove that a state cannot belong to two different classes. In other words, two different classes must be disjoint.  $\square$

The latter implies that the state space  $\Theta$  is divided into a number (finite or infinite) of communicating classes.

### Homework 2.12

Provide an example of a Markov chain with three communicating classes.  $\square$

A communicating class is said to be *closed* if no state outside the class is accessible from a state that belongs to the class.

A Markov chain is said to be *irreducible* if all the states in its state space are accessible from each other. That is, the entire state space is one communicating class.

A state  $i$  has *period*  $m$  if  $m$  is the greatest common divisor of the set  $\{n : P(X_n = i | X_0 = i) > 0\}$ . In this case, the Markov chain can return to state  $i$  only in a number of steps that is a multiple of  $m$ . A state is said to be *aperiodic* if it has a period of one.

### Homework 2.13

Prove that in a communicating class, it is not possible that there are two states that have different periods.  $\square$

Given that the Markov chain is in state  $i$ , define *return time* as the random variable representing the next time the Markov chain returns to state  $i$ . Notice that the return time is a random variable  $R_i$ , defined by

$$R_i = \min\{n : X_n = i \mid X_0 = i\}. \quad (182)$$

A state is called *transient* if, given that we start in it, there is a positive probability that we will never return back to it. In other words, state  $i$  is transient if  $P(R_i < \infty) < 1$ . A state is called *recurrent* if it is not transient. Namely, if  $P(R_i < \infty) = 1$ .

Because in the case of a recurrent state  $i$ , the probability to return to state  $i$  in finite time is one, the process will visit state  $i$  infinitely many number of times. However, if  $i$  is transient, then the process will visit state  $i$  only a geometrically distributed number of times with parameter. (Notice that the probability of “success” is  $1 - P(R_i < \infty)$ .) In this case the number of visits in state  $i$  is finite with probability 1.

**Homework 2.14**

Show that state  $i$  is recurrent if and only if

$$\sum_{n=0}^{\infty} p_{ii}^{(n)} = \infty.$$

**Guide**

This can be shown by showing that if the condition holds, the Markov chain will visit state  $i$  an infinite number of times, and if it does not hold, the Markov chain will visit state  $i$  a finite number of times. Let  $Y_n = J_i(X_n)$ , where  $J_i(x)$  is a function defined for  $x = 0, 1, 2, \dots$ , taking the value 1 if  $x = i$ , and 0 if  $x \neq i$ . Notice that  $E[J_i(X_n) \mid X_0 = i] = P(X_n = i \mid X_0 = i)$ , and consider summing up both sides of the latter.  $\square$

**Homework 2.15**

Prove that if state  $i$  is recurrent then all the states in a class that  $i$  belongs to are recurrent. In other words, prove that recurrence is a class property.

**Guide**

Consider  $m$  and  $n$ , such that  $P_{ji}^{(m)} > 0$  and  $P_{ij}^{(n)} > 0$ , and argue that  $P_{ji}^{(m)} P_{ii}^{(k)} P_{ij}^{(n)} > 0$  for some  $m, k, n$ . Then use the ideas and result of the previous proof.  $\square$

**Homework 2.16**

Provide an example of a Markov chain where  $P(R_i < \infty) = 1$ , but  $E[R_i] = \infty$ .  $\square$

State  $i$  is called *positive recurrent* if  $E[R_i]$  is finite. A recurrent state that is not positive recurrent is called *null recurrent*. In a finite state Markov chain, there are no null recurrent states, i.e., all recurrent states must be positive recurrent. We say that a Markov chain is *stable* if all its states are positive recurrent. This notion of stability is not commonly used for Markov chains or stochastic processes in general and it is different from other definitions of stability. It is however consistent with the notion of stability of queueing systems and this is the reason we use it here.

A Markov chain is said to be aperiodic if all its states are aperiodic.

**2.4.6 Steady State Probabilities**

Consider an irreducible, aperiodic and stable Markov chain. Then the following limit exists.

$$\Pi = \lim_{n \rightarrow \infty} \Pi_n = \lim_{n \rightarrow \infty} \Pi_0 \mathbf{P}^n \quad (183)$$

and it satisfies

$$\Pi = \text{row of } \lim_{n \rightarrow \infty} \mathbf{P}^n \quad (184)$$

where *row of*  $\lim_{n \rightarrow \infty} \mathbf{P}^n$  is any row of the matrix  $\mathbf{P}^n$  as  $n$  approaches  $\infty$ . All the rows are the same in this matrix at the limit. The latter signifies the fact that the limit  $\mathbf{\Pi}$  is independent of the initial distribution. In other words, after the Markov chain runs for a long time, it forgets its initial distribution and converges to  $\mathbf{\Pi}$ .

We denote by  $\pi_j$ ,  $j = 0, 1, 2, \dots$ , the components of the vector  $\mathbf{\Pi}$ . That is,  $\pi_j$  is the steady state probability of the Markov chain to be at state  $j$ . Namely,

$$\pi_j = \lim_{n \rightarrow \infty} \pi_n(j) \quad \text{for all } j. \quad (185)$$

By equation (177), we obtain

$$\pi_n(j) = \sum_{i=0}^{\infty} P_{ij} \pi_{n-1}(i), \quad (186)$$

then by the latter and (185), we obtain

$$\pi_j = \sum_{i=0}^{\infty} P_{ij} \pi_i. \quad (187)$$

Therefore, recalling that  $\pi$  is a proper probability distribution, we can conclude that for an irreducible, aperiodic and stable Markov chain, the steady state probabilities can be obtained by solving the following steady state equations:

$$\pi_j = \sum_{i=1}^{\infty} \pi_i P_{ij} \quad \text{for all } j, \quad (188)$$

$$\sum_{j=1}^{\infty} \pi_j = 1 \quad (189)$$

and

$$\pi_j \geq 0 \quad \text{for all } j. \quad (190)$$

In this case:

$$\pi_j = \frac{1}{E[R_j]}. \quad (191)$$

When the state space  $\Theta$  is finite, one of the equations in (188) is redundant and replaced by (189).

In matrix notation equation (188) is written as:  $\mathbf{\Pi} = \mathbf{\Pi P}$ .

Note that if we consider an irreducible, aperiodic and stable Markov chain, then also a unique non-negative steady state solution vector  $\mathbf{\Pi}$  of the steady equation (188) exists. However, in this case, the  $j$ th component of  $\mathbf{\Pi}$ , namely  $\pi_j$ , is not a probability but it is the proportion of time in steady state that the Markov chain is in state  $j$ .

Note also that the steady state vector  $\mathbf{\Pi}$  is called the stationary distribution of the Markov chain, because if we set  $\mathbf{\Pi}_0 = \mathbf{\Pi}$ ,  $\mathbf{\Pi}_1 = \mathbf{\Pi P} = \mathbf{\Pi}$ ,  $\mathbf{\Pi}_2 = \mathbf{\Pi P} = \mathbf{\Pi}$ ,  $\dots$ , i.e.,  $\mathbf{\Pi}_n = \mathbf{\Pi}$  for all  $n$ .

## Homework 2.17

1. Show that a discrete-time Markov chain (MC) with two states where the rows of the transition probability matrix are identical is a Bernoulli process.

2. Prove that in any finite MC, at least one state must be recurrent.
3. Provide examples of MCs defined by their transition probability matrices that their states (or some of the states) are periodic, aperiodic, transient, null recurrent and positive recurrent. Provide examples of irreducible and reducible (not irreducible) and of stable and unstable MCs. You may use as many MCs as you wish to demonstrate the different concepts.
4. For different  $n$  values, choose an  $n \times n$  transition probability matrix  $\mathbf{P}$  and an initial vector  $\mathbf{\Pi}_0$ . Write a program to compute  $\mathbf{\Pi}_1, \mathbf{\Pi}_2, \mathbf{\Pi}_3, \dots$  and demonstrate convergence to a limit in some cases and demonstrate that the limit does not exist in other cases.
5. Prove equation (191).
6. Consider a binary communication channel between a transmitter and a receiver where  $B_n$  is the value of the  $n$ th bit at the receiver. This value can be either equal to 0, or equal to 1. Assume that the event [a bit to be erroneous] is independent of the value received and only depends on whether of not the previous bit is erroneous or correct. Assume the following:

$$P(B_{n+1} \text{ is erroneous} \mid B_n \text{ is correct}) = 0.0001$$

$$P(B_{n+1} \text{ is erroneous} \mid B_n \text{ is erroneous}) = 0.01$$

$$P(B_{n+1} = 1 \text{ is correct} \mid B_n \text{ is correct}) = 0.9999$$

$$P(B_{n+1} = 0 \text{ is correct} \mid B_n \text{ is erroneous}) = 0.99$$

Compute the steady state error probability.  $\square$

### 2.4.7 Reversibility

Consider an irreducible, aperiodic and stable Markov chain  $\{X_n\}$ . Assume that this Markov chain has been running for a long time to achieve stationarity with transition probability matrix  $\mathbf{P} = [P_{ij}]$ , and consider the process  $X_n, X_{n-1}, X_{n-2}, \dots$ , going back in time. This reversed process is also a Markov chain because  $X_n$  has dependence relationship only with  $X_{n-1}$  and  $X_{n+1}$  and conditional on  $X_{n+1}$ , it is independent of  $X_{n+2}, X_{n+3}, X_{n+4}, \dots$ . Therefore,

$$P(X_{n-1} = j \mid X_n = i) = P(X_{n-1} = j \mid X_n = i, X_{n+1} = i_{n+1}, X_{n+2} = i_{n+2}, \dots).$$

In the following we derive the transition probability matrix, denoted  $\mathbf{Q} = [Q_{ij}]$  of the process  $\{X_n\}$  in reverse. Accordingly Define

$$Q_{ij} = P(X_n = j \mid X_{n+1} = i). \quad (192)$$

By the definition of conditional probability, we obtain,

$$Q_{ij} = \frac{P(X_n = j \cap X_{n+1} = i)}{P(X_{n+1} = i)} \quad (193)$$

or

$$Q_{ij} = \frac{P(X_n = j)P(X_{n+1} = i \mid X_n = j)}{P(X_{n+1} = i)} \quad (194)$$

and if  $\pi_j$  denotes the steady state probability of the Markov chain  $\{X_n\}$  to be in state  $j$ , and let  $n \rightarrow \infty$ , we obtain

$$Q_{ij} = \frac{\pi_j P_{ji}}{\pi_i}. \quad (195)$$

A Markov chain is said to be *time reversible* if  $Q_{ij} = P_{ij}$  for all  $i$  and  $j$ . Substituting  $Q_{ij} = P_{ij}$  in (195), we obtain,

$$\pi_i P_{ij} = \pi_j P_{ji}. \quad (196)$$

Eq. (196) is also a condition for time reversibility.

Intuitively, a Markov chain  $X_n$  is time-reversible if for a large  $k$  (to ensure stationarity) the Markov chain  $X_k, X_{k+1}, X_{k+2} \dots$  is statistically the same as the process  $X_k, X_{k-1}, X_{k-2} \dots$ . In other words, by observing the two processes, you cannot tell which one is going forward and which is going backward.

### Homework 2.18

Provide an example of a Markov chain that is time reversible and another one that is not time reversible.  $\square$

### 2.4.8 Multi-Dimensional Markov Chains

So far, we discussed single dimensional Markov chains. If the state space is made of finite vectors instead of scalars, we can easily convert them to scalars and proceed with the above described approach. For example, if the state space is (0,0) (0,1) (1,0) (1,1) we can simply change the names of the states to 0,1,2,3 by assigning the values 0, 1, 2 and 3 to the states (0,0), (0,1), (1,0) and (1,1), respectively. In fact we do not even have to do it explicitly. All we need to do is to consider a  $4 \times 4$  transition probability matrix as if we have a single dimension Markov chain. Let us now consider an example of a multidimensional Markov chain.

Consider a bit stream transmitted through a channel. Let  $Y_n = 1$ , if the  $n$ th bit is received correctly, and let  $Y_n = 0$  if the  $n$ th bit is received incorrectly. Assume the following

$$\begin{aligned} P(Y_n = i_n \mid Y_{n-1} = i_{n-1}, Y_{n-2} = i_{n-2}) \\ = P(Y_n = i_n \mid Y_{n-1} = i_{n-1}, Y_{n-2} = i_{n-2}, Y_{n-3} = i_{n-3}, Y_{n-4} = i_{n-4}, \dots). \end{aligned}$$

$$\begin{aligned} P(Y_n = 0 \mid Y_{n-1} = 0, Y_{n-2} = 0) &= 0.9 \\ P(Y_n = 0 \mid Y_{n-1} = 0, Y_{n-2} = 1) &= 0.7 \\ P(Y_n = 0 \mid Y_{n-1} = 1, Y_{n-2} = 0) &= 0.6 \\ P(Y_n = 0 \mid Y_{n-1} = 1, Y_{n-2} = 1) &= 0.001. \end{aligned}$$

By the context of the problem, we have

$$P(Y_n = 1) = 1 - P(Y_n = 0)$$

so,

$$\begin{aligned} P(Y_n = 1 \mid Y_{n-1} = 0, Y_{n-2} = 0) &= 0.1 \\ P(Y_n = 1 \mid Y_{n-1} = 0, Y_{n-2} = 1) &= 0.3 \\ P(Y_n = 1 \mid Y_{n-1} = 1, Y_{n-2} = 0) &= 0.4 \\ P(Y_n = 1 \mid Y_{n-1} = 1, Y_{n-2} = 1) &= 0.999. \end{aligned}$$

**Homework 2.19**

Explain why the process  $\{Y_n\}$  is not a Markov chain.  $\square$

Now define the  $\{X_n\}$  process as follows:

$X_n = 0$  if  $Y_n = 0$  and  $Y_{n-1} = 0$ .

$X_n = 1$  if  $Y_n = 0$  and  $Y_{n-1} = 1$ .

$X_n = 2$  if  $Y_n = 1$  and  $Y_{n-1} = 0$ .

$X_n = 3$  if  $Y_n = 1$  and  $Y_{n-1} = 1$ .

**Homework 2.20**

Explain why the process  $\{X_n\}$  is a Markov chain, produce its transition probability matrix, and compute its steady state probabilities.  $\square$

**2.5 Continuous Time Markov Chains****2.5.1 Definitions and Preliminaries**

A continuous time Markov chain is a continuous time stochastic process  $\{X_t\}$ . At any point in time  $t$ ,  $\{X_t\}$  describes the state of the process which is discrete. We will consider only continuous time Markov chain where  $X_t$  takes values that are nonnegative integer. The time between changes in the state of the process is exponentially distributed. In other words, the process stays constant for an exponential time duration before changing to another state.

In general, a continuous time Markov chain  $\{X_t\}$  is defined by the property that for all real numbers  $s \geq 0$ ,  $t \geq 0$  and  $0 \leq v < s$ , and integers  $i \geq 0$ ,  $j \geq 0$  and  $k \geq 0$ ,

$$P(X_{t+s} = j \mid X_t = i, X_v = k_v, v \leq t) = P(X_{t+s} = j \mid X_t = i). \quad (197)$$

That is, the probability distribution of the future values of the process  $X_t$ , represented by  $X_{t+s}$ , given the present value of  $X_t$  and the past values of  $X_t$  denoted  $X_v$ , is independent of the past and depends only on the present.

A general continuous time Markov chain can also be defined as a continuous time discrete space stochastic process with the following properties.

1. Each time the process enters state  $i$ , it stays at that state for an amount of time which is exponentially distributed with parameter  $\delta_i$  before making a transition into a different state.
2. When the process leaves state  $i$ , it enters state  $j$  with probability denoted  $P_{ij}$ . The set of  $P_{ij}$ s must satisfy the following:

$$(1) \quad P_{ii} = 0 \quad \text{for all } i$$

$$(2) \quad \sum_j P_{ij} = 1.$$



### 2.5.2 Examples

An example of a continuous time Markov chain is a Poisson process with rate  $\lambda$ . The state at time  $t$ ,  $\{X_t\}$  can be the number of occurrences by time  $t$  which is the counting process  $N(t)$ . In this example of the Poisson counting process  $\{X_t\} = N(t)$  increases by one after every exponential time duration with parameter  $\lambda$ .

Another example is the so-called *pure birth process*  $\{X_t\}$ . It is a generalization of the counting Poisson process. Again  $\{X_t\}$  increases by one every exponential amount of time but here, instead of having a fixed parameter  $\lambda$  for each of these exponential intervals, this parameter depends of the state of the process and it is denoted  $\delta_i$ . In other words, when  $\{X_t\} = i$ , the time until the next occurrence in which  $\{X_t\}$  increases from  $i$  to  $i + 1$  is exponentially distributed with parameter  $\delta_i$ . If we set  $\delta_i = \lambda$  for all  $i$ , we have the Poisson counting process.

### 2.5.3 Birth and Death Process

In many real life applications, the state of the system sometimes increases by one, and at other times decreases by one. A continuous time Markov chain  $\{X_t\}$  that often models such a system is called a *birth and death process*. In such a process, the time between occurrences in state  $i$  is exponentially distributed, with parameter  $\delta_i$ , and at any point of occurrence, the process increases by one (from its previous value  $i$  to  $i + 1$ ) with probability  $v_i$  and decreases by one (from  $i$  to  $i - 1$ ) with probability  $\vartheta_i = 1 - v_i$ . The transitions from  $i$  to  $i + 1$  are called *births* and the transitions from  $i$  to  $i - 1$  are called *deaths*. Recall that the mean time between occurrences, when in state  $i$ , is  $1/\delta_i$ . Hence, the birth rate in state  $i$ , denoted  $b_i$ , is given by

$$b_i = \delta_i v_i$$

and the death rate ( $d_i$ ) is given by

$$d_i = \delta_i \vartheta_i.$$

Summing up these two equations gives the intuitive result that the total rate at state  $i$  is equal to the sum of the birth and death rates. Namely,

$$\delta_i = b_i + d_i$$

and therefore the mean time between occurrences is

$$\frac{1}{\delta_i} = \frac{1}{b_i + d_i}.$$

### Homework 2.21

Show the following:

$$\vartheta_i = \frac{d_i}{b_i + d_i}$$

and

$$v_i = \frac{b_i}{b_i + d_i}. \quad \square$$

Birth and death processes apply to queueing systems where customers arrive one at a time and depart one at a time. Consider for example a birth and death process with the death rate higher than the birth rate. Such a process could model, for example, a stable single-server queueing system.

### 2.5.4 First Passage Time

An important problem that has applications in many fields, such as biology, finance and engineering, is how to derive the distribution or moments of the time it takes for the process to transit from state  $i$  to state  $j$ . In other words, given that the process is in state  $i$  find the distribution of a random variable representing the time it takes to enter state  $j$  for the first time. This random variable is called the *first passage time from  $i$  to  $j$* . Let us derive the mean of the first passage time from  $i$  to  $j$  in a birth and death process for the case  $i < j$ . To solve this problem we start with a simpler one. Let  $U_i$  be the mean passage time to go from  $i$  to  $i + 1$ . Then

$$U_0 = \frac{1}{b_0}. \quad (198)$$

and

$$U_i = \frac{1}{\delta_i} + \vartheta_i[U_{i-1} + U_i]. \quad (199)$$

### Homework 2.22

Explain equations (198) and (199).  $\square$

Therefore,

$$U_i = \frac{1}{b_i + d_i} + \frac{d_i}{b_i + d_i}[U_{i-1} + U_i] \quad (200)$$

or

$$U_i = \frac{1}{b_i} + \frac{d_i}{b_i}U_{i-1}. \quad (201)$$

Now we have a recursion by which we can obtain  $U_0, U_1, U_2, \dots$ , and the mean first passage time between  $i$  and  $j$  is given by the sum

$$\sum_{k=i}^j U_k.$$

### Homework 2.23

Let  $b_i = \lambda$  and  $d_i = \mu$  for all  $i$ , derive a closed form expression for  $U_i$ .  $\square$

### 2.5.5 Transition Probability Function

Define the *transition probability function*  $P_{ij}(t)$  as the probability that given that the process is in state  $i$  at time  $t_0$ , then a time  $t$  later, it will be in state  $j$ . That is,

$$P_{ij}(t) = P[X(t_0 + t) = j \mid X(t_0) = i]. \quad (202)$$

The continuous time version of the Chapman-Kolmogorov equations are

$$P_{ij}(t + \tau) = \sum_{n=0}^{\infty} P_{in}(t)P_{nj}(\tau) \quad \text{for all } t \geq 0, \tau \geq 0. \quad (203)$$

Using the latter to derive the limit

$$\lim_{\Delta t \rightarrow 0} \frac{P_{ij}(t + \Delta t) - P_{ij}(t)}{\Delta t}$$

we obtain the so called Kolmogorov's Backward Equations:

$$P'_{ij}(t) = \sum_{n \neq i} \delta_i P_{in} P_{nj}(t) - \delta_i P_{ij}(t) \quad \text{for all } i, j \text{ and } t \geq 0. \quad (204)$$

For a birth and death process the latter become

$$P'_{0j}(t) = b_0 \{P_{1j}(t) - P_{0j}(t)\}. \quad (205)$$

and

$$P'_{ij}(t) = b_i P_{i+1,j}(t) + d_i P_{i-1,j}(t) - (b_i + d_i) P_{ij}(t) \quad \text{for all } i > 0. \quad (206)$$

### 2.5.6 Steady State Probabilities

As in the case of the discrete time Markov chain, define a continuous-time Markov chain to be called *irreducible* if there is a positive probability for any state to reach every state, and we define a continuous-time Markov chain to be called *positive recurrent* if the process starts from any state, the random variable that represents the time it returns to that state has finite mean. As in the case of Discrete time Markov chain a Markov chain is said to be *stable* if all its states are positive recurrent.

Henceforth we only consider a continuous-time Markov chains that are irreducible, aperiodic and stable. Then the limit of  $P_{ij}(t)$  as  $t$  approaches infinity exists, and we define

$$\pi_j = \lim_{t \rightarrow \infty} P_{ij}(t). \quad (207)$$

The  $\pi_j$  values are called steady state probabilities or stationary probabilities of the continuous-time Markov chain. In particular,  $\pi_j$  is the steady state probability of the continuous-time Markov chain to be at state  $j$ . We shall now describe how the steady-state probabilities  $\pi_j$ s can be obtained.

We now construct the matrix  $\mathbf{Q}$  which is called the *infinitesimal generator* of the continuous time Markov chain. The matrix  $\mathbf{Q}$  is a matrix of one step infinitesimal rates  $Q_{ij}$  defined by

$$Q_{ij} = \delta_i P_{ij} \quad \text{for } i \neq j \quad (208)$$

and

$$Q_{ii} = - \sum_{j \neq i} Q_{ij}. \quad (209)$$

**Remarks:**

- The state space can be finite or infinite and hence the matrices  $\mathbf{P}$  and  $\mathbf{Q}$  can also be finite or infinite.
- In Eq. (208),  $Q_{ij}$  is the product of the rate to leave state  $i$  and the probability of transition to state  $j$  from state  $i$  which is the rate of transitions from  $i$  to  $j$ .

To obtain the steady state probabilities  $\pi_j$ s, we solve the following set of steady state equations:

$$0 = \sum_i \pi_i Q_{ij} \quad \text{for all } j \quad (210)$$

and

$$\sum_j \pi_j = 1. \quad (211)$$

Denoting  $\mathbf{\Pi} = [\pi_0, \pi_1, \pi_2, \dots]$ , Eq. (210) can be written as

$$0 = \mathbf{\Pi Q}. \quad (212)$$

To explain Eqs. (210), notice that, by (208) and (209), for a particular  $j$ , the equation

$$0 = \sum_i \pi_i Q_{ij} \quad \text{for all } j \quad (213)$$

is equivalent to

$$\pi_j \sum_{i \neq j} Q_{ji} = \sum_{i \neq j} \pi_i Q_{ij} \quad (214)$$

or

$$\pi_j \sum_{i \neq j} \delta_j P_{ji} = \sum_{i \neq j} \pi_i \delta_i P_{ij}. \quad (215)$$

The quantity  $\pi_i Q_{ij}$  which is the steady state probability of being in state  $i$  times the infinitesimal rate of a transition from state  $i$  to state  $j$  is called the *probability flux* from state  $i$  to state  $j$ . Eq. (213) says that the total probability flux from all states into state  $j$  is equal to the total probability flux out of state  $j$  to all other states. To explain this equality, consider a long period of time  $L$ . Assuming the process return to all states infinitely many times, during long time period  $L$ , the number of times the process moves into state  $j$  is equal (in the limit  $L \rightarrow \infty$ ) to the number of times the process moves out of state  $j$ . This leads to Eq. (215) with the factor  $L$  in both sides.

Similar to the case of discrete time Markov chains, the set of equations (210) and (211) is dependent and one of the equations in (210) is redundant in the finite state space case.

Due to the fact that the  $\mathbf{Q}$  matrix may be too large, it may not be possible to solve the steady state equations (210). Actually, the case of a large state space (or large  $\mathbf{Q}$  matrix) is common in practice. Consider for example a 49 cell GSM mobile network, and assume that every cell has 23 voice channels. Assuming Poisson arrivals and exponential holding and cell sojourn times. Then this cellular mobile network can be modeled as a continuous time Markov chain with each state representing the number of busy channels in each cell. In this case, the number of states is equal to  $24^{49}$ , so a numerical solution of the steady state equations is not possible.

### 2.5.7 Simulations

When a numerical solution is not possible, we need to rely on simulations. Fortunately, the special structure of the continuous time Markov chain together with a certain property of the Poisson process called PASTA (Poisson Arrivals See Time Averages), simulations of continuous time Markov chain models can be simplified and expedited so they lead to accurate results. To explain the PASTA property, consider a stochastic process for which steady state probabilities exist. If we are interested in obtaining certain steady state statistical characteristics of the process (like the  $\pi_i$ s in a continuous time Markov chains), we could inspect the entire evolution of the process (in practice, for a long enough time period), or we could use an independent Poisson inspector. (We already discussed the unique property of the Poisson process to see time averages.) The PASTA principle means that we do not need a separate Poisson inspector, but we could inspect the process at occurrences of any given independent Poisson process which is part of the continuous Markov chain. Note that in practice, since we are limited to a finite number of inspections, we should choose a Poisson process that will have sufficient number of occurrences (inspections) during the simulation of the stochastic process we are interested in obtaining its steady state statistics.

In many cases, when we are interested in steady-state statistics of a continuous time Markov chain, we can conveniently find a Poisson process which is part of the continuous time Markov chain we are interested in and use it as a Poisson inspector. For example, if we consider a queueing system in which the arrival process follows a Poisson process, such process could be used for times of arrivals of the inspector if it, at any inspection, does not count (include) its own particular arrival. In other words, we consider a Poisson inspector that arrives just before its own arrival occurrences.

### 2.5.8 Reversibility

We have discussed the **time reversibility** concept in the context of discrete-time Markov chain. In the case of continuous time Markov chain the notion of time reversibility is similar. If you observe the process  $X_t$  for a large  $t$  (to ensure stationarity) and if you cannot tell from its statistical behavior if it is going forward or backward, it is time reversible.

Consider continuous-time Markov chain that has a unique steady-state solution and that its  $[P_{ij}]$  matrix would give rise to a discrete-time Markov chain. This discrete-time Markov chain, called the *embedded chain* of our continuous-time Markov chain, has  $[P_{ij}]$  as its transition probability matrix. This embedded chain is in fact the sequence of states that our original continuous time chain visits where we ignore the time spent in each state during each visit to that state. We already know the condition for time reversibility of the embedded chain, so consider our continuous time chain and assume that it has been running for a long while, and consider its reversed process going backwards in time. In the following we show that also the reversed process spends an exponentially distributed amount of time in each state. Moreover, we will show that as the original process, the reverse process spends an exponentially distributed amount of time with parameter  $\delta_i$  when in state  $i$ .

$$P\{X(t) = i, \text{ for } t \in [u - v, u] \mid X(u) = i\} = \frac{P\{X(t) = i, \text{ for } t \in [u - v, u] \mid X(u) = i\}}{P[X(u) = i]}$$

$$= \frac{P[X(u-v) = i]e^{-\delta_i v}}{P[X(u) = i]} = e^{-\delta_i v}.$$

The last equality is explained by reminding the reader that the process is in steady-state so the probability that the process is in state  $i$  at time  $(u-v)$  is equal to the probability that the process is in state  $i$  at time  $u$ .

