# Machine Learning — Prof. Anderson Rocha.

① Grew out of work in ⒶI

② New capabilities for computers

## Some examples

① Database mining, large datasets, web automation, web click data, medical records, social networks

② Applications can't program by hand
- autonomous car, helicopter
- handwritten recognition
- NLP
- vision

③ Self-customizing
- Amazon
- Netflix product recommendations

④ understanding human learning.

---

## What is ⓂⓁ

① field of study that gives computers the ability to learn without being explicitly programmed. (Arthur Samuel)

② Tom Mitchel : well-posed learning program ⇒ a computer program is said to learn from exp Ⓔ wrt task Ⓣ and some perf. measure Ⓟ if its perf. on Ⓣ as measured by Ⓟ, improves with exp. Ⓔ.

checkers.

---

## Types of ⓂⓁ

- SL (vs) SSL (vs) unsup. learning
- Reinf learning, recommendn systems
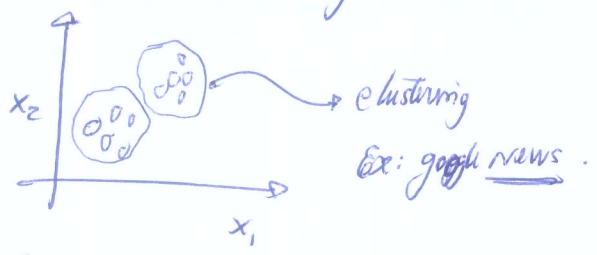
Practical advice : how to use and where to apply ⓂⓁ.

Supervised Learning

- Regression ⟹ predict cont. output.
- classif ⟹ predict discrete value (class) output.

unsupervised Learning    find some structure on the data

$x_2$
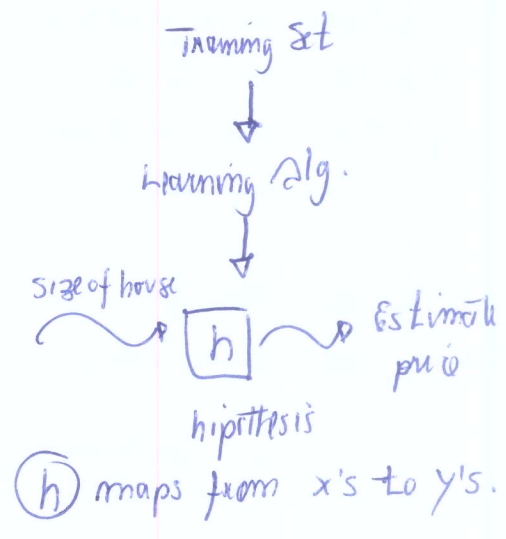
→ clustering

Ex: google news.

Example : - Genes x Individuals
- Datacenters.
- Cohesive groups of people on social networks
- Group customers
- Astronomical data analysis (how galaxies are formed).

Cocktail party problem ⟹ overlapping voices.

---

Languages   - Octave
     - matlab     }   first prototyping then ~→ other languages.
     - R

---
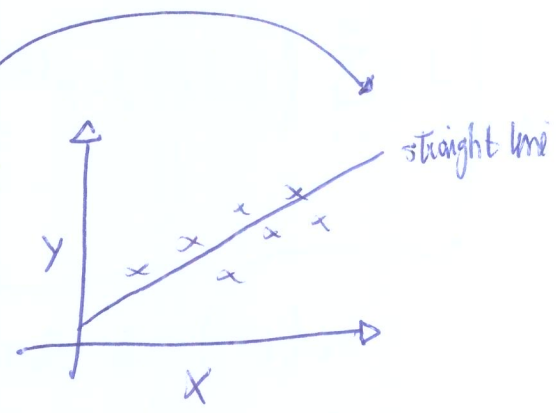
Notation    $m$ : nbr of training examples

$x$'s input variables/features

$y$'s out. variable/target variable

$(x, y)$ one training example.

$(x^{(i)}, y^{(i)})$ $i^{th}$ training example.

Training set
↓
Learning alg.
↓
size of house → $\boxed{h}$ → Estimate price

hipothesis

$h$ maps from x's to y's.

How to represent ⓗ ?