Since the continuous-time Markov chain is composed of two parts, its embedded chain and the time spent in each state, and since we have shown that the reversed process spends time in each state which is statistically the same as the original process, a condition for time reversibility of continuous-time Markov chain is that its embedded chain is time reversible. As we have learned when we discussed reversibility of discrete-time Markov chain, the condition is

$$\hat{\pi}_i P_{ij} = \hat{\pi}_j P_{ji} \text{ for all } i, j. \quad (216)$$

This is equivalent to

$$\pi_i Q_{ij} = \pi_j Q_{ji} \text{ for all } i, j. \quad (217)$$

## Homework 2.24

Derive (217) from (216).  $\square$

It is important to notice that a birth and death process that its embedded chain is time reversible. Consider a very long time  $L$  during that time, the number of transitions from state  $i$  to state  $i+1$ , denoted  $T_{i,i+1}(L)$ , is equal to the number of transitions, denoted  $T_{i+1,i}(L)$ , from state  $i+1$  to  $i$  because every transition from  $i$  to  $i+1$  must eventually follow by a transition from  $i+1$  to  $i$ . Actually, there may be a last transition from  $i$  to  $i+1$  without the corresponding return from  $i+1$  to  $i$ , but since we assume that  $L$  is arbitrarily large, the number of transitions is arbitrarily large and being off by one for an arbitrarily large number is negligible.

Therefore, for arbitrary large  $L$ ,

$$\frac{T_{i,i+1}(L)}{L} = \frac{T_{i+1,i}(L)}{L}. \quad (218)$$

Since for a birth and death process  $Q_{ij} = 0$  for  $|i-j| > 1$  and for  $i = j$ , and since for arbitrarily large  $L$ , we have

$$\pi_i Q_{i,i+1} = \frac{T_{i,i+1}(L)}{L} = \frac{T_{i+1,i}(L)}{L} = \pi_{i+1} Q_{i+1,i}, \quad (219)$$

so our birth and death process is time reversible. This is an important result for the present context because many of our queueing models are special cases of the birth and death process.

## 2.5.9 Multi-Dimensional Continuous Time Markov Chains

The extension discussed earlier to multi-dimensional discrete time Markov chain applies also to the case of continuous time Markov chain. If the state space is made of finite vectors instead of scalars, as discussed, there is a one-to-one correspondence between vectors and scalars, so multi-dimensional continuous-time Markov chain can be converted to a single-dimension continuous-time Markov chain and we proceed with the above described approach that applies to the single dimension.

### 3 General Queueing Concepts

In general, queueing systems may be characterized by complex input process, service distribution and queue disciplines. In practice, such queueing processes and disciplines are often not amenable to analysis. Nevertheless, insight can be often gained using simpler queueing models. This modelling simplification is commonly made in the context of packet switching networks like the Internet that are based on the store and forward principle. Typically, packets on their ways to their destinations arrive at a router where they are stored and further forwarded according to addresses in their headers. One of the most fundamental elements in this process is the single-server queue. One of the aims of telecommunications research is to explain traffic and management processes and their effect on queueing performance. In this section, we briefly cover basic queueing theory concepts. We shall bypass mathematically rigorous proofs and rely instead on simpler intuitive explanations.

#### 3.1 Notation

A commonly used shorthand notation, called Kendall notation [32], for such single queue models describes the arrival process, service distribution, the number of servers and the buffer size (waiting room) as follows:

arrival process / service distribution / number of servers / waiting room

Commonly used characters for the first two positions in this shorthand notation are: D (Deterministic), M (Markovian - Poisson for the arrival process or Exponential for the service time), G (General), GI (General and independent), and Geom (Geometric). The fourth position is used for the number of buffer places in addition to the number of servers and it is usually not used if the waiting room is unlimited.

For example, M/M/1 denotes a single-server queue with Poisson arrival process and exponential service time with infinite buffer. M/G/k/k denotes a k server queue with no additional waiting room except at the servers with the arrival process being Poisson.

#### 3.2 Utilization

An important measure for queueing systems performance is the utilization, denoted  $U$ . It is the proportion of time that a server is busy on average. In many systems, the server is paid for its time regardless if it is busy or not. Normally, the time that transmission capacity is not used is time during which money is spent but no revenue is collected. It is therefore important to design systems that will maintain high utilization.

If you have two identical servers and one is busy 0.4 of the time and the other 0.6. Then the utilization is 0.5. We always have that  $0 \leq U \leq 1$ . If we consider an M/M/ $\infty$  queue (Poisson arrivals, exponentially distributed service times and infinite servers) and the arrival rate is finite, the utilization is zero because the mean number of busy servers is finite and the mean number of idle servers is infinite.

Consider a G/G/1 queue (that is, a single-server queue with arbitrary arrival process and arbitrary service time distribution, with infinite buffer), with  $\lambda$  the mean arrival rate and  $\mu$  the mean service rate. Assume that  $\mu > \lambda$  so that the queue is *stable*, namely, that it will not keep growing forever, and that whenever it is busy, eventually it will reach the state where the system is empty. For a stable G/G/1 queue, we have that  $U = \lambda/\mu$ . To show the latter let  $L$  be a very long period of time. The average number of customers (amount of work) arrived within time period  $L$  is:  $\lambda L$ . The average number of customers (amount of work) that has been served during time period  $L$  is equal to  $\mu U L$ . Since  $L$  is large and the queue is stable, these two values are equal. Thus,  $\mu U L = \lambda L$ . Hence,  $U = \lambda/\mu$ .

Often, we are interested in the distribution of the number (of customers, jobs or packets) in the system. Let  $p_n$  be the probability that there are  $n$  in the system. Having the utilization, we can readily obtain  $p_0$  the probability that the system is empty. For the G/G/1 queue we have,

$$p_0 = 1 - U = 1 - \lambda/\mu. \quad (220)$$

If we have a multi-server queue, e.g. G/G/ $k/k + n$ , then the utilization will be defined as the overall average utilization of the individual servers. That is, each server will have its own utilization defined by the proportion of time it is busy, and the utilization of the entire multi-server system will be the average of the individual server utilization.

### 3.3 Little's Formula

Another important and simple queueing theory result that applies to G/G/1 queue (and to other systems) is known as *Little's Formula* [40, 56, 57]. It has two forms. The first form is:

$$E[Q] = \lambda E[D] \quad (221)$$

where  $E[Q]$  and  $E[D]$  represent the stationary mean queue-size including the customer in service and the mean delay (system waiting time) of a customer from the moment it arrives until its service is complete, respectively. In remainder of this book, when use terms such as *mean queue-size* and *mean delay* we mean to refer to their values in steady state, meaning, stationary mean queue size and delay, respectively.

The second form is:

$$E[N_Q] = \lambda E[W_Q] \quad (222)$$

where  $E[N_Q]$  and  $E[W_Q]$  represent the mean number of customers in the queue in steady-state excluding the customer in service and the mean delay of a customer, in steady-state, from the moment it arrives until its service commences (waiting time in the queue), respectively.

An intuitive (non-rigorous) way to explain Eq. (221) is by considering a customer that just left the system (completed service). This customer sees behind his/her back on average  $E[Q]$  customers. Who are these customers? They are the customers that had been arriving during the time that our customer was in the system. Their average number is  $\lambda E[D]$ .

To obtain (222) from (221), notice that

$$E[Q] = E[N_Q] + U = E[N_Q] + \lambda/\mu \quad (223)$$



and

$$E[D] = E[W_Q] + 1/\mu. \quad (224)$$

Substituting (223) and (224) in (221), (222) follows.

For a graphical proof of Little's Formula for the case of G/G/1 queue see [9]. The arguments there may be summarized as follows. Consider a stable G/G/1 queue that starts at time  $t = 0$  with an empty queue. Let  $A(t)$  be the number of arrivals up to time  $t$ , and let  $D(t)$  be the number of departures up to time  $t$ . The queue-size (number in the system) at time  $t$  is denoted  $Q(t)$  and is given by  $Q(t) = A(t) - D(t)$ ,  $t \geq 0$ . Let  $L$  be an arbitrarily long period of time. Then the mean queue-size  $E[Q]$  is given by

$$E[Q] = \frac{1}{L} \int_0^L Q(t) dt. \quad (225)$$

Also notice that

$$\int_0^L Q(t) dt = \sum_{i=1}^{A(L)} W_i \quad (226)$$

where  $W_i$  is the time spent in the system by the  $i$ th customer. (Notice that since  $L$  is arbitrarily large, there have been arbitrarily large number of events during  $[0, L]$  where our stable G/G/1 queue became empty, so  $A(L) = D(L)$ .) Therefore,

$$\frac{1}{L} \int_0^L Q(t) dt = \frac{1}{L} \sum_{i=1}^{A(L)} W_i \quad (227)$$

and realizing that

$$\lambda = A(L)/L, \quad (228)$$

and

$$E[D] = \frac{1}{A(L)} \sum_{i=1}^{A(L)} W_i, \quad (229)$$

we obtain

$$E[Q] = \frac{1}{L} \int_0^L Q(t) dt = \frac{A(L)}{L} \frac{1}{A(L)} \sum_{i=1}^{A(L)} W_i = \lambda E[D]. \quad (230)$$

Little's formula applies to many systems. Its applicability is not limited to single-server queues, or single queue systems, or systems with infinite buffer.

Interestingly, the result  $U = \lambda/\mu$  for a G/G/1 queue can also be obtained using Little's formula. Let us consider a system to be just the server (excluding the infinite buffer). The mean time a customer spends in this system is  $1/\mu$  because this is the mean service time so this is the mean time spent in the system that includes just the server. The mean arrival rate into that system must be equal to  $\lambda$  because all the customers that arrive at the queue eventually arrive at the server - nothing is lost. Let us now consider the mean number of customers at the server, denoted  $N_s$ . Clearly,  $N_s$  can only take the values zero or one, because no more than one customer can be at the server at any point in time. We also know that the steady state probability  $P(N_s = 0)$  is equal to  $\pi_0$ . Therefore,

$$N_s = 0\pi_0 + 1(1 - \pi_0) = 1 - \pi_0 = U.$$

By Little's formula, we have

$$N_s = \lambda(1/\mu),$$

so

$$U = \lambda/\mu.$$

Another interesting application of Little's formula relates the blocking probability  $P_b$  of a G/G/1/ $k$  queue (a G/G/1 queue with a buffer of size  $k$ ) with its server utilization [27, 50]. Again, consider the server as an independent system. Since the mean number of customers in this system is  $U$ , and the arrival rate into this system is  $(1 - P_b)\lambda$ , we obtain by Little's formula:

$$U = (1 - P_b)\lambda\mu^{-1}, \quad (231)$$

where  $\mu^{-1}$  is the mean service time. Let  $\rho = \lambda/\mu$ , we obtain

$$P_b = 1 - \frac{U}{\rho}. \quad (232)$$

### 3.4 Work Conservation

Another important concept in queueing theory is the concept of *work conservation*. A queueing system is said to be work conservative if no server is idle if there is work to be done. For example, G/G/1 and G/G/1/ $k$  are work conservative. However, G/G/ $k$  and G/G/ $k/n$  ( $n \geq k$ ) are not work conservative because a server can be idle while there are customers waiting and/or being served by other servers.

### 3.5 PASTA

Many of the queueing models we consider in this book involve Poisson arrival processes. The Poisson Arrivals See Time Averages (PASTA) property discussed in the previous section is important for analysis and simulations of such queueing models. Let us further explain and prove this important property.

The PASTA property means that arriving customers in steady state will find the number of customers in the system obeying its steady state distribution. In other words, the statistical characteristics (e.g., mean, variance, distribution) of the number of customers in the system observed by an arrival is the same as those observed by an independent Poisson inspector. This is not true in general. Consider the *lonely person* example of a person lives alone and never has another person comes to his/her house. When this person comes home s/he always finds that there are no people in the house upon its arrival, but if we use a an independent Poisson inspector to evaluate the proportion of time that person is in the house, the inspector will find sometimes that there is one person in the house and in other times that there is no-one in the house. Of course, the arrival process of this person is not a Poisson process as there are no arrivals during the time the person is in the house.

In addition to the Poisson arrival assumption, for PASTA to be valid we also need the condition that arrivals after time  $t$  are independent of the queue size at time  $t$ ,  $Q(t)$ . For example, if we have a single-server queue (SSQ) with Poisson arrivals and the service times have the property

that the service of a customer must always terminate before the next arrival, then the arrivals always see an empty queue, and, of course, an independent arrival does not.

To prove PASTA we consider the limit

$$A_k(t) = \lim_{\Delta t \rightarrow 0} P[Q(t) = k \mid \text{an arrival occurs within } (t, t + \Delta t)].$$

Using Bayes' formula and the condition that arrivals after time  $t$  are independent of  $Q(t)$ , we obtain that

$$A_k(t) = P[Q(t) = k]. \quad (233)$$

Then, by taking the limit of both sides of (233), we complete the proof that the queue size seen by an arrival is statistically identical to the queue size seen by an independent observer.  $\square$

### Homework 3.1

Prove Eq. (233).  $\square$

## 3.6 Queueing Models

In this book we discuss various queueing models that are amenable to analysis. The analysis is simplest for D/D/ type queues where the interarrival and service times are deterministic (fixed values). They will be discussed in the next section. Afterwards, we will consider the so-called Markovian queues. These queues are characterized by the Poisson arrival process, independent exponential service times and independence between the arrival process and the service times. They are denoted by M in the first two positions (i.e., M/M/ · / ·). Because of the memoryless property of Markovian queues, these queues are amenable to analysis. In fact, they are all continuous time Markov chains with the state being the *queue-size* defined as the number in the system  $n$  and the time between state transitions is exponential. The reason that these time periods are exponential is that at any point in time, the remaining time until the next arrival, or the next service completion, is a competition between various exponential random variables.

## 4 Simulations

In many cases, such analytical solutions are not available, so simulations are used to estimate performance measures. Simulations are also used to evaluate accuracy of analytical approximations (e.g. results obtained by the Erlang Fixed-Point Approximation).

### 4.1 Confidence Intervals

Regardless of how long we run a simulation, we will never obtain the exact mathematical result of a steady state measure we are interested in. To assess the error of our simulation, we begin by running a certain number, say  $n$ , of simulation experiments and obtain  $n$  observed values, denoted  $a_1, a_2, \dots, a_n$ , of the measure of interest.

Let  $\bar{a}$  be the observed mean and  $\sigma_a^2$  the observed variance of these  $n$  observations. Their values are given by

$$\bar{a} = \frac{1}{n} \sum_{i=1}^n a_i \quad (234)$$

and

$$\sigma_a^2 = \frac{1}{n-1} \sum_{i=1}^n (a_i - \bar{a})^2. \quad (235)$$

Then the confidence interval of  $\bar{a}$ , with confidence  $\alpha$ ,  $0 \leq \alpha \leq 1$ , is given by  $(\bar{a} - U, \bar{a} + U)$ , where

$$U = \{t_{(1-\alpha)/2, (n-1)}\} \frac{\sigma_a}{\sqrt{n}} \quad (236)$$

where  $t_{(1-\alpha)/2, (n-1)}$  is the appropriate percentage point for Student t-distribution with  $n-1$  degrees of freedom. The  $t_{(1-\alpha)/2, (n-1)}$  values are available in standard tables. For example:  $t_{0.025, 5} = 2.57$  and  $t_{0.025, 10} = 2.23$ . This means that if we are interested in 95% confidence and we have  $n = 6$  observations, we will use  $t_{0.025, 5} = 2.57$  to obtain the confidence interval.

We use here the Student t-distribution (and not Gaussian) because it is the right distribution to use when we attempt to estimate the mean of a population which is normally distributed when we have a small sample size. In fact, the need to estimate such mean based on a small sample gave size to the development of the Student t-distribution. In the next section we will guide the reader on how to write queueing simulations for a G/G/1 queue.

### 4.2 Simulation of a G/G/1 Queue

We will now present an example of how to simulate a G/G/1 queue using an approach called *Discrete Event Simulation* [22]. Although the example presented here is for a G/G/1 queue, the principles can be easily extended to multi server and/or finite buffer queues. The first step is to generate a sequence of inter-arrival times and service times in accordance with the given distributions. (Note the discussion in Section 1.10.1 regarding the generation of random deviates.) In our example, starting at time 0, let us consider the following inter-arrival times: 1, 2, 1, 8, 4, 5,  $\dots$ , and the following sequence of service times: 4, 6, 4, 2, 5, 1,  $\dots$ .

In writing a computer simulation for G/G/1, we aim to fill in the following table for several 100,000s or millions arrivals (rows).

| arrival time | service duration | queue-size on arrival | service starts | service ends | delay |
|--------------|------------------|-----------------------|----------------|--------------|-------|
| 1            | 4                | 0                     | 1              | 5            | 4     |
| 3            | 6                | 1                     | 5              | 11           | 8     |
| 4            | 4                | 2                     |                |              |       |
| 12           | 2                |                       |                |              |       |
| 16           | 5                |                       |                |              |       |
| 21           | 1                |                       |                |              |       |

The following comments explain how to fill in the table.

- The arrival times and the service durations values are readily obtained from the interarrival and service time sequences.
- Assuming that the previous rows are already filled in, the “queue-size on arrival” is obtained by comparing the arrival time of the current arrivals and the values in the “service ends” column of the previous rows. In particular, the queue size on arrival is equal to the number of customers that arrive before the current customer (previous rows) that their “service ends” time values are greater than the arrival time value of the current arrival.
- The “service starts” value is the maximum of the “arrival time” value of the current arrival and the “service end” value of the previous arrival. Also notice that if the queue size on arrival of the current arrival is equal to zero, the service start value is equal to the “arrival time” value of the current arrival and if the queue size on arrival of the current arrival is more than zero the service start value is equal to the “service end” value of the previous arrival.
- The “service ends” value is simply the sum of the “service starts” and the “service duration” values of the current arrival.
- The “delay” value is the difference between the “service ends” and the “arrival time” values.

Using the results obtained in the last column, we can have the delay distribution and moments in steady-state. However, the “queue-size on arrival” values for all the customers do not, in general, provide directly the steady-state queue-size distribution and moments. To estimate accurately the steady-state queue-size distribution, we will need to have inspections performed by an independent Poisson inspector. Fortunately, due to PASTA, for M/G/1 (including M/M/1 and M/D/1) the “queue-size on arrival” values can be used directly to obtain the steady-state queue-size distribution and moments and a separate Poisson inspector is not required. Observing the queue-size just before the arrivals provides the right inspections for steady state queue-size statistics. However, if the arrival process does not follow a Poisson process, a separate independent Poisson inspector is required. In such a case, we generate a Poisson process:  $t_1, t_2, t_3, \dots$ , and for each  $t_i, i = 1, 2, 3, \dots$  we can invoke the queue-size at time  $t_i$ , denoted  $Q_i$ , in a similar way to one we obtained the “queue-size on arrival” values. The  $Q_i$  values are then used to evaluate the queue-size distribution and moments.

**Homework 4.1**

Fill in the above table by hand.     $\square$

**Homework 4.2**

Write a computer simulation for a P/P/1 queue (a single-server queue with Pareto inter-arrival and service time distributions) to derive estimates for the mean delay and mean queue-size. Perform the simulations for a wide range of parameter values. Compute confidence interval as described in Section 4.     $\square$

**Homework 4.3**

Repeat the simulations, for a wide range of parameter values, for a U/U/1 queue (a single-server queue with Uniform inter-arrival and service time distributions) and for an M/M/1 queue. For the M/M/1 queue, verify that your simulation results are consistent with respective analytical results.     $\square$

**Homework 4.4**

Discuss the accuracy of your estimations in the different cases.     $\square$

**Homework 4.5**

Use the principles presented here for a G/G/1 queue simulation to write a computer simulation for G/G/k/k queue. In particular consider the U/U/k/k and M/M/k/k queues to compute results for the blocking probability.     $\square$

There will be many homework assignments in this book that require simulations and in some cases a guide for other queueing simulations will be provided.

## 5 Deterministic Queues

We consider here the simple case where inter-arrival and service times are deterministic. In such queues, we can follow the queue-size process, possibly for some transient period, until we discover a pattern (cycle) that repeats itself. To avoid ambiguity, we assume that if an arrival and a departure occur at the same time, the departure occurs first. Such an assumption is not required for Markovian queues where the queue size process follows a continuous-time Markov chain because the probability of two events occurring at the same time is zero, but it is needed for deterministic queues. Unlike many of the Markovian queues that we study in this book, for deterministic queues steady state queue size distribution does not exist because the queue size deterministically fluctuate according to a certain pattern. Therefore, for deterministic queues we will use the notation  $P(Q = n)$ , normally designating the steady state probability of the queue-size to be equal to  $n$  in cases where such steady state probability exists, for the proportion of time that there are  $n$  customers in the queue, or equivalently,  $P(Q = n)$  is the probability of having  $n$  in the queue at a randomly (uniformly) chosen point in time. Accordingly, the mean queue size  $E[Q]$  will be defined by  $E[Q] = \sum_{n=0}^{\infty} nP(Q = n)$ . We will use the term blocking probability  $P_b$  to designate the proportion of packets that are blocked.

### 5.1 D/D/1

If we consider the case  $\lambda > \mu$ , the D/D/1 queue is unstable. In this case the queue size constantly increases and approach infinity and since there are always packets in the queue waiting for service, the server is always busy, thus the utilization is equal to one.

Let us consider now a stable D/D/1 queue, assuming  $\lambda < \mu$ . Notice that for D/D/1, given our above assumption that if an arrival and a departure occur at the same time, the departure occurs first, the case  $\lambda = \mu$  will also be stable. Assume that the first arrival occurs at time  $t = 0$ . The service time of this arrival will terminate at  $t = 1/\mu$ . Then another arrival will occur at time  $t = 1/\lambda$  which will be completely served at time  $t = 1/\lambda + 1/\mu$ , etc. This gives rise to a deterministic cyclic process where the queue-size takes two values: 0 and 1 with transitions from 0 to 1 in points of time  $n(1/\lambda)$ ,  $n = 0, 1, 2, \dots$ , and transitions from 1 to 0 in points of time  $n(1/\lambda) + 1/\mu$ ,  $n = 0, 1, 2, \dots$ . Each cycle is of length  $1/\lambda$  during which there is a customer to be served  $1/\mu$  of the time and there is no customer  $1/\lambda - 1/\mu$  of the time. Therefore, the utilization is given by  $U = (1/\mu)/(1/\lambda) = \lambda/\mu$  which is consistent with G/G/1.

As all the customers that enter the system are served before the next one arrives, the mean queue-size of D/D/1 must be equal to the mean queue-size at the server, and therefore, it is also equal to the utilization. In other words, the queue-size alternates between the values 1 and 0, spending  $1/\mu$  time at state 1, then  $1/\lambda - 1/\mu$  at state 0, then again  $1/\mu$  time at state 1, etc. If we pick a random point in time, the probability that there is one in the queue is given by  $P(Q = 1) = (1/\mu)/(1/\lambda)$ , and the probability that there are no customers in the queue is given by  $P(Q = 0) = 1 - (1/\mu)/(1/\lambda)$ . Therefore, the mean queue-size is given by  $E[Q] = 0P(Q = 0) + 1P(Q = 1) = (1/\mu)/(1/\lambda) = U$ .

Moreover, we can show that out of all possible G/G/1 queues, with  $\lambda$  being the arrival rate and  $\mu$  the service rate, no-one will have lower mean queue-size than D/D/1. This can be shown using Little's formula  $E[Q] = \lambda E[D]$ . Notice that for each of the relevant G/G/1 queues  $E[D] = 1/\mu + E[W_Q] \geq 1/\mu$ , but for D/D/1  $E[W_Q] = 0$ . Thus,  $E[D]$  for any G/G/1 queue

must be equal or greater than that of D/D/1, and consequently by Little's formula,  $E[Q]$  for any G/G/1 queue must be equal or greater than that of D/D/1.

## 5.2 D/D/ $k$

Here we consider deterministic queues with multiple servers. The interarrival times are again always equal to  $1/\lambda$ , and the service time of all messages is equal to  $1/\mu$ . Again if we consider the case  $\lambda > k\mu$ , the D/D/ $k$  queue is unstable. In this case the queue size constantly increases and approach infinity and since there are always more than  $k$  packets in the queue waiting for service, all  $k$  servers are constantly busy, thus the utilization is equal to one.

Now consider the stable case of  $\lambda < k\mu$ , so that the arrival rate is below the system capacity. Notice again that given our above assumption that if an arrival and a departure occur at the same time, the departure occurs first, the case  $\lambda = k\mu$  will also be stable. Extending the D/D/1 example to a general number of servers, the behavior of the D/D/ $k$  queue is analyzed as follows. As  $\lambda$  and  $\mu$  satisfy the stability condition  $\lambda < k\mu$ , there must exist an integer  $n$ ,  $1 \leq n \leq k$  such that  $(n-1)\mu \leq \lambda \leq n\mu$ . Thus,  $(n-1)/\lambda \leq 1/\mu$ . Meaning that if the first arrival arrives at  $t = 0$ , there will be additional  $n-1$  arrivals before the first customer leaves the system. Therefore, the queue-size increases incrementally taking the value  $j$  at time  $t = (j-1)/\lambda$ ,  $j = 1, 2, 3, \dots, n$ . When the queue reaches  $n$  for the first time, the cyclic behavior starts. Then, at time  $t = 1/\mu$  the queue-size reduces to  $n-1$  when the first customer completes its service. Next, at time  $t = n/\lambda$ , the queue-size increases to  $n$  and reduces to  $n-1$  at time  $t = 1/\lambda + 1/\mu$  when the second customer completes its service. This cyclic behavior continuous forever whereby the queue-size increases from  $n-1$  to  $n$  at time points  $t = (n+i)/\lambda$ , and reduces from  $n$  to  $n-1$  at time points  $t = i/\lambda + 1/\mu$ , for  $i = 0, 1, 2, \dots$ . The cycle length is  $1/\lambda$  during which the queue-size process is at state  $n$   $1/\mu - (n-1)/\lambda$  of the cycle time and it is at state  $n-1$   $n/\lambda - 1/\mu$  of the cycle time. Thus,  $P(Q = n) = \lambda/\mu - (n-1)$  and  $P(Q = n-1) = n - \lambda/\mu$ . The mean queue-size  $E[Q]$ , can be obtained by  $E[Q] = (n-1)P(Q = n-1) + nP(Q = n)$  which after some algebra gives

$$E[Q] = \frac{\lambda}{\mu}. \quad (237)$$

### Homework 5.1

Perform the algebraic operations that lead to (237).  $\square$ .

This result is consistent with Little's formula. As customers are served as soon as they arrive, the time each of them spends in the system is the service time  $1/\mu$  - multiplying it by  $\lambda$ , gives by Little's formula the mean queue size. Since  $E[Q]$  in D/D/ $k$  gives the number of busy servers, the utilization is given by

$$U = \frac{\lambda}{k\mu}. \quad (238)$$

Notice that Equations (237) and (238) applies also to D/D/ $\infty$  for finite  $\lambda$  and  $\mu$ . Eq. (237) gives the mean queue-size of D/D/ $\infty$  (by Little's formula, or by following the arguments that led to Eq. (237)) and for D/D/ $\infty$ , we have that  $U = 0$  by (238). Also notice that in D/D/ $\infty$  there are infinite number of servers and the number of busy servers is finite, so the average utilization per server must be equal to zero.



### 5.3 D/D/k/k

In D/D/k/k there is no waiting room beyond those available at the servers. Recall that to avoid ambiguity, we assume that if an arrival and a departure occur at the same time, the departure occurs first. Accordingly, if  $\lambda \leq k\mu$ , then we have the same queue behavior as in D/D/k as no losses will occur. The interesting case is the one where  $\lambda > k\mu$  and this is the case we focus on. Having  $\lambda > k\mu$ , or  $1/\mu > k/\lambda$ , implies that there exists a positive integer  $n$  such that

$$\frac{(k+n-1)}{\lambda} < \frac{1}{\mu} < \frac{(k+n)}{\lambda}.$$

#### Homework 5.2

Prove the last statement.  $\square$

Again, consider an empty system with the first arrival occurring at time  $t = 0$ . There will be additional  $k - 1$  arrivals before all the servers are busy. Notice that because  $1/\mu > k/\lambda$ , no service completion will occur before the system is completely full. Then  $n$  additional arrivals will be blocked before the first customer completes its service at time  $t = 1/\mu$  at which time the queue-size decreases from  $k$  to  $k - 1$ . Next, at time  $t = (k+n)/\lambda$ , the queue-size increases to  $k$  and reduces to  $k - 1$  at time  $t = 1/\lambda + 1/\mu$  when the second customer completes its service. This behavior of the queue-size alternating between the states  $k$  and  $k - 1$  continues until all the first  $k$  customers complete their service which happens at time  $t = (k-1)/\lambda + 1/\mu$  when the  $k$ th customer completes its service, reducing the queue-size from  $k$  to  $k - 1$ , followed by an arrival at time  $t = (2k+n-1)/\lambda$  increasing the queue-size from  $k - 1$  to  $k$ . Notice that the point in time  $t = (2k+n-1)/\lambda$  is an end-point of a cycle that started at  $t = (k-1)/\lambda$ . This cycle comprises two parts: the first is a period of time where the queue-size stays constant at  $k$  and all the arrivals are blocked, and the second is a period of time during which no losses occur and the queue-size alternates between  $k$  and  $k - 1$ . Then a new cycle of duration  $(k+n)/\lambda$  starts and this new cycle ends at  $t = (3k+2n-1)/\lambda$ . In general, for each  $j = 1, 2, 3, \dots$ , a cycle of duration  $(k+n)/\lambda$  starts at  $t = ((j)k + (j-1)n - 1)/\lambda$  and ends at  $t = ((j+1)k + jn - 1)/\lambda$ .