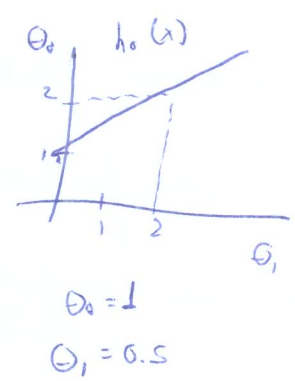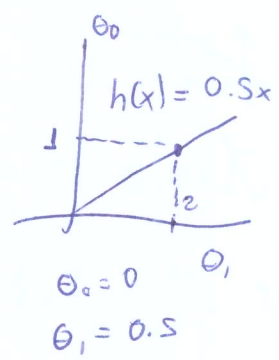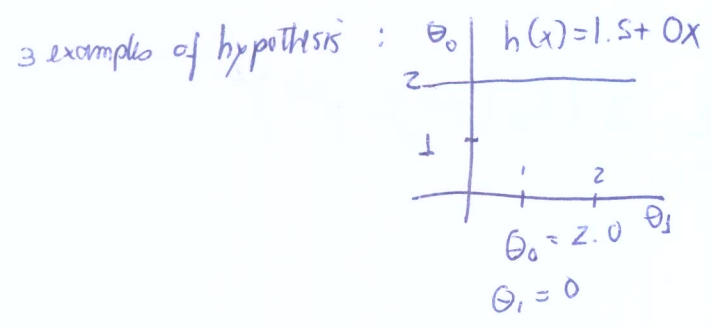Initial choice    $h_\theta(x) = \underline{\theta_0 + \theta_1 x}$

Shorthand    $h(x)$


straight line

why linear function ⇒ simple building block.

Here we have an example of a linear regression with one variable or

$\underbrace{univariate}_{\text{1 variable}}$ linear regression.

---

① Defining a cost function

hypothesis    $h_\theta(x) = \theta_0 + \theta_1 x$    $\theta_s$ parameters

3 examples of hypothesis :    $h(x) = 1.5 + 0x$      $h(x) = 0.5x$      $h_\theta(x)$



$\theta_0 = 2.0$      $\theta_0 = 0$      $\theta_0 = 1$
$\theta_1 = 0$       $\theta_1 = 0.5$     $\theta_1 = 0.5$

In linear Regression


$(\theta_0, \theta_1)$

$h_\theta(x)$ must be $\underline{\underline{close}}$ to the training data.

┌─────────┐
│ More    │
│ formally│
└─────────┘

minimize $\theta_0, \theta_1 \left( h_\theta(x) - y \right)^2$

$\min_{\theta_0, \theta_1} \dfrac{1}{2M} \sum_{i=1}^{M} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

$\dfrac{1}{2M}$  → # Examples
only for make it    in the training.
easier.

$$\min_{\theta_0, \theta_1} \left( \frac{1}{2M} \sum_{i=1}^{M} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right)^2 \right) = J(\theta_0, \theta_1).$$

$$h_\theta \left( x^{(i)} \right) = \theta_0 + \theta_1 x^{(i)}$$

with this, we define a cost function $J(\theta_0, \theta_1)$.

so we can simply say $\boxed{\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)}$ ----- cost function

$\boxed{SSD}$. called squared Error

cost function.

most used for Regression.

---

## Cost function intuition
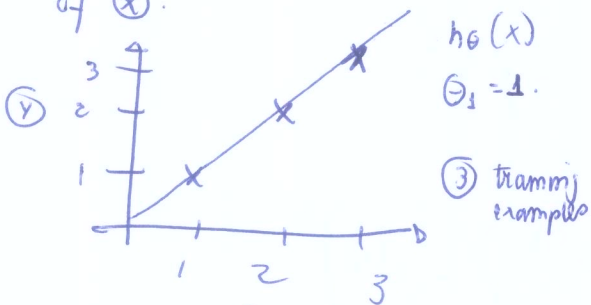
① hypothesis $h_\theta(x) = \theta_0 + \theta_1 \dot{x}$.

② Params $\theta_0, \theta_1$ ⟋ $h(x)$.

③ Cost function $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right)^2$
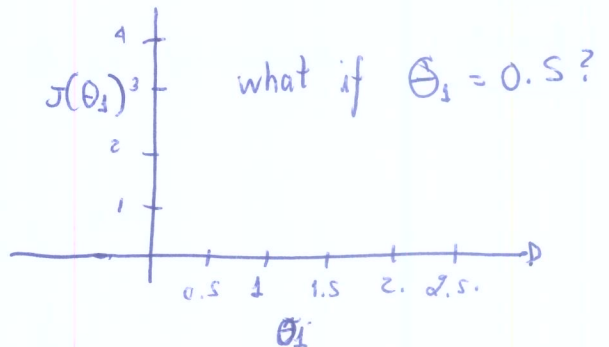
④ Goal: $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

If we simplify $h_\theta(x) = \theta_1 x$.  $\theta_0 = 0$.