In every cycle, there are  $k+n$  arrivals out of which  $n$  are blocked. The blocking probability is therefore  $P_B = n/(k+n)$ . Since,  $(k+n-1)/\lambda < 1/\mu < (k+n)/\lambda$ , we have that  $k+n = \lceil (\lambda/\mu) \rceil$ , where  $\lceil x \rceil$  is the smallest integer greater or equal to  $x$ . Therefore the blocking probability is given by

$$P_b = \frac{\lceil \frac{\lambda}{\mu} \rceil - k}{\lceil \frac{\lambda}{\mu} \rceil}. \quad (239)$$

Let  $A = \lambda/\mu$ , the mean-queue size is obtained using Little's formula to be given by

$$E[Q] = \frac{\lambda}{\mu}(1 - P_b) = \frac{kA}{\lceil A \rceil}. \quad (240)$$

As in D/D/k, since every customer that enters a D/D/k/k system does not wait in a queue, but immediately enters service, the utilization is given by

$$U = \frac{E[Q]}{k} = \frac{A}{\lceil A \rceil}. \quad (241)$$

### Homework 5.3

1. Follow the evolution of the D/D/ $k/k$  queue-size process to derive the queue-size probabilities  $P(Q = k - 1)$  and  $P(Q = k)$  and use these probabilities to confirm (240).
2. Consider a D/D/3/3 queue with  $1/\mu = 5.9$  and  $1/\lambda = 1.1$ . Start with the first arrival at  $t = 0$  and produce a two column table showing the time of every arrival and departure until  $t = 20$ , and the corresponding queue-size values immediately following each one of these events.
3. Write a general simulation program for a D/D/ $k/k$  queue and use it to validate (240) and your results for  $P(Q = k - 1)$  and  $P(Q = k)$ . Use it also to confirm the results you obtained for the D/D/3/3 queue.
4. Describe the evolution of the queue-size process and drive formulae for the mean queue size, mean delay, utilization, and blocking probability of a D/D/1/ $k$  queue. Confirm your results by simulation  $\square$ .

### 5.4 Summary of Results

The following table summarizes the results on D/D/1, D/D/ $k$  and D/D/ $k/k$ . Note that we do not consider the case  $\lambda = k\mu$  which gives the same  $E[Q]$  and  $U$  results as  $\lambda < k\mu$  assuming that if a departure and an arrival occur at the same time, the departure occurs before the arrival.

| Model      | Condition        | $E[Q]$               | $U$                 |
|------------|------------------|----------------------|---------------------|
| D/D/1      | $\lambda < \mu$  | $\lambda/\mu$        | $\lambda/\mu$       |
| D/D/1      | $\lambda > \mu$  | $\infty$             | 1                   |
| D/D/ $k$   | $\lambda < k\mu$ | $A = \lambda/\mu$    | $A/k$               |
| D/D/ $k$   | $\lambda > k\mu$ | $\infty$             | 1                   |
| D/D/ $k/k$ | $\lambda < k\mu$ | $A$                  | $A/k$               |
| D/D/ $k/k$ | $\lambda > k\mu$ | $kA/\lceil A \rceil$ | $A/\lceil A \rceil$ |

## 6 M/M/1

Having considered the straightforward cases of deterministic queues, we will now discuss queues where the interarrival and service times are non-deterministic. We will begin with cases where the inter-arrival and service times are independent and exponentially distributed (memoryless). Here we consider the M/M/1 queue where the arrival process follows a Poisson process with parameter  $\lambda$  and service times are assumed to be IID and exponentially distributed with parameter  $\mu$ , and are independent of the arrival process. As M/M/1 is a special case of G/G/1, all the results that are applicable to G/G/1 are also applicable to M/M/1. For example,  $U = \lambda/\mu$ ,  $p_0 = 1 - \lambda/\mu$  and Little's formula. It is the simplest Markovian queue; it has only a single server and an infinite buffer. It is equivalent to a continuous time Markov chain on the states:  $0, 1, 2, 3, \dots$ . Assuming that the M/M/1 queue-size process starts at state 0, it will stay in state 0 for a period of time that is exponentially distributed with parameter  $\lambda$  then it moves to state 1. The time the process stays in state  $n$ , for  $n \geq 1$ , is also exponentially distributed, but this time, it is a competition between two exponential random variable, one of which is the time until the next arrival - exponentially distributed with parameter  $\lambda$ , and the other is the time until the next departure - exponentially distributed with parameter  $\mu$ . As discussed in Section 1.10.2, the minimum of the two is therefore also exponential with parameter  $\lambda + \mu$ , and this minimum is the time the process stays in state  $n$ , for  $n \geq 1$ . We also know from the discussion in Section 1.10.2 that after spending an exponential amount of time with parameter  $\lambda + \mu$ , the process will move to state  $n + 1$  with probability  $\lambda/(\lambda + \mu)$  and to state  $n - 1$  with probability  $\mu/(\lambda + \mu)$ .

### 6.1 Steady State Queue Size Probabilities

As the M/M/1 queue-size process increases by only one, decreases by only one and stays an exponential amount of time at each state, it is equivalent to a birth and death process. Therefore, by Eqs. (208) and (209), the infinitesimal generator for the M/M/1 queue-size process is given by

$$\begin{aligned} Q_{i,i+1} &= \lambda \text{ for } i=0, 1, 2, 3, \dots \\ Q_{i,i-1} &= \mu \text{ for } i= 1, 2, 3, 4, \dots \\ Q_{0,0} &= -\lambda \\ Q_{i,i} &= -\lambda - \mu \text{ for } i=1, 2, 3, \dots \end{aligned}$$

Substituting this infinitesimal generator in Eq. (210) and performing some simple algebraic operations, we obtain the following steady state equations for the M/M/1 queue.

$$\begin{aligned} \pi_0 \lambda &= \pi_1 \mu \\ \pi_1 \lambda &= \pi_2 \mu \\ &\dots \end{aligned}$$

and in general:

$$\pi_i \lambda = \pi_{i+1} \mu, \text{ for } i = 0, 1, 2, \dots \quad (242)$$

To explain (242) intuitively, Let  $L$  be a very long time. During  $L$ , the total time that the process stays in state  $i$  is equal to  $\pi_i L$ . Since the arrival process is Poisson, the mean number of transitions from state  $i$  to  $i + 1$  is equal to  $\lambda \pi_i L$ . For  $i = 0$ ,  $\lambda \pi_0 L$  is also the mean number of events that occur during  $L$  (because there are no departures at state 0). However, for  $i \geq 1$ ,

the mean number of events that occur during  $L$  in state  $i$  is  $(\lambda + \mu)\pi_i L$  because as soon as the process enters state  $i$  it stays there on average an amount of time equal to  $1/(\mu + \lambda)$  and then it moves out of state  $i$  to either state  $i + 1$ , or to state  $i - 1$ . Since during time  $\pi_i L$  there are on average  $(\lambda + \mu)\pi_i L$  interval times of size  $1/(\mu + \lambda)$ , then  $(\lambda + \mu)\pi_i L$  is also the mean number of events (arrivals and departures) that occur in state  $i$  during  $L$ . Therefore, the mean number of departures that occur in state  $i$  (transitions from  $i$  to  $i - 1$ ) during  $L$  is equal to the product of the mean number of events  $(\lambda + \mu)\pi_i L$  and the probability that an event is a departure which is  $\mu/(\lambda + \mu)$ , namely

$$(\lambda + \mu)\pi_i L \frac{\mu}{\lambda + \mu} = \mu\pi_i L.$$

In a similar way, the mean number of arrivals that occur in state  $i$  during  $L$  is

$$(\lambda + \mu)\pi_i L \frac{\lambda}{\lambda + \mu} = \lambda\pi_i L.$$

Since  $L$  is very long, we must have that the number of transitions from state  $i$  to  $i + 1$  is equal to the number of transitions from state  $i + 1$  to  $i$ . Therefore,  $\pi_i \lambda L = \pi_{i+1} \mu L$ , and dividing both sides by  $L$ , we obtain (242).

Of course, the sum of the steady state probabilities must be equal to one, so we have the additional equation

$$\sum_{j=1}^{\infty} \pi_j = 1. \quad (243)$$

Let  $\rho = \lambda/\mu$ , we obtain,

$$\begin{aligned} \pi_1 &= \rho\pi_0 \\ \pi_2 &= \rho\pi_1 = \rho^2\pi_0 \\ \pi_3 &= \rho\pi_2 = \rho^3\pi_0 \end{aligned}$$

and in general:

$$\pi_i = \rho^i \pi_0 \text{ for } i = 0, 1, 2, \dots \quad (244)$$

As M/M/1 is a special case of G/G/1, we can use Eq. (220) to obtain  $\pi_0 = 1 - \rho$ , so

$$\pi_i = \rho^i (1 - \rho) \text{ for } i = 0, 1, 2, \dots \quad (245)$$

Let  $Q$  be a random number representing the queue-size in steady state. Its mean is obtained by  $E[Q] = \sum_{i=0}^{\infty} i\pi_i$ . This leads to:

$$E[Q] = \frac{\rho}{1 - \rho}. \quad (246)$$

## Homework 6.1

Perform the algebraic operations that lead to (246).  $\square$

## 6.2 Delay Statistics

By (245), and by the PASTA principle, an arriving customer will have to pass a geometric number of IID phases, each of which is exponentially distributed with parameter  $\mu$ , until it

leaves the system. We have already shown that a geometrically distributed sum of an IID exponentially distributed random variables is exponentially distributed (see Eq. (136) in Section 1.13.2). Therefore the total delay of any arriving customer in an M/M/1 system must be exponentially distributed. This can also be intuitively explained. Because both geometric and exponential distributed random variables are memoryless, a geometrically distributed sum of IID exponential random variables is also memoryless. And since the exponential is the only memoryless distribution, the total delay of any arriving customer must be exponentially distributed.

Therefore, to derive the density of the delay, all that is left to do is to obtain its mean which can be derived by (246) invoking Little's formula. Another way to obtain the mean delay is by noticing from (245) that the number of phases is geometrically distributed with mean  $1/(1 - \rho)$ . Observe that this mean must equal  $E[Q] + 1$  which is the mean queue-size observed by an arriving customer plus one more phase which is the service time of the arriving customer. Thus, the mean number of phases is

$$E[Q] + 1 = \frac{\rho}{1 - \rho} + 1 = \frac{1 - \rho + \rho}{1 - \rho} = \frac{1}{1 - \rho}.$$

## Homework 6.2

Prove that the number of phases is geometrically distributed with mean  $1/(1 - \rho)$ .

## Guide

Let  $P_h$  be the number of phases. We know that in steady state an arriving customer will find  $Q$  customers in the system, where

$$P(Q = i) = \pi_i = \rho^i(1 - \rho).$$

Since  $P_h = Q + 1$ , we have

$$P(P_h = n) = P(Q + 1 = n) = P(Q = n - 1) = \rho^{n-1}(1 - \rho).$$

□

The mean delay equals the mean number of phases times the mean service time  $1/\mu$ . Thus,

$$E[D] = \frac{1}{(1 - \rho)\mu} = \frac{1}{\mu - \lambda}. \quad (247)$$

## Homework 6.3

Verify that (246) and (247) are consistent with Little's formula. □

Substituting  $1/E[D] = \mu - \lambda$  as the parameter of exponential density, the density of the delay distribution is obtained to be given by

$$\delta_D(x) = \begin{cases} (\mu - \lambda)e^{(\lambda - \mu)x} & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases} \quad (248)$$

### 6.3 Using Z-Transform

The Z-transform defined in Section 1.13, also known as Probability Generating Function, is a powerful tool to derive statistics of queueing behavior.

As an example, we will now demonstrate how the Z-transform is used to derive the mean queue-size of M/M/1.

Let us multiply the  $n$ th equation of (242) by  $z^n$ . Summing up both sides will give

$$\frac{\Psi(z) - \pi_0}{z} = \rho\Psi(z) \quad (249)$$

where  $\Psi(z) = \sum_{i=0}^{\infty} \pi_i z^i$ . Letting  $z$  approach 1 (from below) gives

$$\pi_0 = 1 - \rho \quad (250)$$

which is consistent with what we know already. Substituting it back in (249) gives after simple algebraic manipulation:

$$\Psi(z) = \frac{1 - \rho}{1 - \rho z}. \quad (251)$$

Taking derivative and substituting  $z = 1$ , after some algebra we obtain

$$E[Q] = \Psi^{(1)}(1) = \frac{\rho}{1 - \rho} \quad (252)$$

which is again consistent with what we know about M/M/1 queue.

### Homework 6.4

1. Derive equations (249) – (252).
2. Derive the variance of the M/M/1 queue-size using Z-transform.  $\square$

### 6.4 Multiplexing

An important and interesting observation we can make by considering the M/M/1 queueing performance results (245)–(248) is that while the queue-size statistics are dependent only on  $\rho$  (the ratio of the arrival rate and service rate), the delay statistics (mean and distribution) are a function of what we call the *spare capacity* (or *mean net input*) which is the difference between the service rate and the arrival rate. To be more specific, it is a linear function of the reciprocal of that difference.

Assume that our traffic model obeys the M/M/1 assumptions. Then if the arrival rate increases from  $\lambda$  to  $N\lambda$  and we increase the service rate from  $\mu$  to  $N\mu$  (maintaining the same  $\rho$ ), the mean queue-size and its distribution will remain the same. However, in this scenario the mean delay does not remain the same. It reduces by  $N$  times to  $1/[N(\mu - \lambda)]$ .

If our QoS of interest is the mean delay, or the probability that the delay exceeds a certain value, and if for a given arrival rate  $\lambda$  there is a service rate  $\mu$  such that our QoS measure (mean delay) is just met, then if the arrival rate increases from  $\lambda$  to  $N\lambda$ , and we aim to find

the service rate  $\mu^*$  such that the QoS measure (mean delay) is just met, we will need to make sure that the spare capacity is maintained, that is

$$\mu - \lambda = \mu^* - N\lambda \quad (253)$$

or

$$\mu^* = \mu + (N - 1)\lambda \quad (254)$$

so by the latter and the stability condition of  $\mu > \lambda$ , we must have that  $\mu^* < N\mu$ . We can therefore define a measure for multiplexing gain to be given by

$$M_{mg} = \frac{N\mu - \mu^*}{N\mu} \quad (255)$$

so by (254), we obtain

$$M_{mg} = \frac{N - 1}{N}(1 - \rho). \quad (256)$$

Recalling the stability condition  $\rho < 1$ , Eq. (256) has four important messages.

1. The multiplexing gain is positive for all  $N > 1$ .
2. The multiplexing gain increases with  $N$ .
3. The multiplexing gain is bounded above by  $1 - \rho$ .
4. In the limiting condition as  $N \rightarrow \infty$ , the multiplexing gain approaches its bound  $1 - \rho$ .

The  $1 - \rho$  bound is consistent with the observation that if  $\rho$  is very close to 1, then the multiplexing gain diminishes because in this case, any significant reduction in the service rate below  $N\mu$  will lead to instability. This intuitively makes sense. Recalling the relationship  $1 - \rho = \pi_0$ , namely,  $\rho$  represents the proportion of time that the server is busy, or its utilization, it is clear that if the utilization is already very high, there is little room for improvement, so the potential multiplexing gain is low. On the other hand, if we have a case where the QoS requirements are strict (requiring very low mean queueing delay) such that the utilization  $\rho$  must be low the potential for multiplexing gain is high.

Let us now apply our general discussion on multiplexing to obtain insight into performance comparison between two commonly used multiple access techniques used in telecommunications. One such technique is called *Time Division Multiple Access* (TDMA) whereby each user is assigned one or more channels (in a form of time slots) to access the network. Another approach, which we call *full multiplexing* (FMUX), is to let all users to separately send the data that they wish to transmit to a switch which then forwards the data to the destination. That is, all the data is stored in one buffer (in the switch) which is served by the entire available link capacity.

To compare between the two approaches, let us consider  $N$  users each transmitting packets at an average rate of  $R_u$  [bits/second]. The average packet size denoted  $S_u$  [bits] is assumed equal for the different users. Let  $\hat{\lambda}$  [packets/second] be the packet rate generated by each of the users. Thus,  $\hat{\lambda} = R_u/S_u$ . Under TDMA, each of the users obtains a service rate of  $B_u$  [bits/sec].  $\hat{\mu}$  [packets/second]. Packet sizes are assumed to be exponentially distributed with mean  $S_u$  [bits], so the service rate in packets/second denoted  $\hat{\mu}$  is given by  $\hat{\mu} = B_u/S_u$ . The

packet service time is therefore exponentially distributed with parameter  $\hat{\mu}$ . Letting  $\hat{\rho} = \hat{\lambda}/\hat{\mu}$ , the mean queue size under TDMA, is given by

$$E[Q_{TDMA}] = \frac{\hat{\rho}}{1 - \hat{\rho}}, \quad (257)$$

and the mean delay is

$$E[D_{TDMA}] = \frac{1}{\hat{\mu} - \hat{\lambda}}. \quad (258)$$

In the FMUX case the total arrival rate is  $N\hat{\lambda}$  and the service rate is  $N\hat{\mu}$ , so in this case, the ratio between the arrival and service rate remains the same, so the mean queue size that only depends on this ratio remains the same

$$E[Q_{FMUX}] = \frac{\hat{\rho}}{1 - \hat{\rho}} = E[Q_{TDMA}]. \quad (259)$$

However, we can observe an  $N$ -fold reduction in the mean delay:

$$E[D_{FMUX}] = \frac{1}{N\hat{\mu} - N\hat{\lambda}} = \frac{E[D_{TDMA}]}{N}. \quad (260)$$

Consider a telecommunication provider that wishes to meet packet delay requirement of its  $N$  customers, assuming that the delay that the customers experienced under TDMA was satisfactory, and assuming that the M/M/1 assumptions hold, such provider does not need a total capacity of  $N\hat{\mu}$  for the FMUX alternative. It is sufficient to allocate  $\hat{\mu} + (N - 1)\hat{\lambda}$ .

### Homework 6.5

Consider a telecommunication provider that aims to serve a network of 100 users each transmits data at a rate of 1 Mb/s. The mean packet size is 1 kbit. Assume that packets lengths are exponentially distributed and that the process of packets generated by each user follows a Poisson process. Further assume that the mean packet delay requirement is 50 millisecond. How much total capacity (bitrate) is required to serve the 100 users under TDMA and under FMUX.  $\square$

## 6.5 The Departure Process

According to the so-called Burke theorem [13], in steady state, the departure process of a stable M/M/1, where  $\rho < 1$ , is a Poisson process with parameter  $\lambda$  and is independent of the number in the queue. To see why this is so all we need is to realize that the queue size process in an M/M/1 queue with  $\rho < 1$  is an irreducible, aperiodic and stable Markov chain. As such it must be reversible. Therefore, the points in time of arrivals in the forward process correspond to points in time that the Markov chain is increased by one. These points represent the arrival process and therefore follow a Poisson process. By reversibility, in steady state, the arrival process of the reversed process must also follow Poisson process with parameter  $\lambda$  and this process is the departure process of the forward process.

Now that we know that in steady state the departure process of M/M/1 is Poisson with parameter  $\lambda$ , we also know that, in steady state, the inter-departure times are also exponentially



distributed with parameter  $\lambda$ . We will now show this fact without using the fact that the departure process is Poisson directly. Instead, we will use it indirectly to induce PASTA for the reversed arrival process to obtain that, following a departure, in steady state, the queue is empty with probability  $1 - \rho$  and non-empty with probability  $\rho$ . If the queue is non-empty, the time until the next departure is exponentially distributed with parameter  $\mu$  – this is the service-time of the next customer. If the queue is empty, we have to wait until the next customer arrival which is exponentially distributed with parameter  $\lambda$  and then we will have to wait until the next departure which will take additional time which is exponentially distributed. All together, if the queue is empty, the time until the next departure is a sum of two exponential random variables, one with parameter  $\lambda$  and the other with parameter  $\mu$ . Let  $U_1$  and  $U_2$  be two independent exponential random variables with parameters  $\lambda$  and  $\mu$ , respectively. Define  $U = U_1 + U_2$ , and let us first derive the density  $f_U(u)$  of  $U$ . In other words,  $U$  has hypo-exponential distribution. Having the density  $f_U(u)$  the density  $f_D(t)$  of a random variable  $D$  representing the inter-departure time will be given by

$$f_D(t) = \rho\mu e^{-\mu t} + (1 - \rho)f_U(t). \quad (261)$$

Knowing that  $f_U(u)$  is a convolution of two exponentials, we obtain

$$\begin{aligned} f_U(t) &= \int_{u=0}^t \lambda e^{-\lambda u} \mu e^{-\mu(t-u)} du \\ &= \frac{\lambda\mu}{\mu - \lambda} (e^{-\lambda t} - e^{-\mu t}). \end{aligned}$$

The by the latter and (261), we obtain

$$f_D(t) = \rho\mu e^{-\mu t} + (1 - \rho) \frac{\lambda\mu}{\mu - \lambda} (e^{-\lambda t} - e^{-\mu t}) \quad (262)$$

which after some algebra gives

$$f_D(t) = \lambda e^{-\lambda t}. \quad (263)$$

This result is consistent with Burke theorem.

## Homework 6.6

Complete all the algebraic details in the derivation of equations (261) – (263).  $\square$

Another way to show consistency with Burke theorem is the following. Let  $d_\epsilon$  be the mean number of departures that leaves the M/M/1 queue during a small interval of time of size  $\epsilon$ , and let  $d_\epsilon(i)$  be the mean number of departures that leaves the M/M/1 queue during a small interval of time of size  $\epsilon$  when there are  $i$  packets in our M/M/1 queue. Then,  $d_\epsilon(i) = 0$  if  $i = 0$ , and  $d_\epsilon(i) = \epsilon\mu$  if  $i > 0$ . Therefore, in steady state,  $d_\epsilon = [0(1 - \rho)\mu\epsilon + 1(\rho)\mu\epsilon] = \epsilon\lambda$ , which is a property consistent with the assertion of Poisson output process with parameter  $\lambda$  in steady state.

## Homework 6.7

So far we have discussed the behaviour of the M/M/1 departure process in steady state. You are now asked demonstrate that the M/M/1 departure process may not be Poisson with parameter

$\lambda$  if we do not assume steady state condition. Consider an M/M/1 system with arrival rate  $\lambda$  and service rate  $\mu$ , assume that  $\rho = \lambda/\mu < 1$  and that there are no customers in the system at time 0. Derive the distribution of the number of customers that leave the system during the time interval  $(0, t)$ . Argue that this distribution is, in most cases, not Poisson with parameter  $\lambda t$  and find a special case when it is.

### Guide

Let  $D(t)$  be a random variable representing the number of customers that leave the system during the time interval  $(0, t)$ . Let  $X_p(\lambda t)$  be a Poisson random variable with parameter  $\lambda t$  and consider two cases: (a) the system is empty at time  $t$ , and (b) the system is not empty at time  $t$ . In case (a),  $D(t) = X_p(\lambda t)$  (why?) and in case (b)  $D(t) = X_p(\lambda t) - 1$  (why?) and use the notation used in Section 2.5  $P_{00}(t)$  to denote the probability that in time  $t$  the system is empty, so the probability that the system is not empty at time  $t$  is  $1 - P_{00}(t)$ . Derive  $P_{00}(t)$  using Eqs. (205) and (206). Then notice that

$$D(t) = P_{00}(t)X_p(\lambda t) + [1 - P_{00}(t)][X_p(\lambda t) + 1]. \quad \square$$

Consider the limit

$$D_k(t) = \lim_{\Delta t \rightarrow 0} P[Q(t) = k \mid \text{a departure occurs within}(t - \Delta t, t)].$$

Considering the fact that the reversed process is Poisson and independence between departures before time  $t$  and  $Q(t)$ , we obtain that

$$D_k(t) = P[Q(t) = k]. \quad (264)$$

Then, by taking the limit of both sides of (264), we show that the queue size seen by a leaving customer is statistically identical to the queue size seen by an independent observer.  $\square$

### Homework 6.8

Write a simulation of the M/M/1 queue by measuring queue size values in two ways: (1) just before arrivals and (2) just before departure. Verify that the results obtain for the mean queue size in steady state are consistent. Use confidence intervals. Verify that the results are also consistent with analytical results. Repeat your simulations and computation for a wide range of parameters values (different  $\rho$  values). Plot all the results in a graph including the confidence intervals (bars).  $\square$

## 6.6 Mean Busy Period and First Passage Time

The *busy period* of a single-server queueing system is defined as the time between the point in time the server starts being busy and the point in time the server stops being busy. In other words, it is the time elapsed from the moment a customer arrives at an empty system until the first time the system is empty again. Recalling the first passage time concept defined in Section 2.5, and that the M/M/1 system is in fact a continuous time Markov chain, the busy period is

also the first passage time from state 1 to state 0. The end of a busy period is the beginning of the so called *idle period* - a period during which the system is empty. We know the mean of the idle period in an M/M/1 queue. It is equal to  $1/\lambda$  because it is the mean time until a new customer arrives which is exponentially distributed with parameter  $\lambda$ . A more interesting question is what is the mean busy period. Let  $T_B$  and  $T_I$  be the busy and the idle periods, respectively. Noticing that  $E[T_B]/(E[T_B] + E[T_I])$  is the proportion of time that the server is busy, thus it is equal to  $\rho$ . Considering also that  $E[T_I] = 1/\lambda$ , we obtain

$$\frac{E[T_B]}{E[T_B] + \frac{1}{\lambda}} = \rho. \quad (265)$$

Therefore,

$$E[T_B] = \frac{1}{\mu - \lambda}. \quad (266)$$

Interestingly, for the M/M/1 queue the mean busy period is equal to the mean delay of a single customer! This may seem counter intuitive. However, we can realize that many busy periods are made of a single customer service time which in this case is equal to its delay which is shorter than the average delay. Furthermore, this interesting fact can be easily proved by considering an M/M/1 queue with a service policy of Last In First Out (LIFO) policy. So far we have considered only queues that their service policy is First In First Out (FIFO). Let us consider an M/M/1 with LIFO with preemptive priority. In such a queue the arrival and service rates  $\lambda$  and  $\mu$ , respectively, are the same as those of the FIFO M/M/1, but in the LIFO queue, the customer just arrived has priority over all other customers that arrived before it and in fact interrupts the customers currently in service.

The two queues we consider, the FIFO and the LIFO, are both birth and death processes with the same parameters so their respective queue size processes are statistically the same. Then by Little's formula their respective mean delays are also the same. Also the delay of a customer in a M/M/1 LIFO queue we consider is equal to the busy period in M/M/1 FIFO queue (why?) so the mean delay must be equal to the busy period in M/M/1 with FIFO service policy.

### Homework 6.9

Derive an expression for the mean first passage time for M/M/1 from state  $n$  to state 0 and from state 0 to state  $n$ , for  $n \geq 3$ .  $\square$

### Homework 6.10

For a wide range of parameter values, simulate an M/M/1 system with FIFO service policy and an M/M/1 system with LIFO service policy with preemptive priority and compare their respective results for the mean delay, the variance of the delay, the mean queue size and the mean busy period.  $\square$

## 6.7 A Markov Chain Simulation of M/M/1

A simulation of an M/M/1 queue can be made as a special case of G/G/1 as described before, or it can be simplified by taking advantage of the M/M/1 Markov chain structure if we are not

interested in performance measures that are associated with times (such as delay distribution). If our aim is to evaluate queue size statistics or blocking probability, we can avoid tracking the time. All we need to do is to collect the relevant information about the process at PASTA time-points without even knowing what is the running time at these points. Generally speaking, using the random walk simulation approach, also called the *Random Walk simulation* approach, we simulate the evolution of the states of the process based on the transition probability matrix and collect information on the values of interest at selective PASTA points without being concerned about the time. We will now explain how these ideas may be applied to few relevant examples.

If we wish to evaluate the mean queue size of an M/M/1 queue, we can write the following simulation.

Variables and input parameters:  $Q$  = queue size;  $\hat{E}(Q)$  = estimation for the mean queue size;  $N$  = number of  $Q$ -measurements taken so far which is also equal to the number of arrivals so far;  $MAXN$  = maximal number of  $Q$ -measurements taken;  $\mu$  = service rate;  $\lambda$  = arrival rate.

Define function:  $I(Q) = 1$  if  $Q > 0$ ;  $I(Q) = 0$  if  $Q = 0$ .

Define function:  $R(01)$  = a uniform  $U(0, 1)$  random deviate. A new value for  $R(01)$  is generated every time it is called.

Initialization:  $Q = 0$ ;  $E[Q] = 0$ ;  $N = 0$ .

1. If  $R(01) \leq \lambda/(\lambda + I(Q)\mu)$ , then  $N = N + 1$ ,  $\hat{E}(Q) = [(N - 1)\hat{E}(Q) + Q]/N$ , and  $Q = Q + 1$ ; else,  $Q = Q - I(Q)$ .

2. If  $N < MAXN$  go to 1; else, print  $\hat{E}(Q)$ .

This signifies the simplicity of the simulation. It has only two If statements: one to check if the next event is an arrival or a departure according to Eq. (50), and the second is merely a stopping criterion.

## Homework 6.11

Simulate an M/M/1 queue using a Markov chain simulation to evaluate the mean queue size for the cases of Section 4.2, and compare the results with the results obtain analytically and with those obtained using the G/G/1 simulation principles. In you comparison consider accuracy (closeness to the analytical results) the length of the confidence intervals and running times.

□

## 7 M/M/ $\infty$

The next queueing system we shall consider is the M/M/ $\infty$  queueing system where the number of servers is infinite. Because the number of servers is infinite, the buffer capacity is unlimited and arrivals are never blocked. We assume that the arrival process is Poisson with parameter  $\lambda$  and each server renders service which is exponentially distributed with parameters  $\mu$ . As in the case of M/M/1 we assume that the service times are independent and are independent of the arrival process.

### 7.1 Steady State Equations

Like M/M/1, the M/M/ $\infty$  system can also be viewed as a continuous time Markov chain with the state being the queue-size (the number of customers in the system). As in M/M/1, the arrival rate is independent of changes in the queue-size. However, unlike M/M/1, in M/M/ $\infty$ , the service rate does change with the queue-size. When there are  $n$  customers in the system, and at the same time,  $n$  servers are busy, the service rate is  $n\mu$ , and the time until the next event is exponentially distributed with parameter  $\lambda + n\mu$ , because it is a competition between  $n + 1$  exponential random variables:  $n$  with parameter  $\mu$  and one with parameter  $\lambda$ .

Equating the probability flux, we obtain the following steady state equations:

$$\pi_0\lambda = \pi_1\mu$$

$$\pi_1\lambda = \pi_2 2\mu$$

...

and in general:

$$\pi_n\lambda = \pi_{n+1}(n+1)\mu, \text{ for } n = 0, 1, 2, \dots \quad (267)$$

Of course, the sum of the steady state probabilities must be equal to one, so we again have the additional equation

$$\sum_{j=1}^{\infty} \pi_j = 1. \quad (268)$$

### 7.2 Solving the Steady State Equations

Let  $A = \lambda/\mu$ . Notice that we use the notation  $A$  here for the  $\lambda/\mu$  while we use the notation  $\rho$  for this ratio in the M/M/1 case. The two serve different roles. In M/M/1, we must have  $\rho < 1$  for stability, while in M/M/ $\infty$ , we usually have  $A > 1$ . In M/M/1,  $\rho$  is the proportion of time the server is busy (the utilization), while in M/M/ $\infty$ , the value of  $A$  is equal to the average number of busy servers, as shown below.

Using the  $A$  notation we obtain

$$\pi_1 = A\pi_0$$

$$\pi_2 = A\pi_1/2 = A^2\pi_0/2$$

$$\pi_3 = A\pi_2/3 = A^3\pi_0/(3!)$$

and in general:

$$\pi_n = \frac{A^n \pi_0}{n!} \text{ for } n = 0, 1, 2, \dots \quad (269)$$

To obtain  $\pi_0$ , we sum up both sides of Eq. (269), and because the sum of the  $\pi_n$ s equals one, we obtain

$$1 = \sum_{n=0}^{\infty} A^n \pi_0 / n! \quad (270)$$

or

$$\pi_0 = e^{-A} \quad (271)$$

Substituting the latter in Eq. (269), we obtain

$$\pi_n = \frac{e^{-A} A^n}{n!} \text{ for } n = 0, 1, 2, \dots \quad (272)$$

By Eq. (272) we observe that the distribution of the number of busy channels (simultaneous calls or customers) in an M/M/ $\infty$  system is Poisson with parameter  $A$ . The quantity  $A$  is therefore the traffic carried by an infinite number of circuits. Interestingly, this result can be obtained by Little's formula. According to Little's formula, the mean number of customers in the system is equal to the arrival rate ( $\lambda$ ) times the mean time a customer spends in the system which in the case of M/M/ $\infty$  is equal to  $1/\mu$ . Thus

$$E[Q] = \lambda(1/\mu) = A. \quad (273)$$

In practice, the number of circuits is limited, and the *offered traffic* is higher than the *carried traffic* because some of the calls are blocked due to call congestion when all circuits are busy. A queueing model which describe this more realistic case is the M/M/ $k/k$  queueing model studied by A. K. Erlang [19] (see next section). To his honor, the quantity of traffic is measured in units of erlangs. Accordingly, if we consider an M/M/ $k/k$  system where the ratio between the arrival rate and the service rate is  $A$ ), we say that, in this case,  $A$  erlangs are offered to the system, as this represents the average number of simultaneous calls in progress assuming  $k$  is unlimited.

### 7.3 Insensitivity

The above results for  $\pi_i$ ,  $i = 0, 1, 2, \dots$  and for the mean number of busy servers are insensitive to the service time (holding time) distribution, all we need is the mean of the distribution and the results are the same regardless of what the distribution is. In other words, the above results apply to an M/G/ $\infty$  system. This is a powerful result because it allows us to use the same results for applications where the service time is not exponential. This insensitivity property is valid also for the M/M/ $k/k$  system.

### 7.4 Applications

An interesting application of the M/M/ $\infty$  system is to the following multi-access problem (see Problem 3.8 in [9]). Consider a stream of packets that their arrival times follow a Poisson process with parameter  $\lambda$ . If the inter-arrival times of any pair of packets (not necessarily a

consecutive pair) is less than the transmission time of the packet that arrived earlier out of the two, these two packets are said to collide. Assume that packets have independent exponentially distributed transmission times with mean  $\mu$ . What is the probability of no collision?

Notice that a packet can collide with any one or more of the packets that arrived before it. In other words, it is possible that it may not collide with its immediate predecessor, but it may collide with a packet that arrived earlier. However, if it does not collide with its immediate successor, it will not collide with any of the packets that arrive after the immediate successor.

Therefore, the probability that an arriving packet will not collide on arrival can be obtained to be the probability of an M/M/ $\infty$  system to be empty, that is,  $e^{-A}$ . While the probability that its immediate successor will not arrive during its transmission time is  $\mu/(\lambda + \mu)$ . The product of the two, namely  $e^{-A}\mu/(\lambda + \mu)$ , is the probability of no collision.

Another application of the M/M/ $\infty$  system (or M/G/ $\infty$  system) is to the following problem. Consider a city with population 3,000,000, and assume that the birth rate  $\lambda$  is constant. Average life time of people in this city is 78 years. How to compute the birth rate? Using the M/M/ $\infty$  model (or actually the M/G/ $\infty$  as human lifetime is not exponentially distributed) with  $E[Q] = 3,000,000$  and  $\mu^{-1} = 78$ , realizing that  $E[Q] = A = \lambda/\mu$ , we obtain,  $\lambda = \mu E[Q] = 3,000,000/78 = 38461$  new births per year or 105 new births per day.

## 8 Erlang B Formula

The next queueing system we shall consider is the M/M/ $k/k$  queueing system where the number of servers is  $k$ . We assume that the arrival process is Poisson with parameter  $\lambda$  and each server renders service which is exponentially distributed with parameters  $\mu$ . As in the other M/M/... cases, we assume that the service times are independent and are independent of the arrival process. We will now discuss Erlang's derivation of the loss probability of an M/M/ $k/k$  system that leads to the well known Erlang's Loss Formula, also known as Erlang B Formula.

### 8.1 Solving the Steady State Equations

The steady state equations for M/M/ $k/k$  are the same as the first  $k$  steady state equations for M/M/ $\infty$ . Accordingly, we obtain for M/M/ $k/k$ :

$$\pi_n = \frac{A^n \pi_0}{n!} \text{ for } n = 0, 1, 2, \dots, k. \quad (274)$$

To obtain  $\pi_0$ , we again sum up both sides of the latter. This leads to

$$\pi_0 = \frac{1}{\sum_{n=0}^k \frac{A^n}{n!}}. \quad (275)$$

Substituting Eq. (275) in Eq. (274), we obtain

$$\pi_n = \frac{\frac{A^n}{n!}}{\sum_{n=0}^k \frac{A^n}{n!}} \text{ for } n = 0, 1, 2, \dots, k. \quad (276)$$

The most important quantity out of the values obtained by Eq. (276) is  $\pi_k$ . It is the probability that all  $k$  circuits are busy. It gives the proportion of time that no new calls can enter the system, namely, they are blocked. The quantity  $\pi_k$  for an M/M/ $k/k$  system loaded by offered traffic  $A$  is usually denoted by  $E_k(A)$  and is given by:

$$E_k(A) = \frac{\frac{A^k}{k!}}{\sum_{n=0}^k \frac{A^n}{n!}}. \quad (277)$$

Eq. (277) is known as Erlang's loss Formula, or Erlang B Formula, published first by A. K. Erlang in 1917 [19].

Due to the special properties of the Poisson process, in addition of being the proportion of time that the calls are blocked,  $E_k(A)$  also gives the proportion of calls blocked due to congestion; namely, it is the *blocking probability*. Eq. (277) has many applications for telecommunications network design. Given its importance, it is necessary to be able to compute Eq. (277) quickly and exactly for large values of  $k$ . We should be able to answer a dimensioning question such as "how many circuits are required so that the blocking probability is no more than 1% given offered traffic of  $A = 1000$ ".

### 8.2 Recursion and Jagerman Formula

Observing Eq. (277), we notice the factorial terms which may hinder such computation for a large  $k$ . We shall now present a short analysis which leads to a recursive relation between



$E_n(A)$  and  $E_{n-1}(A)$  which can be used in a simple recursive algorithm which can compute the blocking probability for large values of  $k$ . By simple algebraic operations using Eq. (277), we obtain

$$\frac{E_n(A)}{E_{n-1}(A)} = \frac{A}{n}(1 - E_n(A)). \quad (278)$$

Isolating  $E_n(A)$ , this leads to

$$E_n(A) = \frac{AE_{n-1}(A)}{n + AE_{n-1}(A)} \text{ for } n = 1, 2, \dots, k. \quad (279)$$

When  $n = 0$ , there are no servers (circuits) available, and therefore all customers (calls) are blocked, namely,

$$E_0(A) = 1. \quad (280)$$

The above two equations give rise to a simple recursive algorithm by which the blocking probability can be calculated for a large  $k$ . An even more computationally stable way to compute  $E_n(A)$  for large values of  $A$  and  $n$  is to use the inverse [29]

$$I_n(A) = \frac{1}{E_n(A)} \quad (281)$$

and the recursion

$$I_n(A) = 1 + \frac{n}{A}I_{n-1}(A) \text{ for } n = 1, 2, \dots, k. \quad (282)$$

with the initial condition  $I_0(A) = 1$ .

A useful formula for  $I_n(A)$  due to Jagerman [30] is:

$$I_n(A) = A \int_0^\infty e^{-Ay}(1+y)^n dy. \quad (283)$$

### 8.3 The Special Case: M/M/1/1

#### Homework 8.1

Derive a formula for the blocking probability of M/M/1/1 in four ways: (1) by Erlang B Formula (277), (2) by the recursion (279), (3) by the recursion (282), and (4) by Jagerman Formula (283).  $\square$

The reader may observe a fifth direct way to obtain a formula for the blocking probability of M/M/1/1 using Little's formula. The M/M/1/1 system can have at most one customer in it. Therefore, its mean queue size is given by  $E[Q] = 0\pi_0 + 1\pi_1 = \pi_1$  which is also its blocking probability. Noticing also that the arrival rate into the system (made only of successful arrivals) is equal to  $\lambda(1 - E[Q])$ , the mean time a customer stays in the system is  $1/\mu$ , and revoking Little's formula, we obtain

$$\frac{\lambda(1 - E[Q])}{\mu} = E[Q]. \quad (284)$$

Isolating  $E[Q]$ , the blocking probability is given by

$$\pi_1 = E[Q] = \frac{A}{1 + A}. \quad (285)$$

## 8.4 Dimensioning and Utilization

Taking advantage of the monotonicity of Erlang formula, we can also solve the dimensioning problem. We simply keep incrementing the number of circuits and calculate in each case the blocking probability. When the desired blocking probability (e.g., 1%) is reached, we have our answer.

### Homework 8.2

Prove that if  $A > A'$  then  $E_n(A) > E_n(A')$ .  $\square$

The mean number of busy circuits in an M/M/k/k system fed by  $A$  erlangs can be obtained in two ways. First, it can be computed by the weighted sum  $\sum_{i=0}^k i\pi_i$ , or using Little's formula. By Little's formula The mean queue size is equal to the arrival rate, which in this case is  $\lambda$  multiplied by  $1 - \pi_k$ , times the mean time in the system which in this case is  $1/\mu$ . Therefore, we obtain

$$E[Q] = (1 - \pi_k)\lambda/\mu = (1 - \pi_k)A. \quad (286)$$

Notice that it is lower than the corresponding mean for M/M/ $\infty$  which is equal to  $A$ .

Accordingly, the utilization of an M/M/k/k system is given by

$$U = \frac{(1 - \pi_k)A}{k}. \quad (287)$$

### Homework 8.3

Prove that  $\sum_{i=0}^k i\pi_i = (1 - \pi_k)A$ .  $\square$

In the following Table, we present the minimal values of  $k$  obtained for various values of  $A$  such that the blocking probability is no more than 1%, and the utilization obtained in each case. It is clearly observed that the utilization increased with the traffic.

| $A$   | $k$  | $E_k(A)$ | Utilization |
|-------|------|----------|-------------|
| 20    | 30   | 0.0085   | 66.10%      |
| 100   | 117  | 0.0098   | 84.63%      |
| 500   | 527  | 0.0095   | 93.97%      |
| 1000  | 1029 | 0.0099   | 96.22%      |
| 5000  | 5010 | 0.0100   | 98.81%      |
| 10000 | 9970 | 0.0099   | 99.30%      |

### Homework 8.4

Reproduce the above Table.  $\square$

We also notice that for the case of  $A = 10,000$  erlangs, to maintain no more than 1% blocking,  $k$  value less than  $A$  is required. Notice however that the carried traffic is not  $A$  but  $A(1 - E_k(A))$ . This means that for  $A \geq 10,000$ , dimensioning simply by  $k = A$  will mean no more than 1% blocking and no less than 99% utilization - not bad for such a simple rule of thumb!

## 8.5 Insensitivity and Many Classes of Customers

We have mentioned in the previous chapter that the distribution and the mean of the number of busy servers of  $M/G/\infty$  and  $M/G/k/k$ . For  $M/G/k/k$ , also the blocking probability is insensitive to the service time distribution. However must make it very clear that this insensitivity property does not go far. We can only use the nice results of the steady state equations such as (274) when we have a continuous time Markov chain, or in the above described very special situations where the insensitivity property applies as in  $M/G/\infty$  and  $M/G/k/k$ . To demonstrate it, let us compare an  $M/M/1/1$  system with a  $D/D/1/1$  system. Suppose that each of these two systems is fed by  $A$  erlangs, and that  $A < 1$ .

Arrivals into the  $D/D/1/1$  system will never experience losses because the inter-arrivals are longer than the service times, so the service of a customer is always completed before the arrival of the next customer. Accordingly, by Little's formula:  $E[Q] = A$ , and since  $E[Q] = 0 \times \pi_0 + 1 \times \pi_1$ , we have that  $\pi_1 = A$  and  $\pi_0 = 1 - A$ . In this case, the blocking probability  $P_b$  is equal to zero and not to  $\pi_1$ . As there are no losses, the utilization will be given by  $U = \pi_1 = A$ .

By contrast, for the  $M/M/1/1$  system,  $P_b = E_1(A) = E[Q] = \pi_1 = A/(1 + A)$ , so  $\pi_0 = 1 - \pi_1 = 1/(1 + A)$ . To obtain the utilization we can either realize that it is the proportion of time our single server is busy, namely it is equal to  $\pi_1 = A/(1 + A)$ , or we can use the above formula for  $U$  in  $M/M/k/k$  system and obtain

$$U = (1 - \pi_k)A = [1 - A/(1 + A)]A = A/(1 + A). \quad (288)$$

This comparison is summarized in the following table:

|         | M/M/1/1     | D/D/1/1 |
|---------|-------------|---------|
| $\pi_0$ | $1/(1 + A)$ | $1 - A$ |
| $\pi_1$ | $A/(1 + A)$ | $A$     |
| $U$     | $A/(1 + A)$ | $A$     |
| $P_b$   | $A/(1 + A)$ | $0$     |
| $E[Q]$  | $A/(1 + A)$ | $A$     |

Clearly, the steady state equations (274) will not apply to a  $D/D/1/1$  system.

We have already mentioned that for  $M/G/k/k$  the distribution of the number of busy servers and therefore also the blocking probability is insensitive to the service time distribution. All we need is to know that the arrival process is Poisson, and the ratio of the arrival rate to the service rate and we can obtain the blocking probability using the Erlang B formula. Let us now consider the following problem.

Consider two classes of customers (packets). Class  $i$  customers arrives at rate of  $\lambda_i$  each of which requires exponentially distributed service with parameter  $\mu_i$ , for  $i = 1, 2$ . There are  $k$  servers without waiting room (without additional buffer). The aim is to derive the blocking probability.

The combined arrival process of all the customers is a Poisson process with parameter  $\lambda = \lambda_1 + \lambda_2$ . Because the probability of an arbitrary customer to belong to the first class is

$$p = \frac{\lambda_1}{\lambda_1 + \lambda_2},$$

the service time of an arbitrary customer has hyperexponential distribution because with probability  $p$  it is exponentially distributed with parameter  $\mu_1$ , and with probability  $1 - p$ , it is exponentially distributed with parameter  $\mu_2$ .

The mean service time (holding time) is therefore given by

$$E(S) = \frac{p}{\mu_1} + \frac{1-p}{\mu_2}$$

so  $A = \lambda E(S)$ , and Erlang B applies.

### Homework 8.5 [9]

Extend the results obtained for two classes of customers to the case of  $n$  classes of customers.  
□

### Homework 8.6 [9]

This assignment applies to both M/M/ $\infty$  and M/M/ $k/k$  models. Consider an M/M/ $\infty$  queueing system with the following twist. The servers are numbered 1, 2, ... and an arriving customer always chooses the server numbered lowest among all the free servers it finds. Find the proportion of time that each of the servers is busy. Why can't you solve it just by considering the probability of having  $k$  or more servers in an M/M/ $\infty$  queueing system?

**Hint:** Use the Erlang B and Little's formulae. □

### Homework 8.7

This assignment also applies to both M/M/ $\infty$  and M/M/ $k/k$  models. Bursts of data of exponential lengths with mean  $1/\mu$  that arrive according to a Poisson process are transmitted through a bufferless optical switch. All arriving bursts compete for  $k$  wavelength channels at a particular output trunk of the switch. If a burst arrives and all  $k$  wavelength channels are busy, the burst is dumped at the wavelength bit-rate. While it is being dumped, if one of the wavelength channels becomes free, the remaining portion of the burst is successfully transmitted through the wavelength channel.

1. Show that the mean loss rate of data  $E[Loss]$  is given by

$$E[Loss] = 1P(X = k + 1) + 2P(X = k + 2) + 3P(X = k + 3) + \dots$$

where  $X$  is a Poisson random variable with parameter  $A = \lambda/\mu$ .

2. Prove that

$$E[Loss] = \frac{A\gamma(k, A)}{\Gamma(k)} - \frac{k\gamma(k + 1, A)}{\Gamma(k + 1)}$$

where  $\Gamma(k)$  is the Gamma function and  $\gamma(k, A)$  is the lower incomplete Gamma function.

## Background information and guide

The Gamma function is defined by

$$\Gamma(a) = \int_0^{\infty} t^{a-1} e^{-t} dt. \quad (289)$$

The lower incomplete Gamma function is defined by

$$\gamma(a, x) = \int_0^x t^{a-1} e^{-t} dt. \quad (290)$$

The upper incomplete Gamma function is defined by

$$\Gamma(a, x) = \int_x^{\infty} t^{a-1} e^{-t} dt. \quad (291)$$

Accordingly,

$$\gamma(a, x) + \Gamma(a, x) = \Gamma(a).$$

For an integer  $k$ , we have

$$\Gamma(k) = (k-1)!. \quad (292)$$

$$\gamma(k, x) = (k-1)! e^{-x} \sum_{m=0}^{k-1} \frac{x^m}{m!}. \quad (293)$$

Therefore,

$$e^{-A} \sum_{m=0}^k \frac{A^m}{m!} = \frac{\Gamma(k+1, A)}{\Gamma(k+1)} \quad (294)$$

so

$$1 - e^{-A} \sum_{m=0}^k \frac{A^m}{m!} = 1 - \frac{\Gamma(k+1, A)}{\Gamma(k+1)} = \frac{\Gamma(k+1) - \Gamma(k+1, A)}{\Gamma(k+1)} = \frac{\gamma(k+1, A)}{\Gamma(k+1)}. \quad (295)$$

□

## 8.6 A Markov Chain Simulation of M/M/k/k

To evaluate the blocking probability of an M/M/k/k queue, we can use the following simulation.

Variables and input parameters:

$Q$  = number of customers in the system (queue size);

$B_p$  = estimation for the blocking probability;

$N_a$  = number of customer arrivals counted so far;

$N_b$  = number of blocked customers counted so far;

$MAXN_a$  = maximal number of customer arrivals (it is used for the stopping condition);

$\mu$  = service rate;

$\lambda$  = arrival rate.

Define function:  $R(01)$  = a uniform  $U(0, 1)$  random deviate. A new value for  $R(01)$  is generated every time it is called.

Initialization:  $Q = 0$ ;  $N_a = 0$ ,  $N_b = 0$ .

1. If  $R(01) \leq \lambda/(\lambda + Q\mu)$ , then  $N_a = N_a + 1$ ; if  $Q = k$  then  $N_b = N_b + 1$ , else  $Q = Q + 1$ ;

else,  $Q = Q - 1$ .

2. If  $N_a < MAXN_a$  go to 1; else, print  $B_p = N_b/N_a$ .

Again, it is a very simple program of two If statements: one to check if the next event is an arrival or a departure, and the other a stopping criterion.

### Homework 8.8

Simulate an M/M/k/k queue based on the Markov chain simulation principles to evaluate the blocking probability for a wide range of parameter values. Compare the results you obtained with equivalent results obtain analytically using the Erlang B Formula and with equivalent M/M/k/k queue blocking probability results obtained using the general simulation principles of Section 4.2. In your comparison consider accuracy (closeness to the analytical results), the length of the confidence intervals and running times.  $\square$

### Homework 8.9

Simulate equivalent U/U/k/k, M/U/k/k and M/M/k/k models. Run these simulations for a wide range of parameter values and compare them numerically. Compare them also with equivalent results obtain analytically using the Erlang B Formula.  $\square$

### Homework 8.10

Use the M/M/k/k model to compare the utilization of an optical switch with full wavelength conversion and without wavelength conversion.

### Background information and guide

Consider a switch with 100 input trunks and 100 output trunks. Each trunk comprises 10 optical fibers each of which comprises 40 wavelengths. Consider a particular input trunk and assume that the traffic directed to it follows a Poisson process with parameter  $\lambda$  and that any packet is of exponential length with parameter  $\mu$ . In the case of full wavelength conversion, every packet from any wavelength can be converted to any other wavelength, so the Poisson traffic with parameter  $\lambda$  can all be directed to the trunk and can use any of the  $10 \times 40 = 400$  links. In the case of no wavelength conversion only traffic on a given wavelength (assume symmetry so model it by a Poisson process with parameter  $\lambda/40$ ) is directed to 10 output links associated with this same wavelength. Compare efficiency that can be achieved if the blocking probability is set limited to 0.001.  $\square$

## 9 M/M/ $k$

The M/M/ $k$  queue is a generalization of the M/M/1 queue to the case of  $k$  servers. As in M/M/1, for an M/M/ $k$  queue, the buffer is infinite and the arrival process is Poisson with rate  $\lambda$ . The service time of each of the  $k$  servers is exponentially distributed with parameter  $\mu$ . As in the case of M/M/1 we assume that the service times are independent and are independent of the arrival process.

### 9.1 Steady State Equations and Their Solution

Letting  $A = \lambda/\mu$ , and assuming the stability condition  $\lambda < k\mu$ , or  $A < k$ , the M/M/ $k$  queue gives rise to the following steady state equations:

$$\begin{aligned} \pi_1 &= A\pi_0 \\ \pi_2 &= A\pi_1/2 = A^2\pi_0/2 \\ \pi_3 &= A\pi_2/3 = A^3\pi_0/(3!) \\ &\dots \\ \pi_k &= A\pi_{k-1}/k = A^k\pi_0/(k!) \\ \pi_{k+1} &= A\pi_k/k = A^{k+1}\pi_0/(k!k) \\ \pi_{k+2} &= A\pi_{k+1}/k = A^{k+2}\pi_0/(k!k^2) \\ &\dots \\ \pi_{k+j} &= A\pi_{k+j-1}/k = A^{k+j}\pi_0/(k!k^j) \text{ for } j = 1, 2, 3, \dots \end{aligned}$$

and in general:

$$\pi_n = \frac{A^n \pi_0}{n!} \text{ for } n = 0, 1, 2, \dots, k-1 \quad (296)$$

and

$$\pi_n = \frac{A^n \pi_0}{k! k^{n-k}} \text{ for } n = k, k+1, k+2, \dots \quad (297)$$

To obtain  $\pi_0$ , we sum up both sides of Eqs. (296) and (297), and because the sum of the  $\pi_n$ s equals one, we obtain an equation for  $\pi_0$ , which its solution is

$$\pi_0 = \left( \sum_{n=0}^{k-1} \frac{A^n}{n!} + \frac{A^k}{k!} \frac{k}{k-A} \right)^{-1}. \quad (298)$$

Substituting the latter in Eqs. (296) and (297), we obtain the steady state probabilities  $\pi_n$ ,  $n = 0, 1, 2, \dots$ .

### 9.2 Erlang C Formula

Of special interest is the so called Erlang C formula. It represents the proportion of time that all  $k$  servers are busy and is given by:

$$C_k(A) = \sum_{n=k}^{\infty} \pi_n = \frac{A^k}{k!} \frac{k}{k-A} \pi_0 = \frac{\frac{A^k}{k!} \frac{k}{k-A}}{\sum_{n=0}^{k-1} \frac{A^n}{n!} + \frac{A^k}{k!} \frac{k}{k-A}}. \quad (299)$$

**Homework 9.1**

Derive Eq. (299).  $\square$

By Eqs. (277) and (299) we obtain the following relationship:

$$C_k(A) = \frac{kE_k(A)}{k - A[1 - E_k(A)]}. \quad (300)$$

**Homework 9.2**

1. Derive Eq. (300);
2. Show that  $C_k(A) \geq E_k(A)$ .  $\square$

Let  $N_Q$  be a random number representing the number of customers waiting in the queue (not including those in service). The mean queue-size is obtained by  $E[N_Q] = \sum_{n=k}^{\infty} (n - k)\pi_n$ . This leads to:

$$E[N_Q] = C_k(A) \frac{A}{k - A}. \quad (301)$$

**Homework 9.3**

1. Derive Eq. (301).
2. Confirm consistence between (301) and (246).  $\square$

In the following table, we add the corresponding  $C_k(A)$  values to the table of the previous section. We can observe the significant difference between  $C_k(A)$  and  $E_k(A)$  as the ratio  $A/k$  increases. Clearly, when  $A/k > 1$ , the M/M/k queue is unstable.

| $A$   | $k$  | $E_k(A)$ | $C_k(A)$ |
|-------|------|----------|----------|
| 20    | 30   | 0.0085   | 0.025    |
| 100   | 117  | 0.0098   | 0.064    |
| 500   | 527  | 0.0095   | 0.158    |
| 1000  | 1029 | 0.0099   | 0.262    |
| 5000  | 5010 | 0.0100   | 0.835    |
| 10000 | 9970 | 0.0099   | unstable |

**Homework 9.4**

Reproduce the results of the above table.  $\square$



### 9.3 Dimensioning and Utilization

One could solve the dimensioning problem of finding the smallest  $k$  such that, for a given  $A$ , the blocking probability will be lower than a given value. Using Eq. (300), and realizing that the value of  $C_k(A)$  decreases as  $k$  increases, the dimensioning problem can be solved in an analogous way to the M/M/k/k dimensioning problem.