$\boxed{h_\theta(x)}$ for fixed $\theta_1$, this is a function of $x$.



$h_\theta(x)$
$\theta_1 = 1$.

③ training example

what is $J(\theta_1)$ ? $\frac{1}{2m} \left[ \theta_1 (x^{(i)} - y^{(i)})^2 \right]$

$= (0^2 + 0^2 + 0^2) \frac{1}{2m} = \boxed{0}$

$\boxed{J(\theta_1)}$ is a function of the param $\theta_1$.



what if $\theta_1 = 0.5$?

$J(\theta_1)$

$\theta_1$

$h_\theta(x)$



$h_\theta(x)$

$\theta_1 = \frac{1}{2}$

line with slope $\frac{1}{2}$.

$y$    $h_\theta(0)$

$x$

$$J(0.5) = \frac{1}{2M}\left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2\right]$$

$$= \frac{1}{2\times3}(3.5) \approx \boxed{0.58}$$

$J(\theta)$



$J(\theta_1)$

$\theta_1$

$J(0)=?$

$$J(0) = \frac{1}{2m}(1^2 + 2^2 + 3^2)$$

$$\frac{1}{2m} \cdot 14 = \boxed{2.3}$$

for each value of $\theta_1$, we have a different hypothesis for $h_\theta(x)$ (different line).

---

Cost function (contour plots)



$\$$    hypothesis

size $M^2$

$h_\theta(x)$



min    same $J(\theta_0, \theta_1)$

$\theta_1$   0   0.15

360   800   $\theta_0$

| $\theta_0,$ | $\theta_1$ |
|---|---|
| 800 | 1.5 |

$\theta_0 = 0$
$\theta_1 = 360$

# Gradient Descent

Used everywhere in (ML)

① Have some function $J(\theta_0, \theta_1)$

② $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

③ outline

   (a) start with some $\theta_0, \theta_1$

   (b) keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$ until we end up at a minimum.

   ⤳ what we need to define is $\alpha$ (the learning rate).

| Implementation | Correct form | Incorrect form |
|---|---|---|

Correct form — partial derivatives:

$t_0 \leftarrow \theta_0 - \alpha \dfrac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$t_1 \leftarrow \theta_1 - \alpha \cdot \dfrac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

$\left. \begin{array}{l} \theta_0 \leftarrow t_0 \\ \theta_1 \leftarrow t_1 \end{array} \right\}$ simultaneously update both $\theta_0$ and $\theta_1$
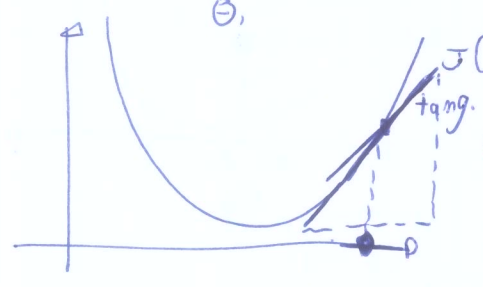
Incorrect form:

~~temp~~ $t_0 \leftarrow \theta_0 - \alpha \dfrac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$\theta_0 \leftarrow t_0$

$\vdots$

# Intuition about the Gradient Descent

Suppose we have $\min\limits_{\theta_1} J(\theta_1)$   $\theta_1 \in \mathbb{R}$.



$J(\theta_1)$   $\theta_1 \in \mathbb{R}$.

$\theta_1 = \theta_1 - \alpha \underbrace{\dfrac{d}{d\theta_1}}_{\text{slope}} J(\theta_1)$.

derivative $\underbrace{\dfrac{\partial}{\partial \theta}}$ partial derivative.

~~tangent~~ of the line tangent to function.
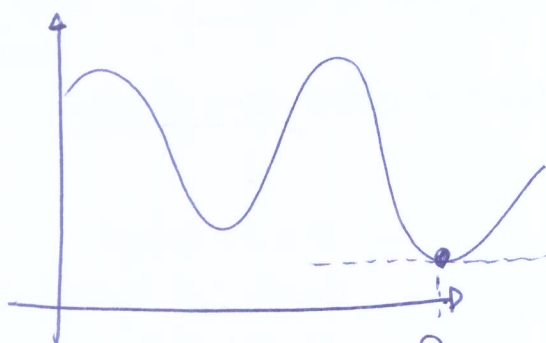
$\geqslant 0$.