The utilization of an M/M/k queue is given by

$$U = \pi_1 \frac{1}{k} + \pi_2 \frac{2}{k} + \pi_3 \frac{3}{k} + \dots + \pi_{k-1} \frac{k-1}{k} + \pi_k + \pi_{k+1} + \pi_{k+2} + \dots$$

or

$$U = 1 - \pi_0 - \pi_1 \frac{k-1}{k} - \pi_2 \frac{k-2}{k} - \pi_3 \frac{k-3}{k} - \dots - \pi_{k-1} \frac{1}{k}. \quad (302)$$

A simpler way to obtain the utilization is again to use Little's formula for the system made of servers. If we consider the system of servers (without considering the waiting room outside the servers), we notice that since there are no losses, the arrival rate into this system is  $\lambda$  and the waiting time of each customer in this system is merely its service time with mean  $1/\mu$ . Therefore, by Little's formula the mean number of busy server is given by  $A = \lambda/\mu$  and therefore the utilization of an M/M/k queue is obtained by

$$U = \frac{\lambda}{k\mu} = \frac{A}{k}. \quad (303)$$

## 10 Engset Loss Formula

Engset Loss Formula has applied to telephony situations where the number of customers is small relative to the number of available circuits. Such situations include: an exchange in a small rural community, PABX, or even a lucrative satellite service to a small number of customers. Let the call holding times be IID exponentially distributed with mean  $1/\mu$  and the time until an idle source attempts to make a call is also exponential with mean  $1/\lambda$ . We also assume that there is not dependence between the holding times and the idle periods of the sources. Let the number of customers (sources of traffic) be  $M$ , the number of available circuits  $k$  and the blocking probability  $P_b$ .

The reader will recall that in  $M/M/1$ , the arrival rate as well as the service rate are independent of the state of the system, and in  $M/M/\infty$ , the arrival rate is also independent of the number of customers in the system, but the service rate is state dependent. In the present case, when the number of customers is limited, we have a case where both the arrival rate and the service rate are state dependent.

As in  $M/M/k/k$ , the service rate is  $n\mu$  when there are  $n$  busy circuits (namely  $n$  customers are making phone calls). However, unlike  $M/M/k/k$ , in the present case, busy customers do not make new phone calls thus they do not contribute to the arrival rate. Therefore, if  $n$  circuits are busy, the arrival rate is  $(M - n)\lambda$ . As a result, considering both arrival and service processes, at any point in time, given that there are  $n$  customers in the system, and at the same time,  $n$  servers/circuits are busy, the time until the next event is exponentially distributed with parameter  $(M - n)\lambda + n\mu$ , because it is a competition between  $M$  exponential random variables:  $n$  with parameter  $\mu$  and  $M - n$  with parameter  $\lambda$ .

An important question we must always answer in any Markov chain analysis is how many states do we have. If  $M > k$ , then the number of states is  $k + 1$ , as in  $M/M/k/k$ . However, if  $M < k$ , the number of states is  $M + 1$  because no more than  $M$  calls can be in progress at the same time. Therefore, the number of states is  $\min\{M, k\} + 1$ .

### 10.1 Steady State Equations and Their Solution

Equating the probability flux, we obtain the following steady state equations:

$$\pi_0 M \lambda = \pi_1 \mu$$

$$\pi_1 (M - 1) \lambda = \pi_2 2 \mu$$

$$\pi_2 (M - 2) \lambda = \pi_3 3 \mu$$

...

and in general:

$$\pi_n (M - n) \lambda = \pi_{n+1} (n + 1) \mu, \text{ for } n = 0, 1, 2, \dots, \min\{M, k\} - 1. \quad (304)$$

Therefore, after standard algebraic manipulations of (304), that are left as an exercise to the reader, we can write  $\pi_n$ , for  $n = 0, 1, 2, \dots, \min\{M, k\}$ , in terms of  $\pi_0$  as follows:

$$\pi_n = \binom{M}{n} \left( \frac{\lambda}{\mu} \right)^n \pi_0, \text{ for } n = 0, 1, 2, \dots, \min\{M, k\}, \quad (305)$$

or

using the notation  $\rho = \lambda/\mu$ , we obtain

$$\pi_n = \binom{M}{n} \rho^n \pi_0, \text{ for } n = 0, 1, 2, \dots, \min\{M, k\}. \quad (306)$$

Of course, the sum of the steady state probabilities must be equal to one, so we again have the additional normalizing equation

$$\sum_{j=1}^{\min\{M,k\}} \pi_j = 1. \quad (307)$$

By (306) together with the normalizing Eq. (307), we obtain

$$\pi_n = \frac{\binom{M}{n} \rho^n}{\sum_{j=0}^{\min\{M,k\}} \binom{M}{j} \rho^j}, \text{ for } n = 0, 1, 2, \dots, \min\{M, k\}. \quad (308)$$

## 10.2 Blocking Probability

Now, what is the blocking probability  $P_b$ ? When  $k \geq M$ , clearly  $P_b = 0$ , as there is never a shortage of circuits.

To derive the blocking probability for the case when  $k < M$ , we first realize that unlike in the case of Erlang Formula,  $\pi_k$  does not give the blocking probability. Still,  $\pi_k$  is the probability of having  $k$  busy circuits, or the proportion of time that all circuits are busy which is the so-called *time-congestion*, but it is not the probability that a call is blocked – the so-called *call-congestion*. Unlike the case of Erlang Formula, here, call-congestion is not equal to time congestion. This is because in the Engset model, the arrival rate is dependent on the state of the system. When the system is full the arrival rate is much lower than when the system is empty.

In particular when  $i$  circuits are busy, the arrival rate is  $\lambda(M-i)$ , therefore to find the proportion of calls blocked, or the blocking probability denoted  $P_b$ , we compute the ratio between calls arrive when there are  $k$  circuits busy and the total calls arrive. This gives

$$P_b = \frac{\lambda(M-k)\pi_k}{\lambda \sum_{i=0}^k (M-i)\pi_i}. \quad (309)$$

After some derivations, which are left as an exercise to the reader, we obtain the Engset loss formula that gives the blocking probability for the case  $M > k$  as follows.

$$P_b = \frac{\binom{M-1}{k} \rho^k}{\sum_{i=0}^k \binom{M-1}{i} \rho^i}. \quad (310)$$

Notice that  $\rho$ , defined above by  $\rho = \lambda/\mu$ , is the intensity of a **free** customer. An interesting interpretation of (310) is that the call congestion, or the blocking probability, when there are  $M$  sources is equal to the time congestion when there are  $M - 1$  sources. This can be intuitively explained as follows. Consider an arbitrary tagged source (or customer). For this particular customer, the proportion of time it cannot access is equal to the proportion of time the  $k$  circuits are all busy by the other  $M - 1$  customers. During the rest of the time our tagged source can successfully access a circuit.

### 10.3 Obtaining the Blocking Probability by a Recursion

Letting  $B_i$  be the blocking probability given that the number of available circuits is  $i$ , the Engset loss formula can be solved numerically by the following recursion:

$$B_i = \frac{\rho(M-i)B_{i-1}}{i + \rho(M-i)B_{i-1}} \quad i = 1, 2, 3, \dots, k \quad (311)$$

with the initial condition

$$B_0 = 1. \quad (312)$$

### 10.4 Insensitivity

In his original work [18], Engset assumed that the idle time as well as the holding time are exponentially distributed. These assumptions have been relaxed over the years and now it is known that Engset formula applies also to arbitrary idle and holding time distributions [28].

### 10.5 Load Classifications and Definitions

An important feature of Engset setting is that a customer already engaged in a conversation does not originate calls. This leads to an interesting peculiarity that if we fix the number of customers (assuming  $M > k$ ) and reduce  $k$ , the offered traffic increases because reduction in  $k$  leads to increase in  $P_b$  and reduction in the average number of busy customers which in turn leads to increase in idle customers the offer more calls.

Let us now discuss the concept of the so-called *intended* offered load [6] under the Engset setting. We know that  $1/\lambda$  is the mean time until a free customer will make a call (will attempt to seize a circuit). Also,  $1/\mu$  is the mean holding time of a call. If a customer is never blocked, it is behaving like an on/off source, alternating between on and off states, being on for an exponentially distributed period of time with mean  $1/\mu$ , and being off for an exponentially distributed period of time with mean  $1/\lambda$ . For each cycle of average length  $1/\lambda + 1/\mu$ , a source will be busy, on average, for a period of  $1/\mu$ . Therefore, in steady state, the proportion of time a source is busy is  $\lambda/(\lambda + \mu)$ , and since we have  $M$  sources, the *intended* offered load is given by

$$T = M \frac{\lambda}{\lambda + \mu} = \frac{\rho M}{(1 + \rho)}. \quad (313)$$

This *intended* offered load is equal to the offered traffic load and the carried traffic load if  $M \leq k$ , namely, when  $P_b = 0$ . However, when  $M > k$  (thus  $P_b > 0$ ), the offered traffic load and the carried traffic load are not equal. Let  $T_c$  and  $T_o$  be the *carried* and the *offered* traffic load respectively. The carried traffic is the mean number of busy circuits and it is given by

$$T_c = \sum_{i=0}^k i\pi_i. \quad (314)$$

The offered traffic is obtained by averaging the arrival rates weighted by their corresponding probabilities as follows.

$$T_o = \sum_{i=0}^k \rho(M - i)\pi_i. \quad (315)$$

To compute the values for  $T_c$  and  $T_o$  in terms of the blocking probability  $P_b$ , we first realize that

$$T_c = T_o(1 - P_b), \quad (316)$$

and also,

$$T_o = \sum_{i=0}^k \rho(M - i)\pi_i = \rho M - \rho \sum_{i=0}^k i\pi_i = \rho(M - T_c) \quad (317)$$

and by (314) – (317) we obtain

$$T_c = \frac{\rho M(1 - P_b)}{[1 + \rho(1 - P_b)]} \quad (318)$$

$$T_o = \frac{\rho M}{[1 + \rho(1 - P_b)]}. \quad (319)$$

Notice that when  $P_b = 0$ , we have  $T = T_o = T_c$ . Notice also that the above three measures may be divided by  $k$  to obtain the relevant traffic load per server.

## 10.6 The Many Sources Limit

Let  $M$  approach infinity and  $\lambda$  approach zero in a way that maintains the intended offered load constant. In this case, since  $\lambda + \mu \rightarrow \mu$ , the limit of the intended load will take the form

$$\lim T = M \frac{\lambda}{\mu} = \rho M. \quad (320)$$

Furthermore, under this limiting condition, the terms  $\rho(M - i)$ ,  $i = 1, 2, 3, \dots, k$ , in (311) can be substituted by  $\rho M$  which is the limit of the intended traffic load. It is interesting to observe that if we substitute  $A = \rho M$  for the  $\rho(M - i)$  terms in (311), equations (279) and (311) are equivalent. This means that if the number of sources increases and the arrival rate of each source decreases in a way that the intended load stays fixed, the blocking probability obtained by Engset loss formula approaches that of Erlang B formula.

## 10.7 Obtaining the Blocking Probability by Successive Iterations

In many cases, it is convenient to obtain the blocking probability  $P_b$  in terms of the offered load  $T_o$ . By Eq. (319) we obtain,

$$\rho = \frac{T_o}{M - T_o(1 - P_b)}. \quad (321)$$

The latter can be used together with Eq. (310) or (311) to obtain  $P_b$  by an iterative process. One begin by setting an initial estimate value to  $P_b$  (e.g.  $P_b = 0.1$ ). Then this initial estimate is substituted into Eq. (321) to obtain an estimate for  $\rho$  then the value you obtain for  $\rho$  is substituted in Eq. (310), or use the recursion (311), to obtain another value for  $P_b$  which is then substituted in Eq. (321) to obtain another estimate for  $\rho$ . This iterative process continues until the difference between two successive estimations of  $P_b$  is arbitrarily small.

### Homework 10.1

Consider the case  $M = 20$ ,  $k = 10$ ,  $\lambda = 2$ ,  $\mu = 1$ . Compute  $P_b$  using the recursion Eq. (311). The Compute  $P_b$  using the iterative processes starting with various initial estimations. Compare the results and the running time of the program.  $\square$

## 11 A Queue with State Dependent Arrivals and Service

So far we have considered Markovian queues where the arrival and service rates are independent of the number of customers in the system. It is common, however, to have situations where this is not the case. In many systems capacity is added (service rate increases) and/or traffic is throttled back as queue size increases.

In this section we study a model of a single-server queue in which the arrival process is state dependent Poisson process. This is a Poisson process that its rate  $\lambda(i)$  fluctuates based on the queue size  $i$ . The service rate  $\mu(i)$  also fluctuates based on  $i$ . That is, when there are  $i$  customers in the system, the service is exponentially distributed with parameter  $\mu(i)$ . If during service, before the service is complete, the number of customers changes from  $i$  to  $j$  ( $j$  could be either  $i + 1$  or  $i - 1$ ) then the remaining service time changes to exponentially distributed with parameter  $\mu(j)$ . We assume that the number of customer in the queue is limited by  $B$ . In telecommunications jargon, we say that  $B$  is the buffer size in packets assuming that all packets are of equal size.

This model gives rise to a birth and death process described in Section 2.5. The state dependent arrival and service rates  $\lambda$  and  $\mu$  are equivalent to the birth and death rates  $a_i$  and  $b_i$ , respectively.

### 11.1 Steady State Queue Size Probabilities

Following the birth and death model of Section 2.5 the infinitesimal generator for our Markovian state dependent queue-size process is given by

$$\begin{aligned} Q_{i,i+1} &= \lambda(i) \text{ for } i = 0, 1, 2, 3, \dots, B \\ Q_{i,i-1} &= \mu \text{ for } i = 1, 2, 3, 4, \dots, B \\ Q_{0,0} &= -\lambda(0) \\ Q_{i,i} &= -\lambda(i) - \mu(i) \text{ for } i = 1, 2, 3, \dots, B - 1 \\ Q_{B,B} &= -\mu(B). \end{aligned}$$

Then the steady state equations  $0 = \mathbf{\Pi Q}$ , can be written as:

$$0 = -\lambda_0\pi_0 + \pi_1\mu_1 \quad (322)$$

and

$$0 = \pi_{i-1}\lambda_{i-1} - \pi_i(\lambda_i + \mu_i) + \pi_{i+1}\mu_{i+1} \text{ for } i = 1, 2, 3, \dots, B - 1. \quad (323)$$

There is an additional last dependent equation

$$0 = \pi_{B-1}\lambda_{B-1} - \pi_B(\mu_B) \quad (324)$$

which is redundant. The normalizing equation

$$\sum_{i=0}^B \pi_i = 1 \quad (325)$$

must also be satisfied.

## 11.2 Solving the Steady State Equations

Equations 322, 323 and 325 are  $B + 1$  equations for  $B + 1$  unknowns. One way to solve this set of equations is by a standard method such as Cramer's Rule or the Inverse Matrix Method.

An alternative way, that normally works well for solving a set of equations such as 322, 323 and 325 is the so called *successive substitution* method (it is also known as Gauss-Seidel, successive approximation or iterations) [17]. It can be described as follows. First, isolate  $\pi_0$  in Eq. 322. Next, isolate  $\pi_i$  in each of the  $B - 1$  equations 323 for  $i = 1, 2, \dots, B - 1$ , respectively. Then, isolate  $\pi_B$  in Eq. 324. This leads to the following vector equation for  $\mathbf{\Pi}$

$$\mathbf{\Pi} = \mathbf{f}(\mathbf{\Pi}). \quad (326)$$

Then perform the successive substitution operations by setting initial values to the vector  $\mathbf{\Pi}$ ; substitute them in the right hand side of 326 obtain different values at the left hand side which are then substituted in the right hand side, etc. For example, the initial setting can be  $\mathbf{\Pi} = \mathbf{1}$  without any regards to the normalization equation. When the values obtain for  $\mathbf{\Pi}$  are sufficiently close, say, within a distance no more than  $10^{-6}$ , stop. Then normalize the vector  $\mathbf{\Pi}$  obtained in the last iteration. This is the desired solution.

As in the case of M/M/k/k and M/M/1/k, the blocking probability is given by  $\pi_B$ .

### Homework 11.1

Consider a single-server Markovian queue with state dependent arrivals and service. You are free to choose the  $\lambda(i)$  and  $\mu(i)$  rates, but make sure they are different for different  $i$  values. Set the buffer size at  $B = 200$ . Solve the steady state equations using the successive relaxation method and using a standard method. Compare the results and the computation time. Then obtain the blocking probability by simulation and compare with the equivalent results obtained by solving the state equations. Repeat the results for a wide range of parameters by using various  $\lambda(i)$  vectors obtained by multiplying the original  $\lambda(i)$  vector by various constants.  $\square$



## 12 Three Queueing Models with Finite Buffers

So far we have considered queueing models with either infinite buffer, or no waiting room at all. Let us consider now consider three cases of finite buffers: the M/M/1/ $k$ , the MMPP(2)/M/1/ $k$  and the M/E<sub>n</sub>/1/ $k$  Queues.

### 12.1 M/M/1/ $k$

As in the M/M/1 case, the M/M/1/ $k$  queue-size process increases by only one and decreases by only one, it is a birth and death process. However, unlike the case of the M/M/1 birth and death process the state space is infinite, in the case of the M/M/1/ $k$  birth and death process, the state space is finite limited by the buffer size. As  $k$  is the buffer size, the infinitesimal generator for the M/M/1/ $k$  queue-size process is given by

$$\begin{aligned} Q_{i,i+1} &= \lambda \text{ for } i=0, 1, 2, 3, \dots, k-1. \\ Q_{i,i-1} &= \mu \text{ for } i= 1, 2, 3, 4, \dots, k. \\ Q_{0,0} &= -\lambda \\ Q_{i,i} &= -\lambda - \mu \text{ for } i=1, 2, 3, \dots, k. \end{aligned}$$

Substituting this infinitesimal generator in Eq. (210) and performing some simple algebraic operations, we obtain the following steady state equations for the M/M/1/ $k$  queue.

$$\pi_0 \lambda = \pi_1 \mu$$

$$\pi_1 \lambda = \pi_2 \mu$$

...

and in general:

$$\pi_i \lambda = \pi_{i+1} \mu, \text{ for } i = 0, 1, 2, \dots, k-1. \quad (327)$$

The normalizing equation is:

$$\sum_{j=1}^k \pi_j = 1. \quad (328)$$

Setting  $\rho = \lambda/\mu$ , so we obtain,

$$\pi_1 = \rho \pi_0$$

$$\pi_2 = \rho \pi_1 = \rho^2 \pi_0$$

$$\pi_3 = \rho \pi_2 = \rho^3 \pi_0$$

and in general:

$$\pi_i = \rho^i \pi_0 \text{ for } i = 0, 1, 2, \dots, k. \quad (329)$$

Summing up both sides of (329), we obtain

$$1 = \sum_{i=0}^k \rho^i \pi_0 = \pi_0 \frac{1 - \rho^{k+1}}{1 - \rho}. \quad (330)$$

Therefore,

$$\pi_0 = \frac{1 - \rho}{1 - \rho^{k+1}}. \quad (331)$$

Substituting the latter in (329), we obtain

$$\pi_i = \rho^i \frac{1 - \rho}{1 - \rho^{k+1}} \text{ for } i = 0, 1, 2, \dots, k. \quad (332)$$

Of particular interest is the blocking probability  $p_k$  given by

$$\pi_k = \rho^k \frac{1 - \rho}{1 - \rho^{k+1}} = \frac{\rho^k - \rho^{k+1}}{1 - \rho^{k+1}} = \frac{\rho^k(1 - \rho)}{1 - \rho^{k+1}}. \quad (333)$$

Notice that since M/M/1/k has a finite state-space, stability is assured even if  $\rho > 1$ .

### Homework 12.1

Consider an M/M/1/k queue with  $k = \rho = 1000$ , estimate the blocking probability with error bounded by  $10^{-4}$ .  $\square$

### Homework 12.2

A well known approximate formula that links TCP's flow rate  $R_{TCP}$  [packets/sec], its round trip time (RTT) denoted  $RTT$  and TCP packet loss rate  $L_{TCP}$  is [42]:

$$R_{TCP} = \frac{1.22}{RTT \sqrt{L_{TCP}}}. \quad (334)$$

Consider a model of TCP over an M/M/1/k. That is, consider many TCP connections with a given RTT all passing through a bottleneck modeled as an M/M/1/k queue. Assuming that packet size are exponentially distributed, estimate TCP throughput, using Equations (333) and (334) for a given RTT, mean packet size and service rate of the M/M/1/k queue. Compare your results with those obtained by ns2 simulations [43].

### Guide

Use the method of iterative fixed point solution.  $\square$

## 12.2 MMPP(2)/M/1/k

The MMPP(2)/M/1/k Queue is an SSQ with buffer size  $k$  characterized by an MMPP(2) arrival process with parameters  $\lambda_0$ ,  $\lambda_1$ ,  $\delta_0$ , and  $\delta_1$ , and exponentially distributed service time with parameter  $\mu$ . The service times are mutually independent and are independent of the arrival process. Unlike the Poisson arrival process, the interarrival times in the case of the MMPP(2) process are not independent. As will be discussed, such dependency has implication on queueing performance, packet loss and utilization.

The MMPP(2)/M/1 queue process is a continuous time Markov chain, but its states are two-dimensional vectors and not scalars. Each state is characterized by two scalars: the mode  $m$  of the arrival process that can be either  $m = 0$  or  $m = 1$  and the queue size. Notice that all the other queueing systems we considered so far were based on a single dimensional state space.

Let  $p_{im}$  for  $i = 0, 1, 2, \dots, k$  be the probability that the arrival process is in mode  $m$  and that there are  $i$  packets in the system. After we obtain the  $\pi_{im}$  values, the steady-state queue size probabilities can then be obtained by

$$\pi_i = \pi_{i0} + \pi_{i1} \quad \text{for } i = 0, 1, 2, \dots, k.$$

Note that the mode process itself is a two-state continuous-time Markov chain, so the probabilities of the arrival mode being in state  $j$ , denoted  $P(m = j)$ , for  $j = 0, 1$ , can be solved using the following equations:

$$P(m = 0)\delta_0 = P(m = 1)\delta_1$$

and the normalizing equation

$$P(m = 0) + P(m = 1) = 1.$$

Solving these two equations gives

$$P(m = 0) = \frac{\delta_1}{\delta_0 + \delta_1} \tag{335}$$

$$P(m = 1) = \frac{\delta_0}{\delta_0 + \delta_1}. \tag{336}$$

Because the probability of the arrival process to be in mode  $m$  (for  $m = 0, 1$ ) is equal to  $\sum_{i=0}^k \pi_{im}$ , we obtain by (335) and (336)

$$\sum_{i=0}^k \pi_{im} = \frac{\delta_{1-m}}{\delta_{1-m} + \delta_m} \quad \text{for } m = 0, 1. \tag{337}$$

The average arrival rate denoted  $\lambda_{av}$  is given by

$$\lambda_{av} = P(m = 0)\lambda_0 + P(m = 1)\lambda_1 = \frac{\delta_1}{\delta_0 + \delta_1}\lambda_0 + \frac{\delta_0}{\delta_0 + \delta_1}\lambda_1. \tag{338}$$

Denote  $\rho = \frac{\lambda_{av}}{\mu}$ . The MMPP(2)/M/1/ $k$  queueing process is a stable, irreducible and aperiodic continuous time Markov chain with finite state space (because the buffer size  $k$  is finite). We again remind the reader that the condition  $\rho < 1$  is not required for stability in a finite buffer queueing system, or more generally, in any case of a continuous time Markov chain with finite state space. Such a system is stable even if  $\rho < 1$ .

An important performance factor in queues with MMPP(2) input is the actual time the queue stays in each mode. Even if the apportionment of time between the mode stays fixed, the actual time can make a big difference. This is especially true for the case  $\rho_1 = \lambda_1/\mu > 1$  and  $\rho_2 = \lambda_2/\mu < 1$ , or vice versa. In such a case, if the actual time of staying in each mode is long, there will be a long period of overload when a long queue is built up and/or many packets lost, followed by long periods of light traffic during which the queues are cleared. In

such a case we say that the traffic is *bursty* or strongly correlated. (As mentioned above here interarrival times are not independent.) On the other hand, if the time of staying in each mode is short; i.e., the mode process exhibits frequent fluctuations, the overall traffic process is smoothed out and normally long queues are avoided. To see this numerically one could set initially  $\delta_0 = \delta_0^*$   $\delta_1 = \delta_1^*$  where, for example,  $\delta_0 = 1$  and  $\delta_1^* = 2$ , or  $\delta_0^* = \delta_1^* = 1$ , and then set  $\delta_m = \psi \delta_m^*$  for  $m = 0, 1$ . Letting  $\psi$  move towards zero will mean infrequent fluctuations of the mode process that may lead to bursty traffic (long stay in each mode) and letting  $\psi$  move towards infinity means frequent fluctuations of the mode process. The parameter  $\psi$  is called *mode duration parameter*. In the exercises below the reader is asked to run simulations and numerical computations to obtain blocking probability and other measures for a wide range of parameter values. Varying  $\psi$  is one good way to gain insight into performance/burstiness effects.

Therefore, the  $\pi_{im}$  values can be obtain by solving the following finite set of steady state equations:

$$0 = \mathbf{\Pi} \mathbf{Q} \quad (339)$$

where  $\mathbf{\Pi} = [\pi_{00}, \pi_{01}, \pi_{10}, \pi_{11}, \pi_{20}, \pi_{21}, \dots, \pi_{k-1,0}, \pi_{k-1,1}, \pi_{k0}, \pi_{k1}]$ , and the infinitesimal generator  $2k \times 2k$  matrix is  $\mathbf{Q} = [Q_{\mathbf{i},\mathbf{j}}]$ , where  $\mathbf{i}$  and  $\mathbf{j}$  are two-dimensional vectors. Its non-zero entries are:

$$Q_{00,00} = -\lambda_0 - \delta_0; \quad Q_{00,01} = \delta_0; \quad Q_{00,10} = \lambda_0;$$

$$Q_{01,00} = \delta_1; \quad Q_{01,01} = -\lambda_1 - \delta_1; \quad Q_{01,11} = \lambda_1;$$

For  $k > i > 0$ , the non-zero entries are:

$$Q_{i0,i0} = -\lambda_0 - \delta_0 - \mu; \quad Q_{i0,i1} = \delta_0; \quad Q_{i0,(i+1,0)} = \lambda_0;$$

$$Q_{i1,i0} = \delta_1; \quad Q_{i1,i1} = -\lambda_1 - \delta_1 - \mu; \quad Q_{i1,(i+1,1)} = \lambda_1;$$

and

$$Q_{k0,(k-1,0)} = \mu; \quad Q_{k0,k0} = -\delta_0 - \mu; \quad Q_{k0,k1} = \delta_0;$$

$$Q_{k1,(k-1,1)} = \mu; \quad Q_{k1,k1} = -\delta_1 - \mu; \quad Q_{k1,k0} = \delta_1.$$

Eq. (339) can be solved by successive substitutions and then the results are normalized by the equation

$$\sum_{i=0}^k \sum_{m=0}^1 \pi_{im} = 1. \quad (340)$$

After obtaining the solution, one may verify that (337) holds.

As an example, we hereby provide the infinitesimal generator for  $k = 3$ :

|    | 00                      | 01                      | 10                            | 11                | 20                            | 21                |
|----|-------------------------|-------------------------|-------------------------------|-------------------|-------------------------------|-------------------|
| 00 | $-\lambda_0 - \delta_0$ | $\delta_0$              | $\lambda_0$                   | 0                 | 0                             | 0                 |
| 01 | $\delta_1$              | $-\lambda_1 - \delta_1$ | 0                             | $\lambda_1$       | 0                             | 0                 |
| 10 | $\mu$                   | 0                       | $-\lambda_0 - \delta_0 - \mu$ | $\delta_0$        | $\lambda_0$                   | 0                 |
| 11 | 0                       | $\mu$                   | $\delta_1$                    | $-\delta_1 - \mu$ | 0                             | $\lambda_1$       |
| 20 | 0                       | 0                       | $\mu$                         | 0                 | $-\lambda_0 - \delta_0 - \mu$ | $\delta_0$        |
| 21 | 0                       | 0                       | 0                             | $\mu$             | $\delta_1$                    | $-\delta_1 - \mu$ |

### Homework 12.3

Consider an MMPP(2)/M/1/1 queue with  $\lambda_0 = \delta_0 = 1$  and  $\lambda_1 = \delta_1 = 2$  and  $\mu = 2$ .

1. Without using a computer solve the steady state equations by standard methods to obtain  $\pi_{00}, \pi_{01}, \pi_{10}, \pi_{11}$  and verify that (337) holds.
2. Obtain the blocking Probability.
3. Find the proportion of time that the server is idle.
4. Derive an expression and a numerical value for the utilization.
5. Find the mean queue size.  $\square$

### Homework 12.4

Consider an MMPP(2)/M/1/200 queue with  $\lambda_0 = 1$ ,  $\delta_0 = 10^{-3}$ ,  $\lambda_1 = 2$ ,  $\delta_1 = 2 \times 10^{-3}$  and  $\mu = 1.9$ .

1. solve the steady state equations by successive substitutions to obtain the  $\pi_{im}$  values and verify that (337) holds.
2. Obtain the blocking Probability.
3. Find the proportion of time that the server is idle.
4. Obtain numerical value for the utilization.
5. Find the mean queue size.
6. Compare the results obtained with those obtained before for the case  $k = 1$  and explain the differences.  $\square$

### Homework 12.5

Consider an MMPP(2)/M/1/200 queue. Using successive substitutions, obtain the mean queue size for a wide range of parameter values and discuss differences. Confirm your results by simulations with confidence intervals. Compare the results with those obtained by successive substitution and simulation of an equivalent M/M/1/200 queue that has the same service rate and its arrival rate is equal to  $\lambda_a v$  of the MMPP(2)/M/1/200. Provide interpretations and explanations to all your results.  $\square$

## Homework 12.6

Consider again the MMPP(2)/M/1/200 queue and its MMPP(2)/M/1/200 equivalence. For a wide range of parameter value, compute the minimal service rate  $\mu$  obtained such that the blocking probability is no higher than  $10^{-4}$  and observe the utilization. Plot the utilization as a function of the mode duration parameter  $\psi$  to observe the effect of burstiness on the utilization. Confirm your results obtained by successive substitution by simulations using confidence intervals. Demonstrate that as  $\psi \rightarrow \infty$  the performance (blocking probability and utilization) achieved approaches that of the M/M/1/200 equivalence. Discuss and explain all the results you obtained.  $\square$

## 12.3 M/E<sub>n</sub>/1/k

We consider here an M/E<sub>n</sub>/1/k queue model characterized by a Poisson arrival process with parameter  $\lambda$ , buffer size size of  $k$ , and service time that has Erlang distribution with  $n$  phases (E<sub>n</sub>) with mean  $1/k\mu$ . Such service time model arises in situations when the standard deviation to mean ratio of the service time is lower than one (recall that for the exponential random variable this ratio is equal to one).

## Homework 12.7

Derive and plot the standard deviation to mean ratio as a function of  $n$  for an E<sub>n</sub> random variable.  $\square$

Let  $\pi_0$  be the probability that the queue is empty. Also, let  $\pi_{im}$  be the steady state distribution of having  $i$  customers in the system  $i = 1, 2, 3, \dots, k$  and that the customer in service is in phase  $m$ ,  $m = 1, 2, 3, \dots, n$ .

This model is a continuous time Markov chain, so the steady state probabilities  $\pi_0$  and  $\pi_{im}$  satisfy the following steady state equations.

$$\begin{aligned}
 0 &= -\lambda\pi_0 + k\mu\pi_{1k} \\
 0 &= -(\lambda + k\mu)\pi_{11} + k\mu\pi_{2k} + \lambda\pi_0 \\
 0 &= -(\lambda + k\mu)\pi_{1m} + k\mu\pi_{1,m+1} \quad \text{for } m = 2, 3, \dots, k \\
 0 &= -(\lambda + k\mu)\pi_{i1} + k\mu\pi_{i+1,k} + \lambda\pi_{i-1,1} \quad \text{for } i = 2, 3, \dots \\
 0 &= -(\lambda + k\mu)\pi_{im} + k\mu\pi_{i,m-1} + \lambda\pi_{i-1,m} \quad \text{for } i = 2, 3, \dots \text{ and } m = 2, 3, \dots, k.
 \end{aligned}$$

The first equation equates the probability flux of leaving state 0 (to state 11) with the probability flux of entering state 0 only from state 1k - where there is only one customer in the system who is in its last service phase. The second equation equates the probability flux of leaving state 11 with the probability flux of entering state 11 (only from states 0 and 2k). The third equation equates the probability flux of leaving state 1m ( $m > 1$ ) with the probability flux of entering state 1m ( $m > 1$ ) (only from states 0 and 1,  $m + 1$ ). The fourth equation equates the probability flux of leaving state  $i1$  ( $i > 1$ ) with the probability flux of entering state  $i1$  ( $i > 1$ ) (only from states  $i + 1, k$  and  $i - 1, 1$ ). The last equation equates the probability flux of leaving

state  $im$  ( $i > 1, m > 1$ ) with the probability flux of entering state  $im$  ( $i > 1, m > 1$ ) (only from states  $i, m - 1$  and  $i - 1, m$ ).

The probability of have  $i$  in the system denoted  $\pi_i$  is obtained by

$$\pi_i = \sum_{m=1}^n \pi_{im}.$$

The blocking probability is the probability that the buffer is full namely  $\pi_k$ . The mean queue size is obtained by

$$E[Q] = \sum_{m=1}^k i\pi_i.$$

The mean delay is obtained by little's formula:

$$E[D] = \frac{E[Q]}{\lambda}.$$

### Homework 12.8

Consider an  $M/E_n/1/k$  queue. For a wide range of parameter values (varying  $\lambda, \mu, n, k$ ) using successive substitutions, obtain the mean queue size, mean delay and blocking probability and discuss the differences. Confirm your results by simulations using confidence intervals. Provide interpretations and explanations to all your results.  $\square$

## 13 Discrete Time Queue

To complement the considerable attention we have given to continuous time queues, we will now provide an example of a discrete time queueing system. Discrete-time models are very popular in a study of computer and telecommunications systems because in some cases time is divided into fixed length intervals (time-slots) and packets of information called cells are of fixed length, such that exactly one cell can be transmitted during a time-slot. Examples of such cases include technologies such as ATM and the IEEE 802.6 Metropolitan Area Network (MAN) standard.

Let the number of cells that joins the queue at different time-slots be IID random variable. Let  $a_i$  be the probability of  $i$  cells joining the queue at the beginning of any time slot. Assume that at any time slot, if there are cells in the queue, one cell is served, namely, removed from the queue. Assuming that arrivals occur at the beginning of a time-slot means that if a cell arrives during a time-slot it can be served in the same time slot.

In this case, the queue size process follows a discrete time Markov chain with state space  $\Theta$  composed of all the nonnegative integers, and a Transition Probability Matrix  $\mathbf{P} = [P_{ij}]$  given by

$$P_{i,i-1} = a_0 \quad \text{for } i \geq 1 \quad (341)$$

and

$$P_{0,0} = a_0 + a_1$$

$$P_{i,i} = a_1 \quad \text{for } i \geq 1$$

$$P_{i,i+1} = a_2 \quad \text{for } i \geq 0$$

and in general

$$P_{i,i+k} = a_{k+1} \quad \text{for } i \geq 0, k \geq 1. \quad (342)$$

Defining the steady state probability vector by  $\mathbf{\Pi} = [\pi_0, \pi_1, \pi_2, \dots]$ , it can be obtained by solving the steady state equations:

$$\mathbf{\Pi} = \mathbf{\Pi P}.$$

together with the normalizing equation

$$\sum_{i=0}^{\infty} \pi_i = 1.$$

To solve for the  $\pi_i$ s, we will begin by writing down the steady state equations as follows



$$\pi_0 = \pi_0 P_{00} + \pi_1 P_{10}$$

$$\pi_1 = \pi_0 P_{01} + \pi_1 P_{11} + \pi_2 P_{21}$$

$$\pi_2 = \pi_0 P_{02} + \pi_1 P_{12} + \pi_2 P_{22} + \pi_3 P_{32}$$

$$\pi_3 = \pi_0 P_{03} + \pi_1 P_{13} + \pi_2 P_{23} + \pi_3 P_{33} + \pi_4 P_{43}$$

and in general

$$\pi_n = \sum_{i=0}^{n+1} \pi_i P_{i,n} \text{ for } n \geq 0.$$

Substituting (341) and (342) in the latter, we obtain

$$\pi_0 = \pi_0[a_0 + a_1] + \pi_1 a_0 \quad (343)$$

$$\pi_1 = \pi_0 a_2 + \pi_1 a_1 + \pi_2 a_0 \quad (344)$$

$$\pi_2 = \pi_0 a_3 + \pi_1 a_2 + \pi_2 a_1 + \pi_3 a_0 \quad (345)$$

and in general

$$\pi_n = \sum_{i=0}^{n+1} \pi_i a_{n+1-i} \text{ for } n \geq 1. \quad (346)$$

Defining  $\Pi(z)$  the Z-transform of the  $\Pi$  vector and  $A(z)$  as the Z-Transform of  $[a_0, a_1, a_2, \dots]$ , multiplying the  $n$ th equation of the set (343) – (346) by  $z^n$ , we obtain after some algebraic operations

$$\Pi(z) = \pi_0 a_0 + z^{-1} a_0 [\Pi(z) - \pi_0] + z^{-1} \Pi(z) [A(z) - a_0] \quad (347)$$

which leads to

$$\pi_0 = \frac{\Pi(z)[1 - z^{-1}A(z)]}{a_0(1 - z^{-1})}. \quad (348)$$

Then deriving the limit as  $z \rightarrow 1$  by applying L'Hopital rule, and denoting  $A'(1) = \lim_{z \rightarrow 1} A(z)$ , we obtain,

$$\pi_0 = \frac{1 - A'(1)}{a_0}. \quad (349)$$

This equation is somewhat puzzling. We already know that the proportion of time the server is idle must be equal to one minus the utilization. We also know that  $A'(1)$  is the mean arrival rate of the number of arrivals per time-slot and since the service rate is equal to one,  $A'(1)$  is also the utilization; so what is wrong with Eq. (349)? The answer is that nothing wrong with it. What we call  $\pi_0$  here is not the proportion of time the server is idle. It is the probability that the queue is empty at the slot boundary. There may have been one cell served in the previous slot and there may be an arrival or more in the next slot which keep the server busy.

The proportion of time the server is idle is in fact  $\pi_0 a_0$  which is the probability of empty queue at the slot boundary times the probability of no arrivals in the next slot, and the consistency of Eq. (349) follows.

### Homework 13.1

1. Provide in detail all the algebraic operations and the application of L'Hopital rule to derive equations (347), (348) and (349).

2. Derive the mean and variance of the queue size using the z-transform method and verify your result by simulations over a wide range of parameter values using confidence intervals.

□

## 14 M/G/1

The M/G/1 queue is a generalization of the M/M/1 queue where the service time is no longer exponential. We now assume that the service times are IID with mean  $1/\mu$  and standard deviation  $\sigma_s$ . The arrival process is assumed to be Poisson with rate  $\lambda$  and we will use the previously defined notation:  $\rho = \lambda/\mu$ . As in the case of M/M/1 we assume that the service times are independent and are independent of the arrival process. In addition to M/M/1, another commonly used special case of the M/G/1 queue is the M/D/1 queue where the service time is deterministic.

The generalization from M/M/1 to M/G/1 brings with it a significant increase in complexity. No longer can we use the Markov chain structure that was so useful in the previous analyzes where both service and inter-arrival times are memoryless. Without the convenient Markov chain structure, we will use different methodologies as described in this section.

### 14.1 Pollaczek Khinchin Formula: Residual Service Approach [9]

The delay of an arriving customer to an M/G/1 queue is the remaining service time of the customer in service plus the sum of the service times of all the customers in the queue ahead of the arriving customer. Therefore, the mean waiting time in the queue is given by

$$E[W_Q] = E[R] + E[N_Q]/\mu \quad (350)$$

where  $E[R]$  denotes the mean residual service time. Note that for M/M/1,  $E[R] = \rho/\mu$ , which is the probability of having one customer in service, which is equal to  $\rho$ , times the residual service time of that customer, which is equal to  $1/\mu$  due to the memoryless property of the exponential distribution, plus the probability of having no customer in service (the system is empty), which is  $1 - \rho$ , times the residual service time if there is no customer in service, which is equal to zero.

#### Homework 14.1

Verify that Eq. (350) holds for M/M/1.  $\square$

By Little's formula, (350) gives

$$E[W_Q] = \frac{E[R]}{1 - \rho}. \quad (351)$$

It remains to obtain  $E[R]$  to obtain results for the mean values of waiting time and queue-size.

Now that as the service time is generally distributed, we encounter certain interesting effects. Let us ask ourselves the following question. If we randomly inspect an M/G/1 queue, will the mean of the remaining (residual) service time of the customer in service be longer or shorter than the mean service time? A hasty response may be: shorter. Well, let us consider the following example. There are two types of customers. Each of the customers of the first type requires  $10^6$  service units, while each of the customers of the second type requires  $10^{-6}$  service units. Assume that the proportion of the customers of the first type is  $10^{-7}$ , so the proportion of the customers of the second type is  $1 - 10^{-7}$ . Assume that the capacity of the server to render service is one service unit per time unit and that the mean arrival rate is one customer

per time unit. As the mean service time is of the order of  $10^{-1}$ , and the arrival rate is one, although the server is idle 90% of the time, when it is busy it is much more likely to be busy serving a customer of the first type despite the fact that these are very rare, so the residual service time in this case is approximately  $0.1 \times 10^6/2 = 50,000$  which is much longer than the  $10^{-1}$  mean service time. Intuitively, we may make the conjecture that the residual service time is affected significantly by the variance of the service time.

Notice that what we have computed above is the unconditional mean residual service time which is our  $E[R]$ . If we condition on the event that the server is busy, the mean residual service time will be 10 times longer. We know that if the service time is exponentially distributed, the conditional residual service time of the customer in service has the same distribution as the service time due to the memoryless property of the exponential distribution. Intuitively, we may expect that if the variance of the service time is greater than its exponential equivalence (an exponential random variable with the same mean), then the mean residual service time will be longer than the mean service time. Otherwise, it will be shorter. For example, if the service time is deterministic of length  $d$ , the conditional mean residual service time is  $d/2$ , half the size of its exponential equivalence.

To compute the (unconditional) mean residual service time  $E[R]$ , consider the process  $\{R(t), t \geq 0\}$  where  $R(t)$  is the residual service time of the customer in service at time  $t$ . And consider a very long time interval  $[0, T]$ . Then

$$E[R] = \frac{1}{T} \int_0^T R(t) dt. \quad (352)$$

Following [9], let  $S(T)$  be the number of service completions by time  $T$  and  $S_i$  the  $i$ th service time. Notice that the function  $R(t)$  takes the value zero during times that there is no customer in service and jumps to the value of  $S_i$  at the point of time the  $i$ th service time commences. During a service time it linearly decreases with rate of one and reaches zero at the end of a service time. Therefore, the area under the curve  $R(t)$  is equal to sum the areas of  $S(T)$  isosceles right triangles where the side of the  $i$ th triangle is  $S_i$ . Therefore, for large  $T$ , we can ignore the last possibly incomplete triangle, so we obtain

$$E[R] = \frac{1}{T} \sum_{i=1}^{S(T)} \frac{1}{2} S_i^2 = \frac{1}{2} \frac{S(T)}{T} \frac{1}{S(T)} \sum_{i=1}^{S(T)} S_i^2. \quad (353)$$

Letting  $T$  approach infinity, the latter gives

$$E[R] = \frac{1}{2} \lambda \overline{S^2} \quad (354)$$

where  $\overline{S^2}$  is the second moment of the service time.

By (351) and (354), we obtain

$$E[W_Q] = \frac{\lambda \overline{S^2}}{2(1 - \rho)}. \quad (355)$$

Thus, considering (224), we obtain

$$E[D] = \frac{\lambda \overline{S^2}}{2(1 - \rho)} + 1/\mu. \quad (356)$$

Using Little's formula and recalling that  $\sigma_s^2 = \overline{S^2} - (1/\mu)^2$ , Eq. (356) leads to the well known Pollaczek Khinchin Formula for the mean number of customers in an M/G/1 system:

$$E[Q] = \rho + \frac{\rho^2 + \lambda^2 \sigma_s^2}{2(1 - \rho)}. \quad (357)$$

## 14.2 Pollaczek Khinchin Formula: by Kendall's Recursion [33]

Let us now derive (357) in a different way. Letting  $q_i$  be the number of customers in the system immediately following the departure of the  $i$ th customer, the following recursive relation, is obtained.

$$q_{i+1} = q_i + a_{i+1} - I(q_i) \quad (358)$$

where  $a_i$  is the number of arrivals during the service time of the  $i$ th customer, and  $I(x)$  is a function defined for  $x \geq 0$ , taking the value 1 if  $x > 0$ , and the value 0 if  $x = 0$ . This recursion was first introduced by Kendall [33], so we will call it Kendall's Recursion. Some call it a "Lindley's type Recursion" in reference to an equivalent recursion for the G/G/1 waiting time in [39]. Along with Little's and Erlang B formulae, and Pollaczek-Khinchin equation, the Kendall's and Lindley's recursions are key foundations of queueing theory.

Squaring both sides of (358) and taking expectations, we obtain

$$E[q_{i+1}^2] = E[q_i^2] + E[I(q_i)^2] + E[a_{i+1}^2] - 2E[q_i I(q_i)] + 2E[q_i a_{i+1}] - 2E[I(q_i) a_{i+1}] \quad (359)$$

Notice that in steady state  $E[q_{i+1}^2] = E[q_i^2]$ ,  $I(q_i)^2 = I(q_i) = \rho$ , and that for any  $x \geq 0$ ,  $xI(x) = x$ . Also notice that because of the independence between  $a_{i+1}$  and  $q_i$ , and because (by (79)) the mean number of arrivals during service time in M/G/1 is equal to  $\rho$ , we obtain in steady state that  $E[I(q_i) a_{i+1}] = \rho^2$  and  $E[q_i a_{i+1}] = E[q_i] \rho$ . Therefore, considering (359), and setting the steady state notation  $E[a] = E[a_i]$  and  $E[Q] = E[q_i]$ , we obtain after some algebra

$$E[Q] = \frac{\rho + E[a^2] - 2\rho^2}{2(1 - \rho)}. \quad (360)$$

To obtain  $E[a^2]$ , we notice that by EVVE,

$$\text{var}[a] = E[\text{var}[a | S]] + \text{var}[E[a | S]] = \lambda E[S] + \lambda^2 \sigma_s^2 = \rho + \lambda^2 \sigma_s^2 \quad (361)$$

recalling that  $S$  is the service time and that  $\sigma_s^2$  is its variance.

Therefore,

$$E[Q] = \frac{2\rho + \lambda^2 \sigma_s^2 - 2\rho^2}{2(1 - \rho)} \quad (362)$$

or

$$E[Q] = \rho + \frac{\rho^2 + \lambda^2 \sigma_s^2}{2(1 - \rho)} \quad (363)$$

which is identical to (357) - the Pollaczek-Khinchin Formula.

### Homework 14.2

Re-derive the Pollaczek-Khinchin Formula in the two ways presented above with attention to all the details (some of which are skipped in the above derivations).  $\square$

### 14.3 Special Cases: M/M/1 and M/D/1

Now let us consider the special case of exponential service time. That is, the M/M/1 case. To obtain  $E[Q]$  for M/M/1, we substitute  $\sigma_s^2 = 1/\mu^2$  in (357), and after some algebra, we obtain

$$E[Q] = \frac{\rho}{1 - \rho} \quad (364)$$

which is consistent with (246).

Another interesting case is the M/D/1 queue where  $\sigma_s^2 = 0$ . Substituting the latter in (357), we obtain after some algebra

$$E[Q] = \frac{\rho}{1 - \rho} \times \frac{2 - \rho}{2}. \quad (365)$$

Because the second factor of (365), namely  $(2 - \rho)/2$ , is less than one for the range  $0 < \rho < 1$ , we clearly see that the mean number of customers in an M/M/1 queue is higher than that of an M/D/1 queue.

### 14.4 Busy Period

We have defined and discussed the concept of busy period in Section 6.6 in the context of the M/M/1 queue. The same analysis applied in the case of the M/G/1 system, and we obtain:

$$E[T_B] = \frac{1}{\mu - \lambda}. \quad (366)$$

What we learn from this is that the mean busy period is insensitive to the service time distribution. In other words, the mean busy periods of M/M/1 and M/G/1 systems are the same if the mean arrival rate and service rates are the same.

#### Homework 14.3

1. Prove that  $E[T_B]/(E[T_B] + E[T_I])$  is the proportion of time that the server is busy.
2. Show that Equation (366) also applies to a G/G/1 queue.  $\square$

#### Homework 14.4

Consider an M/G/1 queueing system with the following twist. When a new customer arrives at an empty system, the server is not available immediately. The customer then rings a bell and the server arrives an exponentially distributed amount of time with parameter  $\zeta$  later. As usual customers arrive in accordance with a Poisson process with rate  $\lambda$  and the mean service time is  $1/\mu$ . Service times are mutually independent and independent of the interarrival times. Find the mean busy period defined as a continuous period that the server is busy.  $\square$

## 14.5 M/G/1 with Priorities

We have already considered a non-FIFO service policy. We mentioned the LIFO policy in our discussion of the M/M/1 queue. We will now discuss non FIFO service disciplines in the context of the M/G/1 [9].

Let us consider an M/G/1 queueing system  $m$  priority classes. Let  $\lambda_j$  and  $\mu_j$  be the arrival and service rate of customers belonging to the  $j$ th priority class for  $j = 1, 2, 3, \dots, m$ . The mean service time of customers belonging to the  $j$ th priority class is therefore equal to  $1/\mu_j$ . The second moment of the service time of customers belonging to the  $j$ th priority class is denoted  $\overline{S^2(j)}$ . We assume that that priority class  $j$  has higher priority than priority class  $j+1$ , so Class 1 represents the highest priority class and Class  $m$  the lowest. For each class  $j$ , the arrival process is assumed to be Poisson with parameter  $\lambda_j$ , and the service times are assumed mutually independent and independent of any other service times of customers belonging to the other classes, and are also independent of any interarrival times. Let  $\rho_j = \lambda_j/\mu_j$ . We assume that  $\sum_{j=1}^m \rho_j < 1$ . We will consider two priority policies: *nonpreemptive* and *preemptive resume*.

## 14.6 Nonpreemptive

Under this regime, a customer in service will complete its service even if a customer of a higher priority class arrive while it is being served. Let  $E[N_Q(j)]$  and  $E[W_Q(j)]$  represent the mean number of class  $j$  customers in the queue excluding the customer in service and the mean waiting time of a class  $j$  customer in the queue (excluding its service time), respectively. Further let  $R$  be the residual service time (of all customers of all priority classes). In similar way we derived (354), we obtain:

$$E[R] = \frac{1}{2} \sum_{j=1}^m \lambda_j \overline{S^2(j)}. \quad (367)$$

### Homework 14.5

Derive Eq. (367).  $\square$

As in Eq. (350), we have for the highest priority,

$$E[W_Q(1)] = E[R] + \frac{E[N_Q(1)]}{\mu_1} \quad (368)$$

and similar to (350) we obtain

$$E[W_Q(1)] = \frac{E[R]}{1 - \rho_1}. \quad (369)$$

Regarding the second priority,  $E[W_Q(2)]$  is the sum of the mean residual service time  $E[R]$ , the mean time it takes to serve the Class 1 customers in the queue  $E[N_Q(1)]/\mu_1$ , the mean time it takes to serve the Class 2 customers in the queue  $E[N_Q(2)]/\mu_2$ , and the mean time it takes to serve all the Class 1 customers that arrives during the waiting time in the queue for the Class 2 customer  $E[W_Q(2)]\lambda_1/\mu_1 = E[W_Q(2)]\rho_1$ . Putting it together

$$E[W_Q(2)] = R + \frac{E[N_Q(1)]}{\mu_1} + \frac{E[N_Q(2)]}{\mu_2} + E[W_Q(2)]\rho_1. \quad (370)$$

By the latter and Little's formula for Class 2 customers, namely,

$$E[N_Q(2)] = \lambda_2 E[W_Q(2)],$$

we obtain

$$E[W_Q(2)] = \frac{E[R] + \rho_1 E[W_Q(1)]}{1 - \rho_1 - \rho_2}. \quad (371)$$

By Eqs. (371) and (369), we obtain

$$E[W_Q(2)] = \frac{E[R]}{(1 - \rho_1)(1 - \rho_1 - \rho_2)}. \quad (372)$$

### Homework 14.6

Show that for  $m = 3$ ,

$$E[W_Q(3)] = \frac{E[R]}{(1 - \rho_1 - \rho_2)(1 - \rho_1 - \rho_2 - \rho_3)}. \quad (373)$$

and that in general

$$E[W_Q(j)] = \frac{E[R]}{(1 - \sum_{i=1}^{j-1} \rho_i)(1 - \sum_{i=1}^j \rho_i)}. \quad \square \quad (374)$$

The mean delay for a  $j$ th priority class customer, denoted  $E(D(j))$ , is given by

$$E[D(j)] = E[W_Q(j)] + \frac{1}{\mu_j} \text{ for } j = 1, 2, 3, \dots, m. \quad (375)$$

### Homework 14.7

Consider the case of  $m = 2$ ,  $\lambda_1 = \lambda_2 = 0.5$  with  $\mu_1 = 2$  and  $\mu_2 = 1$ . Compute the average delay for each class and the overall average delay. Then consider the case of  $m = 2$ ,  $\lambda_1 = \lambda_2 = 0.5$  with  $\mu_1 = 1$  and  $\mu_2 = 2$  and compute the average delay for each class and the overall average delay. Explain the difference between the two cases and draw conclusions. Can you generalize your conclusions?  $\square$

## 14.7 Preemptive Resume

In this case an arriving customer of priority  $j$  never waits for a customer of a lower priority class (of Class  $i$  for  $i > j$ ) to complete its service. Therefore, when we are interested in deriving the delay of a customer of priority  $j$ , we can ignore all customers of class  $i$  for all  $i > j$ . Therefore the mean delay of a priority  $j$  customer satisfies the following equation

$$E[D(j)] = \frac{1}{\mu_j} + \frac{R(j)}{1 - \sum_{i=1}^j \rho_i} + E[D(j)] \sum_{i=1}^{j-1} \rho_i \quad (376)$$

where  $R(j)$  is the mean residual time of all customers of classes  $i = 1, 2, \dots, j$  given by

$$R(j) = \frac{1}{2} \sum_{i=1}^j \lambda_i \overline{S^2(i)}.$$



The first term of Eq. (376) is simply the mean service time of a  $j$ th priority customer. The second term in the mean time it takes to clear all the customers of priority  $j$  or higher that are already in the system when a customer of Class  $j$  arrives. It is merely Eq. (351) that gives the mean time of waiting in the queue in an M/G/1 queueing system where we replace  $\rho$  of (351) by  $\sum_{i=1}^j \rho_i$  which is the total utilization of all the customers of priority  $j$  or higher. From the point of view of the  $j$ th priority customer the order of the customers ahead of it will not affect its mean delay, so we can “mix” all these customers up and consider the system as M/G/1. The first term of Eq. (376) is the mean total work introduced to the system by customers of priorities higher than  $j$  that arrive during the delay time of our  $j$  priority customer. Notice that we use the  $\rho_i$ s there because  $\rho_i = \lambda_i(1/\mu_i)$  representing the product of the mean rate of customer arrivals and the mean work they bring to the system for each priority class  $i$ .

Eq. (376) leads to

$$E[D(1)] = \frac{(1/\mu_1)(1 - \rho_1) + R(1)}{1 - \rho_1}, \quad (377)$$

and

$$E[D(j)] = \frac{(1/\mu_j)(1 - \sum_{i=1}^j \rho_i) + R(j)}{(1 - \sum_{i=1}^{j-1} \rho_i)(1 - \sum_{i=1}^j \rho_i)}. \quad (378)$$

### Homework 14.8

Derive Eqs. (377)(378).  $\square$

## 15 G/G/1

In many situations where there is non-zero correlation between inter-arrival times, the Poisson assumption for the arrival process which makes queueing models amenable to analysis does not apply. In this case, we consider the most general single-server queue - the G/G/1 queue.

We have already covered some results applicable to G/G/1. We already know that for G/G/1, the utilization  $U$  representing the proportion of time the server is busy satisfies  $U = \lambda/\mu$ . We know that G/G/1 is work conservative, and we also know that Little's formula

$$E[Q] = \lambda E[D] \quad (379)$$

is applicable to G/G/1.

### 15.1 Reich's Formula

We would like to introduce here a new and important concept the *virtual waiting time*, and a formula of wide applicability in the study of G/G/1 queues known as *Reich's formula* [8, 16, 49].

The virtual waiting time, denoted  $W_q(t)$ , is the time that a packet has to wait in the queue (not including its own service) if it arrives at time  $t$ . It is also known as *remaining workload*; meaning, the amount of work remains in the queue at time  $t$  where work is measured in time it needed to be served. We assume nothing about the interarrival times or the service process. The latter is considered as an arbitrary sequence representing the workload that each packet brings with it to the system, namely, the time required to serve each packet. For simplicity, we assume that the system is empty at time  $t = 0$ . Let  $W_a(t)$  be a function of time representing the total work arrived during the interval  $[0, t)$ . Then Reich's formula says that

$$W_q(t) = \sup_{0 \leq s < t} \{W_a(t) - W_a(s) - t + s\}. \quad (380)$$

If the queue is not empty at time  $t$ , the  $s$  value that maximizes the right hand side of (380) corresponds to the point in time where the current (at time  $t$ ) busy period started. If the queue is empty at time  $t$ , then that  $s$  value is equal to  $t$ .

#### Homework 15.1

Consider the arrival process and the corresponding service duration requirements in the following Table.

| Arrival time | Service duration (work requirement) | $W_q(t)$ | optimal $s$ |
|--------------|-------------------------------------|----------|-------------|
| 1            | 3                                   |          |             |
| 3            | 4                                   |          |             |
| 4            | 3                                   |          |             |
| 9            | 3                                   |          |             |
| 11           | 2                                   |          |             |
| 11.5         | 1                                   |          |             |
| 17           | 4                                   |          |             |

Plot the function  $W_q(t)$  for every  $t$ ,  $0 \leq t \leq 25$  and fill in the right values for  $W_q(t)$  and the optimal  $s$  for each time point in the Table.  $\square$

We have often considered the queue size probability function  $P(Q = n)$ , for  $n = 0, 1, 2, \dots$ . Its complementary distribution function is given by  $P(Q > n)$ . Note that for a G/D/1 queue we have [50]

$$Q = \lceil W_q(t) \rceil, \quad (381)$$

so if we consider  $n$  integer, then for a G/D/1 queue we have the following equality for the complementary distribution functions of the virtual waiting time  $P(W_q(t) > n)$  and the queue size [50]

$$P(Q > n) = P(W_q(t) > n), \text{ for } n = 0, 1, 2, \dots \quad (382)$$

## 15.2 The G/GI/1 Queue and Its G/GI/1/ $k$ Equivalent

Let us consider special cases of the G/G/1 and G/G/1/ $k$  queues which we call them G/GI/1 and G/GI/1/ $k$ , respectively. The GI notation indicates that the service times are mutually independent and independent of the arrival process and the state of the queue. we consider two queueing system a G/GI/1 queue and a G/GI/1/ $k$  queue that are statistically equal in every aspect except for the fact that the first has an infinite buffer and the second has a finite buffer. They both have the same arrival process the distribution of their service times and the relationship of service times to interarrival times are all statistically the same. In queueing theory there are many cases where it is easier to obtain overflow probability estimations of the unlimited buffer queue G/GI/1, namely, the steady state probability that the queue size  $Q$  exceeds a threshold  $k$ ,  $P(Q > k)$ , than to obtain the blocking probability of its G/GI/1/ $k$  equivalent denoted  $P_{loss}$ . In practice, no buffer is of unlimited size, so the more important problem in applications is the blocking probability of a G/GI/1/ $k$  queue.

This gives rise to the following problem. Given  $P(Q > k)$  for a G/GI/1 queue, what can we say about the blocking probability of the G/GI/1/ $k$  equivalent. Let us begin with two examples. First, consider a discrete-time single-server queueing model where time is divided into fixed-length intervals called slots. Assume that the service time is deterministic and is equal to a single slot. Let the arrival process be described as follows:  $10^9$  packets arrive at the first time-slot and no packets arrived later. Consider the case of  $k = 1$ . In the finite buffer case, almost all the  $10^9$  packets that arrived are lost because the buffer can store only one packet. Therefore,  $P_{loss} \approx 1$ . However, for the case of  $P(W_q(t) > k)$  the case is completely the opposite. After the  $10^9$  time-slots that it takes to serve the initial burst the queue is empty forever, so in steady-state  $P(W_q(t) > k) = 0$ .

On the other hand, consider another discrete-time queueing model with  $k = 10^9$  and a server that may serve  $10^9$  customers in a time-slot with probability  $1 - 10^{-9}$  and may serve  $10^{90}$  customers in a time-slot with probability  $10^{-9}$ . (The rare high service rate is just to assure stability.) Assume that at a beginning of every time-slot,  $10^9 + 1$  customers arrives at the buffer before service takes place. This means that one out of the arriving  $10^9 + 1$  customers is lost, thus  $P_{loss} \approx 10^{-9}$ , while  $P(W_q(t) > k) \approx 1$ . We conclude that  $P_{loss}$  and  $P(W_q(t) > k)$  can be very different.

Wong [59] considered this problem in the context of an ATM multiplexer fed by multiple deterministic flows (a queueing model denoted N\*D/D/1 and its finite buffer equivalent) and

obtained the following relationship.

$$\rho P_{loss} \leq P(Q > k) \quad (383)$$

Roberts et al. [50] argued that it can be generalized to G/D/1 and its G/D/1/ $k$  equivalent. This can be further generalized. The arguments are analogous to those made in [59]. Let  $\lambda$  be the arrival rate and  $\mu$  the service rate in both the G/GI/1 queue and its G/GI/1/ $k$  equivalent, with  $\rho = \lambda/\mu$ . Consider a continuous period of time, in our G/GI/1 queue, during which  $Q > k$  and that just before it begins and just after it ends  $Q \leq k$ , and define such time period as *overflow period*. Since the queue size at the beginning is the same as at the end of the overflow period, the number of customers that joined the queue during an overflow period must be equal to the number of customers served during the overflow period, as the server is continuously busy during an overflow period.

Now consider a G/GI/1/ $k$  queue that has the same realization of arrivals and their work requirements as the G/GI/1 queue. Let us argue that in the worst case, the number of lost customers in the G/GI/1/ $k$  queue is maximized if all customers that arrive during an overflow period of the are lost. If for a given G/GI/1 overflow period, not all customers in the G/GI/1/ $k$  queue are lost, the losses are reduced from that maximum level without increasing future losses because at the end of a G/GI/1 overflow period, the number of customers in the G/GI/1/ $k$  queue can never be more than  $k$ .

Consider a long period of time  $L$ , the mean number of lost customers the G/GI/1/ $k$  queue during  $L$  is  $\lambda L P_{loss}$ . This must be lower or equal to the number of customers that arrived during the G/GI/1 overflow period which is equal, as discussed above, to the number of customers served during the G/GI/1 overflow period. Therefore,

$$\lambda L P_{loss} \leq \mu L P(Q > k)$$

and Eq. (383) follows.  $\square$

## Homework 15.2

Show that (383) applies to a M/M/1 queue and its M/M/1/ $k$  Equivalent, and discuss how tight is the bound in this case for the complete range of parameter values.

## Guide

Recall that for M/M/1/ $k$ ,

$$P_{loss} = \frac{\rho^k(1 - \rho)}{1 - \rho^{k+1}},$$

and for M/M/1,

$$P(Q > k) = \rho^{k+1}(1 - \rho) + \rho^{k+2}(1 - \rho) + \rho^{k+3}(1 - \rho) + \dots = \rho^{k+1}. \quad \square$$

## Homework 15.3

Using the UNIX command *netstat* collect a sequence of 100,000 numbers representing the number of packets arriving recorded every second for consecutive 100,000 seconds. Assume

that these numbers represent the amount of work, measured in packets, which arrive at an SSQ during 100,000 consecutive time-intervals. Write a simulation of an SSQ fed by this arrival process, assume that all the packets are of equal length and compute the Packet Loss Ratio (PLR) for a range of buffer sizes and the overflow probabilities for a range of thresholds. PLRs are relevant in the case of a finite buffer queue and overflow probabilities represent the probability of exceeding a threshold in an infinite buffer queue. Plot the results in two curves one for the PLR and the other for the overflow probabilities times  $\rho^{-1}$  and observe and discuss the relationship between the two.  $\square$

#### Homework 15.4

Consider the sequence of 100,000 numbers you have collected. Let  $E[A]$  be their average. Generate a sequence of 100,000 independent random numbers governed by a Poisson distribution with mean  $\lambda = E[A]$ . Use your SSQ simulation, and compute the PLR for a range of buffer sizes, and the overflow probabilities for a range of thresholds. Compare your results to those obtained in the previous Assignment, and try to explain the differences.  $\square$

#### Homework 15.5

In this exercise the reader is asked to repeat the previous homework assignment for the Bernoulli process. Again, consider the sequence of 100,000 numbers you have collected. Let  $E[A]$  be their average. Generate a sequence of 100,000 independent random numbers governed by the Bernoulli distribution with mean  $p = E[A]$ . Use your SSQ simulation from Exercise 1, and compute the PLR for a range of buffer sizes, and the overflow probabilities for a range of thresholds. Compare your results to those obtained previously, and discuss the differences.  $\square$

## 16 Queueing Networks

So far we have considered various queueing systems, but in each case we have considered a single queueing system in isolation. Very important and interesting models involve networks of queues. One important application is the Internet itself. It may be viewed as a network of queueing systems where all network elements such as routers and computers are connected and where the packets are the customers served by the various network elements and are often queued there waiting for service.

Queueing network models can be classified into two groups: (1) open queueing networks, and (2) closed queueing networks. In closed queueing networks the same customers stay in the network all the time. No new customers join and no customer leaves the network. Customers that complete their service in one queueing system goes to another and then to another and so forth, and never leaves the network. In open queueing systems new customers from the outside of the network can join any queue, and when they complete their service in the network obtaining service from an arbitrary number of queueing system they may leave the network. In this section we will only consider open queueing networks.

### 16.1 Jackson Networks

Consider a network made of Markovian queues ( $M/M/1$ ,  $M/M/k$  and  $M/M/\infty$ ). An issue that is very important for such Markovian queueing networks is what is the output of such queues because in queueing networks output of one queue may be the input of another. The so-called Burke's Theorem answers this question. Burke's Theorem states that, in steady state, the output (departure) process of  $M/M/1$ ,  $M/M/k$  or  $M/M/\infty$  queue is Poisson. Because no traffic is lost in such queues, the arrival rate must be equal to the departure rate, then any  $M/M/1$ ,  $M/M/k$ , or  $M/M/\infty$  queue with arrival rate of  $\lambda$  will have a Poisson departure process with rate  $\lambda$  in steady state.

Having information about the output processes, we will now consider an example of a very simple queueing network made of two identical single-server queues in series, in steady state. Let us assume that all the traffic arrives into the first queue following a Poisson process with parameter  $\lambda$ . The service times required by each of the arriving customers at the two queues are independent and exponentially distributed with parameter  $\mu$ . This means that the amount of time a customer requires in the first queue is independent of the amount of time a customer requires in the second queue and they are both independent of the arrival process into the first queue. Since the output process of the first queue is Poisson with parameter  $\lambda$ , and since the first queue is clearly an  $M/M/1$  queue, we have here nothing but two identical  $M/M/1$  queues in series. This is an example of a network of queues where Burke's theorem [13] leads immediately to a solution for queue size and waiting time statistics. A class of networks that can be easily analyzed this way is the class of the so-called acyclic networks. These networks are characterized by the fact that a customer never goes to the same queue twice for service.

If the network is not acyclic, the independence between inter arrival times and between inter arrival and service times do not hold any longer. This means that the queues are no longer Markovians. To illustrate this let us consider a network of two single-server queues called Q1 and Q2. The service times in Q1 and Q2 are exponentially distributed with parameters  $\mu_1 = 1$  and  $\mu_2 = 1$ , respectively. Assume that the arrival rate from the outside into Q1 follows Poisson

process with rate  $r = 10^{-8}$ . Further assume that all the packets that leave Q1 enter Q2 and that every packet that leaves Q2 leaves the network with probability  $10^{-3}$  and return to Q1 with probability  $1 - 10^{-3}$ . The total arrival process into Q1 will include the original Poisson stream with rate  $r = 10^{-8}$  plus all the feedback from Q2. This results in a process based on very infrequent original arrivals each of which brings with it a burst of mean of some thousand feedback arrivals from Q2. Clearly this is not a Poisson process. Furthermore, the inter-arrivals of packets within a burst, most of which are feedback from Q2, are very much dependent on the service times of Q1 and Q2, so clearly we have dependence between inter-arrival times and service times.

Nevertheless, the so-called Jackson's Theorem extends the above result to networks that are not acyclic. In other words, although the queues are not M/M/1 (or M/M/ $k$  or M/M/ $k$ ), they behave in terms of their queue-size and waiting time statistics as if they are.

Consider a network of  $N$  single-server queues in steady state. The result can easily be extended to multi-server queues such as M/M/ $k$  and M/M/ $\infty$ , but let us consider single-server queues for now. For queue  $i$ ,  $i = 1, 2, 3, \dots, N$ , the arrival process is Poisson with rate  $A_i$ . We allow for  $A_i = 0$  for some queues, but there must be at least one queue  $j$ , such that  $A_j > 0$ . Once a customer completes its service in queue  $i$ , it continues to queue  $j$  with probability  $P_{ij}$ , or leaves the system with probability  $1 - \sum_{j=1}^N P_{ij}$ . Notice that we allow for  $P_{ii} > 0$  for some queues. That is, we allow for positive probability for customers to return to the same queue they just exited.

Let  $\lambda_j$  be the total arrival rate into queue  $j$ . These arrival rates can be computed by solving the following set of equations.

$$\lambda_j = A_j + \sum_{i=1}^N \lambda_i P_{ij}, \quad j = 1, 2, 3, \dots, N. \quad (384)$$

The above set of equations can be solved uniquely, if every customer eventually leaves the network. This means that the routing probabilities  $P_{ij}$  must be such that there is a sequence of positive routing probabilities and a final exit probability that create an exit path of positive probability from each node.

The service times at the  $j$ th queue are assumed mutually independent and independent of the arrival process at each queue. The service times at the  $j$ th queue are assumed exponentially distributed with parameter  $\mu_j$  and mutually independent and are also assumed independent of the arrival process at that queue. Let  $\rho_j$  be defined by

$$\rho_j = \frac{\lambda_j}{\mu_j} \quad \text{for } j = 1, 2, 3, \dots, N. \quad (385)$$

Let  $Q_j$  be the queue-size of queue  $j$ . Then according to Jackson's Theorem, in steady state, we have that

$$P(Q_1 = k_1, Q_2 = k_2, \dots, Q_N = k_N) = P(k_1)P(k_2)P(k_3) \cdot \dots \cdot P(k_N) \quad (386)$$

where  $P(k_i) = \rho_i^{k_i}(1 - \rho_i)$ , for  $i = 1, 2, 3, \dots, N$ . In other words, the queues behave as M/M/1 queues despite the fact that the network may be acyclic in which case the queues are not M/M/1 queues. Therefore, the mean queue-size of the  $j$ th queue is given by

$$E[Q_j] = \frac{\rho_j}{1 - \rho_j}. \quad (387)$$

The mean delay of a customer in the  $j$ th queue  $E[D_j]$  can be obtain by Little's formula as follows.

$$E[D_j] = \frac{Q_j}{\lambda_j}. \quad (388)$$

Using Little's formula, by considering the entire queueing network as our system, we can also derive the mean delay of an arbitrary customer  $E[D]$ :

$$E[D] = \frac{\sum_{j=1}^N Q_j}{\sum_{j=1}^N A_j}. \quad (389)$$

Let us now consider the above-mentioned two-queue network example. Using our notation, we have  $A_1 = 10^{-8}$  and  $A_2 = 0$ ;  $\mu_1 = \mu_2 = 1$ ; further,

$$\lambda_1 = A_1 + (1 - 10^{-3})\lambda_2$$

and

$$\lambda_2 = \lambda_1.$$

Thus,

$$\lambda_1 = 10^{-8} + (1 - 10^{-3})\lambda_1,$$

so

$$\lambda_1 = \lambda_2 = 10^{-5}$$

and

$$\rho_1 = \rho_2 = 10^{-5},$$

so

$$E[Q1] = E[Q2] = \frac{10^{-5}}{1 - 10^{-5}} \approx 10^{-5}$$

and

$$E[D1] = E[D2] \approx \frac{10^{-5}}{10^{-5}} = 1.$$

This means that negligible queueing delay is expected. (The use word “negligible” instead of zero is used because of the approximation  $1 - 10^{-5} \approx 1$  made above.) This result makes sense intuitively. Although the feedbacked traffic is more bursty than Poisson we are considering here the same packet that returns over and over again and it is impossible for the same packet to wait in the queue for itself to be served.

An open network of M/M/1 or M/M/ $k$  or M/M/ $\infty$  queues described above is called a Jackson Network. For such network an exact solution is available. However, in most practical cases, especially when we have to deal with the so-called loss networks that comprise queues such as M/M/ $k/k$  where traffic is lost, we have to make additional modelling assumptions and to rely on approximations to evaluate measures such as blocking probability or carried traffic. One useful approximation is the so-called Reduced-Load Erlang Fixed-Point Approximation which is reasonably accurate and very useful for loss networks.



### Homework 16.1

Consider a 6-node network of M/M/1 queues, the service rate of all the queues is equal to one, i.e.,  $\mu_i = 1$  for  $i = 1, 2, 3, \dots, 6$ . The arrival rates from the outside into the different queues is given by  $r_1 = 0.6$ ,  $r_2 = 0.5$ , and  $r_i = 0$  for  $i = 3, 4, 5, 6$ . The routing matrix is as follows

|   | 1 | 2   | 3   | 4   | 5   | 6   |
|---|---|-----|-----|-----|-----|-----|
| 1 | 0 | 0.4 | 0.6 | 0   | 0   | 0   |
| 2 | 0 | 0.1 | 0   | 0.7 | 0.2 | 0   |
| 3 | 0 | 0   | 0   | 0.3 | 0.7 | 0   |
| 4 | 0 | 0   | 0   | 0   | 0   | 0.6 |
| 5 | 0 | 0   | 0   | 0.3 | 0   | 0.2 |
| 6 | 0 | 0   | 0.3 | 0   | 0   | 0   |

1. Find the mean delay in each of the queues.
2. Find the mean time a packet spends in the network from the moment it enters the network until it leaves the network.
3. Find the probability that the entire network is empty.  $\square$

## 16.2 Erlang Fixed-Point Approximation

Let us consider a circuit switched network made of nodes (switching centers) that are connected by links. Each link has a fixed number of circuits. In order to make a call between two nodes, a user should reserve a free circuit in each consecutive link of a path between the two nodes. Such reservation is successful if and only if there exists a free circuit on each of the links of that path.

To evaluate the probability that a circuit reservation is blocked we first make the following simplifying assumptions:

1. all the links are independent,
2. the arrival process of calls for each origin destination is Poisson, and
3. the arrival process seen by each link is Poisson.

Having made these assumptions, we now consider each link as an independent M/M/k/k system for which the blocking probability is readily available. Now that we have means to obtain the blocking probability on each link, we will explain how they can be used to compute the blocking probability of a call made on a given route. Let  $B_R$  be the blocking probability of a call made on route  $R$ . The route  $R$  can be viewed as an ordered set of links. Let  $B(j)$  be the blocking probability on link  $j$  obtained by Erlang formula where the traffic load on that link may include traffic from many other traffic routes not only the one we consider. Recalling that multiplexing of Poisson processes give another Poisson process which its rate is the sum of the individual rates we can compute the total traffic in Erlang offered to any link. There are several ways to implement

In particular, using the notation of Section 16.1, in a similar way to (384), the  $\lambda_j$  values (designating the total arrival rate into queue  $j$ ) can be computed by iteratively solving the following set of equations.

$$\lambda_j = A_j + \sum_{i=1}^N \lambda_i (1 - B(i)) P_{ij}, \quad j = 1, 2, 3, \dots, N. \quad (390)$$

To solve these equations, we start with an initial vector of  $\lambda_j$  values, computing the corresponding values of  $B(j)$  using for example the Erlang B formula and a new set of  $\lambda_j$  values by (390). We continue iteratively until two consecutive sets of  $\lambda_j$  values are close enough. This results also in a set of  $B(j)$  values. The  $P_{ij}$  values are initially set based on the traffic in routes assuming no losses. After solving (390) by the fixed point iterations just described, a new set of the  $P_{ij}$  values can be obtained considering the resulting blocking probability values ( $B(j)$ ). The process will repeat until two consecutive sets of  $B(j)$  values are close enough, and the fixed point solution is achieved for a given set of routes.

Then, by the independence assumption, we obtain an evaluation of the blocking probability by

$$B_R = 1 - \prod_{j \in R} (1 - B(j)). \quad (391)$$

The above solution based on the principles of the Reduced-Load and Erlang Fixed-Point Approximations can be applied to modelling a cellular mobile networks where each cell is equivalent to one M/M/k/k (or M/G/k/k) system, so the cellular mobile network is modelled by a network of M/M/k/k queues. Generation of new calls in a cell is equivalent to arrivals into an M/M/k/k queue and handover between cells is equivalent to traffic that completes service in one M/M/k/k system and moves to another. Another application is an Optical Burst Switching (OBS) network [51] where bursts of data are moving between OBS nodes each of which is modelled as an M/M/k/k system.

Another important applications of the Erlang Fixed-Point Approximation (also called Reduced-Load approximation) are circuit switched networks where a further step is required. For circuit switching networks having the blocking probability  $B_R$ , for each route  $R$ , implies that the traffic on the entire route  $R$  will be reduced by a factor of  $(1 - B_R)$ . Updating all the traffic loads on all the routes will give a new set of values for  $B(j)$ ,  $j \in R$ . We update these values in (391) to obtain a new value for  $B_R$ . We repeat this process until the updated value of  $B_R$  is arbitrarily close to its previous value.

### 16.3 A Markov Chain Simulation of a Mobile Cellular Network

A mobile cellular network can be modelled as a network of M/M/k/k systems by assuming that the number of channels in each cell is fixed and equal to  $k$ , that new call generations in each cell follows a Poisson process, that call holding times are exponentially distributed and that times until handover occurs in each cell are also exponentially distributed. In the following we describe how to simulate such a network.

Variables and input parameters:

$N$  = total of M/M/k/k Systems (cells) in the network;

$Q(i)$  = number of customers (queue size) in cell  $i$  ;

$B_p$  = estimation for the blocking probability;

$N_a(i)$  = number of customer arrivals counted so far in cell  $i$ ;

$N_b(i)$  = number of blocked customers counted so far in cell  $i$ ;

$MAXN_a$  = maximal number of customers - used as a stopping criterion;

$\mu(i)$  = service rate in cell  $i$ ;

$\lambda(i)$  = arrival rate of new calls in cell  $i$ ;

$P(i, j)$  = Matrix of routing probabilities;

$\delta(i)$  = handover rate in cell  $i$

$P_B$  = Blocking probability.

Again, we will repeatedly consider  $R(01)$  a uniform  $U(0, 1)$  random deviate. A new value for  $R(01)$  is generated every time it is called.

To know if the next event is an arrival, we use the following **if** statement.

If

$$R(01) \leq \frac{\sum_{i=1}^N \lambda(i)}{\sum_{i=1}^N \lambda(i) + \sum_{i=1}^N Q(i)\mu(i) + \sum_{i=1}^N \delta(i)}$$

then the next event is an arrival. Else, to find out if it is a departure (it could also be a handover) we use the following **if** statement. If

$$R(01) \leq \frac{\sum_{i=1}^N \lambda(i) + \sum_{i=1}^N Q(i)\mu(i)}{\sum_{i=1}^N \lambda(i) + \sum_{i=1}^N Q(i)\mu(i) + \sum_{i=1}^N \delta(i)}$$

then the next event is a departure; else, it is a handover.

If the next event is an arrival, we need to know in which of the  $N$  cells it occurs. To find out, we use the following loop.

For  $i = 1$  to  $N$ , do: If

$$R(01) \leq \frac{\sum_{j=1}^i \lambda(j)}{\sum_{j=1}^N \lambda(j)},$$

stop the loop. The arrival occurs in cell  $i$ , so if  $\sum_{j=1}^N S_a(j) = MAXN_a$ , the simulation ends, so we compute the blocking probabilities as follows.

$$P_B = \frac{\sum_{i=1}^N N_b(i)}{\sum_{i=1}^N \lambda(i)}.$$

Else,  $S_a(i) = S_a(i) + 1$  and if  $Q(i) < k$  then  $Q(i) = Q(i) + 1$ , else the number of lost calls needs to be incremented, namely,  $N_b(i) = N_b(i) + 1$ .

If the next event is a departure, we need to know in which of the  $N$  cells it occurs. To find out we use the following loop.

For  $i = 1$  to  $N$ , do: If

$$R(01) \leq \frac{\sum_{j=1}^i Q(j)\mu(j)}{\sum_{j=1}^N Q(j)\mu(j)}.$$

Then stop the loop. The departure occurs in System  $i$ , so  $Q(j) = Q(j) - 1$ . Note that we do not need to verify that  $Q(j) > 0$  (why?).

If the next event is a handover, we need to know out of which of the  $N$  cells it handovered. To find out we use the following loop.

For  $i = 1$  to  $N$ , do: If

$$R(01) \leq \frac{\sum_{j=1}^i Q(j)\mu(j)}{\sum_{j=1}^N Q(j)\delta(j)}.$$

Then stop the loop. The handover occurs out of cell  $i$ , so  $Q(j) = Q(j) - 1$ . Note that again we do not need to verify that  $Q(j) > 0$ .

## 17 Stochastic Processes as Traffic Models

In general, the aim of traffic modelling is to provide the network designer with relatively simple means to characterize traffic load on a network. Ideally, such means can be used to estimate performance and to enable efficient provisioning of network resources. Modelling a traffic stream emitted from a source, or a traffic stream that represents a multiplexing of many Internet traffic streams, is part of traffic modelling. It is normally reduced to finding a stochastic process that behaves like the real traffic stream from the point of view of the way it affects network performance or provides QoS to customers.

### 17.1 Parameter Fitting

One way to choose such a stochastic process is by fitting its statistical characteristics to those of the real traffic stream. Consider time to be divided into fixed length consecutive intervals, and consider the number of packets arriving during each time interval as the real traffic stream. Then, the model of this traffic stream could be a stationary discrete time stochastic process  $\{X_n, n \geq 0\}$ , with similar statistical characteristics as those of the real traffic stream. In this case,  $X_n$  could be a random variable representing the number of packets that arrive in the  $n$ th interval. Let  $S_n$  be a random variable representing the number of packets arriving in  $n$  consecutive intervals. We may consider the following for fitting between the statistics of  $\{X_n, n \geq 0\}$  and those of the real traffic stream:

- The mean  $E[X_n]$ .
- The variance  $var[X_n]$ .
- The the AVR discussed in Section 2.1. The AVR is related to the so-called Index of Dispersion for Counts (IDC) [25] as follows: the AVR is equal to  $E[X_n]$  times the IDC.

A stationary stochastic process  $\{X_n, n \geq 0\}$ , where autocorrelation function decays slower than exponential is said to be Long Range Dependent (LRD). Notice that if the autocovariance sum  $\sum_{k=1}^{\infty} Cov(X_1, X_k)$  is infinite the autocorrelation function must decay slower than exponential, so the process is LRD. In such processes the use of AVR (or IDC) may not appropriate because it is not finite, so a time dependent version of the IDC, i.e.,  $IDC(n) = var[S_n]/E[X_n]$  may be considered. Another statistical characteristic that is suitable for LRD processes is the so-called *Hurst parameter* denoted by  $H$  for the range  $0 \leq H < 1$  that satisfies

$$\lim_{n \rightarrow \infty} \frac{var[S_n]}{\alpha n^{2H}} = 1. \quad (392)$$

Each of these statistical parameters have their respective continuous time counterparts. As the concepts are equivalent, we do not present them here. We will discuss now a few examples of stochastic processes (out of many more available in the literature) that have been considered as traffic models.

## 17.2 Poisson Process

For many years the Poisson process has been used as a traffic model for the arrival process of phone calls at a telephone exchange. The Poisson process is characterized by one parameter  $\lambda$ , and  $\lambda t$  is the mean as well as the variance of the number of occurrences during any time interval of length  $t$ . Its memoryless nature makes it amenable to analysis as noticed through the analyzes of the above-mentioned queueing systems. Its ability to characterize telephone traffic well, being characterized by a single parameter, and its memoryless nature which makes it so amenable to analysis have made the Poisson process very useful in design and dimensioning of telephone networks.

By its nature, the Poisson process can accurately model events generated by a large number of independent sources each of which generating relatively sparsely spaced events. Such events could include phone calls or generation of Internet traffic flows. For example, a download of a page could be considered such a traffic flow. However, it cannot accurately model a packet traffic stream generated by a single user or a small number of users. It is important to note here that many textbooks and practitioners do consider the Poisson process as a model of a packet traffic stream (despite the inaccuracy it introduces) due to its nice analytical properties.

Normally, the Poisson process is defined as a continuous time process. However, in many cases, it is used as a model for a discrete sequence of a traffic stream by considering time to be divided into fixed length intervals each of size one (i.e.,  $t = 1$ ), and simply to generate a sequence of independent random numbers which are governed by a Poisson distribution with mean  $\lambda$  where  $\lambda$  is equal to the average of the sequence we try to model. As we fit only one parameter here, namely the mean, such model will not have the same variance, and because of the independence property of the Poisson process, it will not mimic the autocorrelation function of the real process. In an assignment below, you will be asked to demonstrate that such process does not lead to a similar queueing curves as the real traffic stream.

## 17.3 Markov Modulated Poisson Process (MMPP)

Traffic models based on MMPP have been used to model bursty traffic. Due to its Markovian structure together with its versatility, the MMPP can capture bursty traffic statistics better than the Poisson process and still be amenable to queueing analysis. The simplest MMPP model is MMPP(2) with only four parameters:  $\lambda_0$ ,  $\lambda_1$ ,  $\delta_0$ , and  $\delta_1$ .

Queueing models involving MMPP input have been analyzed in the 70s and 80s using Z-transform [60, 61, 62, 63]. Neuts developed matrix methods to analyse such queues [46]. For applications of these matrix methods for Queueing models involving MMPP and the use of MMPP in traffic modelling and its related parameter fitting of MMPP the reader is referred to [21, 25, 37, 45].

## 17.4 Autoregressive Gaussian Process

A traffic model based on a Gaussian process can be described as a traffic process where the amount of traffic generated within any time interval has a Gaussian distribution. There are several ways to represent a Gaussian process. The Gaussian auto-regressive is one of them.

Also, in many engineering applications, the Gaussian process is described as a continuous time process. In this section, we shall define the process as a discrete time.

Let time be divided into fixed length intervals. Let  $X_n$  be a continuous random variable representing the amount of work entering the system during the  $n$ th interval.

According to the Gaussian Autoregressive model we assume that  $X_n$ ,  $n = 1, 2, 3 \dots$  is the so-called  $k$ th order autoregressive process, defined by

$$X_n = a_1 X_{n-1} + a_2 X_{n-2} + \dots + a_k X_{n-k} + b \tilde{G}_n, \quad (393)$$

where  $\tilde{G}_n$  is a sequence of IID Gaussian random variables each with mean  $\eta$  and variance 1, and  $a_i$  ( $i = 1, 2, \dots, k$ ) and  $b$  are real numbers with  $|a| < 1$ .

In order to characterize real traffic, we will need to find the best fit for the parameters  $a_1, \dots, a_k, b$ , and  $\eta$ . On the other hand, it has been shown in [3], [4], [5] that in any Gaussian process only three parameters are sufficient to estimate queueing performance to a reasonable degree of accuracy. It is therefore sufficient to reduce the complexity involved in fitting many parameters and use only the 1st order autoregressive process, also called the AR(1) process. In this case we assume that the  $X_n$  process is given by

$$X_n = a X_{n-1} + b \tilde{G}_n, \quad (394)$$

where  $\tilde{G}_n$  is again a sequence of IID Gaussian random variables with mean  $\eta$  and variance 1, and  $a$  and  $b$  are real numbers with  $|a| < 1$ . Let  $\lambda = E[X_n]$  and  $\sigma^2 = \text{var}[X_n]$ . The AR(1) process was proposed in [41] as a model of a VBR traffic stream generated by a single source of video telephony.

The  $X_n$ s can be negative with positive probability. This may seem to hinder the application of this model to real traffic processes. However, in modeling traffic, we are not necessarily interested in a process which is similar in every detail to the real traffic. What we are interested in is a process which has the property that when it is fed into a queue, the queueing performance is sufficiently close to that of the queue fed by the real traffic.

Fitting of the parameters  $a$ ,  $b$  and  $\eta$  with measurable (estimated) parameters of the process  $\lambda$ ,  $\sigma^2$  and  $S$ , are provided based on [4]:

$$a = \frac{S}{S + \sigma^2} \quad (395)$$

$$b = \sigma^2(1 - a^2) \quad (396)$$

$$\eta = \frac{(1 - a)\lambda}{b} \quad (397)$$

where  $S$  is the autocovariance sum given by Eq. (153).

## 17.5 Exponential Autoregressive (1) Process

In the previous section we considered an autoregressive process which is Gaussian. What made it a Gaussian process was that the so-called *innovation process*, which in the case of the previous section was the sequence  $b\tilde{G}_n$ , was a sequence of Gaussian random variables. Letting  $D_n$  be

a sequence of inter-arrival times, here we consider another AR(1) process called *Exponential Autoregressive (1)* (EAR(1)) [23], defined as follows:

$$D_n = aD_{n-1} + I_n E_n, \quad (398)$$

where  $D_0 = I_0$ ,  $\{I_n\}$  is a sequence of IID random variables in which  $P(I_n = 1) = 1 - a$  and  $P(I_n = 0) = a$ , and  $\{E_n\}$  is a sequence of IID exponential random variables with parameter  $\lambda$ .

The EAR(1) has many nice and useful properties. The  $\{D_n\}$  process is a sequence of exponential random variables with parameter  $\lambda$ . These are IID only for the case  $a = 0$ . That is, when  $a = 0$ , the  $\{D_n\}$  is a sequence of inter-arrival times of a Poisson process. The autocorrelation function of  $\{D_n\}$  is given by

$$C_{EAR1}(k) = a^k. \quad (399)$$

It is very easy to simulate the  $\{D_n\}$  process, so it is useful in demonstrating by simulation the relationship between correlation in the arrival process and queueing performance.

### Homework 17.1

Prove that  $D_n$  is exponentially distributed for all  $n \geq 0$ .

### Guide

Knowing that the statement is true for  $D_0$ , prove that the statement is true for  $D_1$ . Let  $\mathcal{L}_X(s)$  be the Laplace transform of random variable  $X$ . By definition,  $\mathcal{L}_X(s) = E[e^{-sX}]$ , so  $\mathcal{L}_{I_1 E_1}(s) = E[e^{-sI_1 E_1}]$ . Thus, by (79),  $\mathcal{L}_{I_1 E_1}(s) = P(I = 1)E[e^{-sE_1}] + P(I = 0)E[e^{-0}] = (1 - a)\lambda/(\lambda + s) + a$ . By definition,  $\mathcal{L}_{D_1}(s) = E[e^{-s(aD_0 + I_1 E_1)}] = \mathcal{L}_{D_0}(as)\mathcal{L}_{I_1 E_1}(s)$ . Recall that  $D_0$  is exponentially distributed with parameter  $\lambda$ , so  $\mathcal{L}_{D_0}(as) = \lambda/(\lambda + as)$ . Use the above to show that  $\mathcal{L}_{D_1}(s) = \lambda/(\lambda + s)$ . This proves that  $D_1$  is exponentially distributed. Use the recursion to prove that  $D_n$  is exponentially distributed for all  $n > 1$ .  $\square$

## 17.6 Poisson Pareto Burst Process

Unlike the previous models, the Poisson Pareto Burst Process (PPBP) is Long Range Dependent (LRD). The PPBP has been proposed as a more realistic model for Internet traffic than its predecessors. According to this model, bursts of data (e.g. files) are generated in accordance with a Poisson process with parameter  $\lambda$ . The size of any of these bursts has a Pareto distribution, and each of them is transmitted at fixed rate  $r$ . At any point in time, we may have any number of sources transmitting at rate  $r$  simultaneously because according to the model, new sources may start transmission while others are active. If  $m$  sources are simultaneously active, the total rate equals  $mr$ . A further generalization of this model is the case where the burst lengths are generally distributed. In this case, the amount of work introduced by this model as a function of time is equivalent to the evolution of an M/G/ $\infty$  queueing system. Having  $m$  sources simultaneously active is equivalent to having  $m$  servers busy in an M/G/ $\infty$  system. M/G/ $\infty$  which is a name of a queueing system is also often use to describe the above describe traffic model. The PPBP is sometimes called M/Pareto/ $\infty$  or simply M/Pareto [2].



Again, let time be divided into fixed length intervals, and let  $X_n$  be a continuous random variable representing the amount of work entering the system during the  $n$ th interval. For convenience, we assume that the rate  $r$  is the amount transmitted by a single source within one time interval if the source was active during the entire interval. We also assume that the Poisson rate  $\lambda$  is per time interval. That is, the total number of transmissions to start in one time interval is  $\lambda$ .

To find the mean of  $X_n$  for the PPBP process, we consider the total amount of work generated in one time interval. The reader may notice that the mean of the total amount of work generated in one time interval is equal to the mean of the amount of work transmitted in one time interval. Hence,

$$E[X_n] = \lambda r / (\gamma - 1). \quad (400)$$

Also, another important relationship for this model, which is provided here without proof, is

$$\gamma = 3 - 2H, \quad (401)$$

where  $H$  is the Hurst parameter.

Having the last two equations, we are able to fit the overall mean of the process ( $E[X_n]$ ) and the Hurst parameter of the process with those measured in a real life process, and generate traffic based on the M/Pareto/ $\infty$  model.

## Homework 17.2

Use the 100,000 numbers representing the number of packets arriving recorded every second for consecutive 100,000 seconds you have collected in the assignments of Section 15 Using the UNIX command *netstat*. Again assume that these numbers represent the amount of work, measured in packets, which arrive at an SSQ during 100,000 consecutive time-intervals. Let  $E[A]$  be their average. Use your SSQ simulation of the assignments of Section 15, and compute the PLR, the correlation and the variance of the amount of work arrive in large intervals (each of 1000 packet-transmission times) for the various processes you have considered and discuss the differences.  $\square$

## Homework 17.3

Compare by simulations the effect of the correlation parameter  $a$  on the performance of the queues EAR(1)/EAR(1)/1 versus their EAR(1)/M/1, M/EAR(1)/1 and M/M/1 equivalence. Demonstrate the effect of  $a$  and  $\rho$  on mean delay. Use the ranges  $0 \leq a \leq 1$  and  $0 \leq \rho \leq 1$ .

$\square$

## The End of the Beginning

It is appropriate now to recall Winston Churchill's famous quote: "Now this is not the end. It is not even the beginning of the end. But it is, perhaps, the end of the beginning." In this book, the reader has been introduced to certain fundamental theories, techniques and numerical methods of queueing theory and related stochastic models as well as to certain practical telecommunications applications. However, for someone who is interested to pursue a research career in this field, there is a scope for far deeper and broader study of both the theory of queues as well as the telecommunications applications. For the last half a century, advances in telecommunications technologies have provided queueing theorists with a wealth of interesting problems and research challenges and it is often said that the telecommunications and information technologies actually revived queueing theory. However, this is only part of the story. The fact is that many exceptional queueing theorists also became experts in the technologies and have made tremendous contributions to the design, operations and understanding of telecommunications systems and networks. This dual relationship between the two fields will likely to continue, so it is very much encouraged to develop expertise in both fields. And if the aim is to become expert in both, it is not the end of the beginning, but merely the beginning.

## References

- [1] R. G. Addie, T. Neame and M. Zukerman, "Performance evaluation of a queue fed by a Poisson Pareto Burst process", *Computer Networks*, vol. 40, no. 3, October 2002, pp. 377–397.
- [2] R. G. Addie, M. Zukerman and T. D. Neame, "Broadband traffic modeling: simple solutions to hard problems", *IEEE Communication Magazine*, August 1998, pp. 88–95.
- [3] R. G. Addie and M. Zukerman, "An approximation for performance evaluation of stationary single server queues," *IEEE Transactions on Communications*, vol. 42, no. 12, December 1994, pp. 3150–3160.
- [4] R. G. Addie and M. Zukerman, "Performance evaluation of a single server autoregressive queue," *Australian Telecommunication Research*, vol. 28, no. 1, 1994, pp. 25–32.
- [5] R. G. Addie and M. Zukerman, "Queues with total recall - application to the B-ISDN," in J. Labetoulle and J. W. Roberts Editors, *The Fundamental Role of Teletraffic in the Evolution of Telecommunications Networks, Vol. 1a in the series: Teletraffic Science and Engineering*, North Holland - Elsevier Science Publishers, Amsterdam 1994, pp. 45–54.
- [6] H. Akimaru and K. Kawashima, *Teletraffic - theory and application*, ISBN 3-540-19805-9, 2nd ed., Springer, London 1999.
- [7] O. D. Anderson, *Time Series analysis and forecasting - The Box Jenkins approach*, Butterworths, London, 1976.
- [8] V. E. Benes, *General Stochastic Processes in the Theory of Queues*, Addison Wesley, 1963.
- [9] D. Bertsekas and R. Gallager, *Data networks*, Prentice Hall, Englewood Cliff, New Jersey 1992.
- [10] D. Bertsekas and J. N. Tsitsiklis, *Introduction to Probability*, Athena Scientific, Belmont, Massachusetts 2002.
- [11] S. K. Bose, *An Introduction to queueing systems*, Kluwer Academic/Plenum Publisher, New York 2002.
- [12] G. E. P. Box and G. M. Jenkins, *Time series analysis forecasting and control*, Holden Day San Francisco, 1970.
- [13] P. J. Burke, "The output of a queueing system," *Operations Research*, vol. 4, (1956), pp. 699–704.
- [14] R. L. Cruz, "A calculus for network delay, part I: Network elements in isolation," *IEEE Trans. on Information Theory*, vol. 37, no. 1, January 1991, pp. 114–131.
- [15] R. L. Cruz, "A calculus for network delay, part II: Network analysis," *IEEE Trans. on Information Theory*, vol. 37, no. 1, January 1991, pp. 132–141.
- [16] J. W. Cohen, *The Single Server Queue*, North-Holland Publ. Cy., rev. ed., Amsterdam, 1982.

- [17] R. B. Cooper, *Introduction to Queueing Theory*, North Holland 1981.
- [18] T. Engset, "Die wahrscheinlichkeitsrechnung zur bestimmung der wähleranzahl in automatischen fernsprechamtern." *Elektrotechnische zeitschrift*, vol. 39, no. 31, Aug. 1918, pp. 304–306.
- [19] A. K. Erlang, "Solution of some problems in the theory of probabilities of significance in automatic telephone exchanges", *The Post Office Engineer's Journal*, vol. 10, 1917, pp. 189–197.
- [20] W. Feller, *An introduction to probability theory and its applications*, Wiley, New York, 1968.
- [21] W. Fischer and K. Meier-Hellstern, "The Markov Modulated Poisson process (MMPP) cookbook," *Performance Evaluation*, vol. 18, pp. 149–171, 1992.
- [22] G. S. Fishman, "Discrete-event simulation," Springer-Verlag, London 2001.
- [23] D. P. Gaver and P. A. W. Lewis, "First-Order Autoregressive Gamma Sequences and Point Processes," *Advances in Applied Probability*, vol. 12, no. 3, pp. 727–745, Sep. 1980.
- [24] D. Gross and C. M. Harris, *Fundamentals of queueing theory*, Wiley Interscience, 1974.
- [25] H. Heffes and D. Lucantoni, "A Markov modulated characterization of packetized voice and data traffic and related statistical multiplexer performance," *IEEE Journal on Selected Areas in Communications*, vol. 4, no. 6, pp. 856–868, Sep 1986.
- [26] P. G. Hoel, S. C. Port and C. J. Stone, "Introduction to Stochastic Processes," Houghton Mifflin, 1972.
- [27] F. Hübner, "Analysis of a finite-capacity asynchronous multiplexer with deterministic traffic sources", *proceedings of 7th ITC Specialist Seminar*, Moristoen, October 1990.
- [28] J. Hui, *Switching and traffic theory for integrated broadband networks*, Kluwer 1990.
- [29] V. B. Iversen, *Teletraffic engineering and network planning*, Technical University of Denmark, Revised January 2004, available:  
<http://www.com.dtu.dk/education/34340/material/telenook.pdf>
- [30] D. L. Jagerman, "Some properties of the Erlang loss function," *Bell Syst. Tech. J.*, vol. 53, pp. 525–511, 1974.
- [31] F. P. Kelly, *Reversibility and stochastic networks*, John Wiley, 1979.
- [32] D. G. Kendall, "Stochastic processes occurring in the theory of queues and their analysis by the method of the imbedded Markov chain", *The Annals of Mathematical Statistics*, vol. 24, no. 3, Sep., 1953, pp. 338–354.
- [33] Kendall D. G., "Some problems in the theory of queues", *J. R. Statist. Soc. B*, vol. 13, pp. 151–173, 1951.
- [34] L. Kleinrock, *Queueing Systems, Volume 1 - Theory*, John Wiley and Sons 1975.

- [35] L. Kleinrock, *Queueing Systems, Volume 2 - Computer Applications*, John Wiley and Sons 1976.
- [36] A. Leon Garcia, *Probability and random processes for electrical engineering*, Addison Wesley, Reading, Mass., 1989.
- [37] G. Latouche and V. Ramaswami, *Introduction to matrix analytic methods in stochastic modeling*, SIAM, 1999.
- [38] H. J.-Y., Le Boudec and P. Thiran, *Network calculus*, available at URL: [http://lrcwww.epfl.ch/PS\\_files/NetCal.htm](http://lrcwww.epfl.ch/PS_files/NetCal.htm)
- [39] D. V. Lindley, "The theory of queues with a single server", *Proc. Camb. Phil. Soc.*, vol. 48, pp. 277-289, 1952.
- [40] J. D. C. Little, "A proof of the queueing formula:  $L = \lambda W$ ", *Operations Research*, vol. 9, pp. 383-387, 1961.
- [41] B. Maglaris, D. Anastassiou, P. Sen, G. Karlsson, and J. D. Robbins, "Performance models of statistical multiplexing in packet video communications", *IEEE Transactions on Communications*, vol. 36, no. 7, pp. 834-844, July 1988.
- [42] J. Mahdavi and S. Floyd, "Tcp-friendly unicast rate-based flow control", January 1997, technical note sent to the end2endinterest mailing list.
- [43] USC/ISI, Los Angeles, "The ns simulator and the documentation", <http://www.isi.edu/nsnam/ns/>.
- [44] C. R. Nelson, *Applied time series analysis*, Holden Day, San Francisco, 1973.
- [45] M. F. Neuts, *Structured stochastic matrices of M/G/1 type and their applications*, Marcel Dekker, NY 1989.
- [46] M. F. Neuts, *Matrix-geometric solutions in stochastic models - algorithmic approach*, The John Hopkin's University Press, Baltimore (1981).
- [47] A. Papoulis, *Probability, random variables and stochastic processes*, 2nd Edition, McGraw Hill Book Company, New York, 1984.
- [48] E. Parzen, *Stochastic processes*, Holden Day, San Francisco, 1962.
- [49] E. Reich, "On the integrodifferential equation of takacs. 1." *Ann. Math. Stat.*, vol. 29, pp. 563-570, 1958.
- [50] J. Roberts, U. Mocci and J. Virtamo, *Broadband Network Teletraffic*, Final Report of Action COST 242, Lecture Notes in Computer Science, vol. 1155, Springer, 1996.
- [51] Z. Rosberg, H. L. Vu, M. Zukerman and J. White, "Performance analyses of optical burst switching networks", *IEEE journal of selected areas in communications*, vol. 21, no. 7, Sept. 2003, pp. 1187-1197.
- [52] K. W. Ross, *Multiservice Loss Models for Broadband Telecommunication Networks*, ISBN 3-540-19918-7 Springer-Verlag Berlin Heidelberg New York, 1995.

- [53] S. M. Ross, *Introduction to probability models*, Academic Press, New York, 1993.
- [54] S. M. Ross, *A first course in probability*, Macmillan, New York, 1976.
- [55] S. M. Ross, *Stochastic processes*, Wiley, New York, 1996.
- [56] S. Stidham, “ $L = \lambda W$ : a discounted analogue and a new proof”, *Operations Research* vol. 20, pp. 1115–1126, 1972.
- [57] S. Stidham, “A last word on  $L = \lambda W$ ”, *Operations Research*, vol. 22, pp. 417–421, 1974.
- [58] J. S. Turner, “New directions in communications (or which way to the information age?)”, *IEEE Communications Magazine*, vol. 24, no. 10, October 1986, pp. 8–15.
- [59] A. K. Wong, “Queueing analysis for ATM switching of continuous-bit-rate traffic – a recursion computation method,” *Proceedings of GLOBECOM '90*, vol. 3, San Diego, CA, Dec. 2-5, 1990, pp. 1438–1444.
- [60] U. Yechiali and P. Naor, “Queueing Problems with Heterogeneous Arrivals and Service,” *Operations Research*, vol. 19, pp. 722–734, 1971.
- [61] M. Zukerman and I. Rubin, “Performance of flow-controlled communications systems under bursty traffic,” *Proceedings of IEEE GLOBECOM '86*, vol. 3, Houston, December 1986, pp. 1266–1271.
- [62] M. Zukerman and I. Rubin, “Queueing performance of demand-assigned multi access communication systems under bursty traffic conditions,” *Proceedings of IEEE ICC '86*, vol. 3, no. 57.2, Toronto, Canada, June 1986, pp. 1827–1832.
- [63] M. Zukerman and I. Rubin, “On multi channel queueing systems with fluctuating parameters,” *Proceedings of IEEE INFOCOM '86*, Miami, Florida, April 1986, pp. 600–608.

## Exam/Midtest Formula Sheet

$$P(X = i) = e^{-\lambda} \frac{\lambda^i}{i!} \quad i = 0, 1, 2, 3, \dots \quad f(x) = \begin{cases} \mu e^{-\mu x} & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

$$F(x) = \int_0^x \mu e^{-\mu s} ds = 1 - e^{-\mu x} \quad x \geq 0. \quad \bar{F}(x) = e^{-\mu x} \quad x \geq 0.$$

$$E[Q] = \lambda E[D] \quad p_0 = 1 - U = 1 - \lambda/\mu \quad E[Q] = \frac{\rho}{1-\rho} \quad E[D] = \frac{1}{\mu(1-\rho)} = \frac{1}{\mu-\lambda}$$

$$\delta_D(x) = \begin{cases} (\mu - \lambda)e^{(\lambda-\mu)x} & \text{if } x \geq 0 \\ 0 & \text{otherwise.} \end{cases} \quad f_X(x) = \frac{\lambda^k x^{k-1} e^{-\lambda x}}{(k-1)!}.$$

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m)^2/2\sigma^2} \quad -\infty < x < \infty. \quad P(X > x) = \begin{cases} \left(\frac{x}{\delta}\right)^{-\gamma}, & x \geq \delta \\ 1, & \text{otherwise.} \end{cases}$$

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx. \quad E[Z] = \int_0^{\infty} P(Z > z) dz = \int_0^{\infty} [1 - F_Z(z)] dz.$$

$$E[Z] = \sum_{n=0}^{\infty} P(Z > n) = \sum_{n=0}^{\infty} [1 - F_Z(n)]. \quad \text{var}[X] = \sum_{\{k: P(k) > 0\}} (k - E[X])^2 P_X(k)$$

$$\text{var}[X] = \int_{-\infty}^{\infty} (x - E[X])^2 f_X(x) dx \quad P(X > C) \leq \frac{E[X]}{C}. \quad P(|X - E[X]| > C) \leq \frac{\text{var}[X]}{C^2}.$$

| random variable | parameters                        | mean                              | variance     |
|-----------------|-----------------------------------|-----------------------------------|--------------|
| Bernoulli       | $0 \leq p \leq 1$                 | $p$                               | $p(1-p)$     |
| binomial        | $n$ and $0 \leq p \leq 1$         | $np$                              | $np(1-p)$    |
| Poisson         | $\lambda > 0$                     | $\lambda$                         | $\lambda$    |
| uniform         | $a$ and $b$                       | $(a+b)/2$                         | $(b-a)^2/12$ |
| exponential     | $\mu > 0$                         | $1/\mu$                           | $1/\mu^2$    |
| Gaussian        | $m$ and $\sigma$                  | $m$                               | $\sigma^2$   |
| Pareto          | $\delta > 0$ and $1 < \gamma < 2$ | $\frac{\delta\gamma}{(\gamma-1)}$ | $\infty$     |

$$E_k(A) = \frac{\frac{A^k}{k!}}{\sum_{n=0}^k \frac{A^n}{n!}} \quad \text{var}[X] = E[\text{var}[X | Y]] + \text{var}[E[X | Y]]$$

$$P_X(x) = E_Y[P(X = x | Y = y)] \quad I_n(A) = A \int_0^{\infty} e^{-Ay} (1+y)^n dy \quad U = \frac{(1-\pi_k)A}{k} \quad U = \frac{A}{k}$$

$$\pi_0 = \left( \sum_{n=0}^{k-1} \frac{A^n}{n!} + \frac{A^k}{k!} \frac{k}{k-A} \right)^{-1} \quad C_k(A) = \sum_{n=k}^{\infty} \pi_n = \frac{A^k}{k!} \frac{k}{k-A} \pi_0 = \frac{\frac{A^k}{k!} \frac{k}{k-A}}{\sum_{n=0}^{k-1} \frac{A^n}{n!} + \frac{A^k}{k!} \frac{k}{k-A}}$$

$$C_k(A) = \frac{k E_k(A)}{k-A[1-E_k(A)]} \quad \pi_k = \rho^k \frac{1-\rho}{1-\rho^{k+1}} = \frac{\rho^k - \rho^{k+1}}{1-\rho^{k+1}} = \frac{\rho^k(1-\rho)}{1-\rho^{k+1}}$$

$$E[Q] = \rho + \frac{\rho^2 + \lambda^2 \sigma_s^2}{2(1-\rho)} \quad W_q(t) = \sup_{0 \leq s < t} \{W_a(t) - W_a(s) - t + s\}$$

$$E[D(j)] = E[W_Q(j)] + \frac{1}{\mu_j} \text{ for } j = 1, 2, 3, \dots, m \quad E[D(j)] = \frac{(1/\mu_j)(1 - \sum_{i=1}^j \rho_i) + R(j)}{(1 - \sum_{i=1}^{j-1} \rho_i)(1 - \sum_{i=1}^j \rho_i)}$$

$$E[X_n] = \lambda r / (\gamma - 1) \quad \gamma = 3 - 2H \quad D_n = a D_{n-1} + I_n E_n$$

$$X_n = a_1 X_{n-1} + a_2 X_{n-2} + \dots + a_k X_{n-k} b \tilde{G}_n \quad B_R = 1 - \prod_{j \in R} (1 - B(j)) \quad \rho P_{\text{loss}} \leq P(Q > k)$$

$$\lambda_j = A_j + \sum_{i=1}^N \lambda_i P_{ij}, \quad j = 1, 2, 3, \dots, N$$

$$P_b = \frac{\binom{M-1}{k} \rho^k}{\sum_{i=0}^k \binom{M-1}{i} \rho^i}. \quad B_i = \frac{\rho^{(M-i)B_{i-1}}}{i + \rho^{(M-i)B_{i-1}}} \quad i = 1, 2, 3, \dots, k$$