$\theta_1 \leftarrow \theta_1 - \alpha \cdot (\text{positive number}) \downarrow \theta$.

$$\frac{\partial J(\theta_1)}{d\theta_1}$$

$$\leq 0.$$

$$\theta_1 \leftarrow \theta_1 - \alpha \left(\text{neg. number}\right) \uparrow \theta.$$

↝ **what if (GD) is already at a _minimum_ ?**



tangent line
slope = 0 (zero)

$\theta_1$

Current value of $\theta_1$

$$\theta_1 \Leftarrow \theta_1 - \alpha \left| \frac{d}{d\theta_1} J(\theta_1) \right|$$

zero

$$\theta_1 \leftarrow \theta_1 - \alpha \cdot 0$$

$$\theta_1 \leftarrow \theta_1.$$

$\boxed{\theta_1 \text{ stays unchanged}}$ that's why GD can converge to a local _min_ even with $\alpha$ fixed.

GD automatically takes smaller steps as we approach a local min due to the derivative term. So, no need for decreasing $\alpha$ over time.

---

**GD for linear regression**

how to derive $\dfrac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \dfrac{\partial}{\partial \theta_j} \cdot \dfrac{1}{2m} \cdot \sum\limits_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2.$

$$= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^{m} \left( \theta_0 + \theta_1 \cdot x^{(i)} - y^{(i)} \right)^2$$

$$\frac{\partial}{\partial x} \left( x^2 - z \right)^2$$

$$2 \left( x^2 - z \right) \cdot 2x$$

$$4x \left( x^2 - z \right)$$

$$\theta_0 \implies j = 0: \quad \frac{1}{m} \cdot \sum_{i=1}^{m} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right)$$

$$\theta_1 \implies j = 1: \quad \frac{1}{m} \cdot \sum_{i=1}^{m} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right) \cdot x^{(i)}$$

↝ • The cost function for ~~$\theta\theta$~~ linear regression is <u>always</u> a bow shape function

↝ • The technical term for this is that it is called <u>convex function</u>

                  ↓

                   — is a bow shape

                   — doesn't have local opt except for the global one.

The algorithm we just defined is what we call <u>Batch Gradient Descent</u>.

    <u>Batch</u> : each step of ⊙ⅅ uses all training examples.

↝ • there other version that use only ~~a few~~ subset of the training.

↝ • for the case of linear regression, there is a <u>closed form solution</u> (normal equations) but Gradient descent will scale better for <u>large</u> datasets than normal equations.

↝ • The derivative is just the <u>slope</u> of the cost function ⊙.

# Generalization of Gradient Descent

Recalling, the GD algorithm is given by

repeat until convergence {

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right)$$

$$\theta_1 \leftarrow \theta_1 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^{m} \left( h_\theta \left( x^{(i)} \right) - y^{(i)} \right) - x^{(i)}}_{\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1).}$$

}

→ It can be susceptible to <u>local</u> optima.
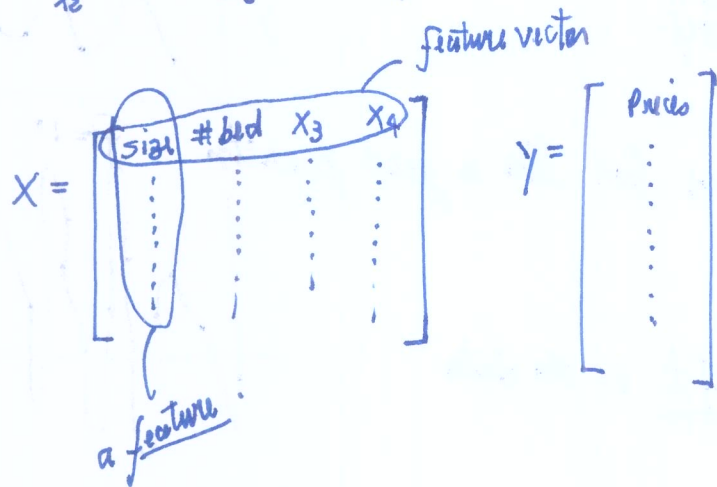
→ But the cost function for the linear regression is always <u>convex</u> (bow-shape)

→ only the global optimum.

---- × ----

There are two extensions of GD that we must consider

① In min $J(\theta_0, \theta_1)$, solve for $\theta_0, \theta_1$ exactly, without iterative GD algorithm
One advantage is that $\alpha$ is not necessary anymore but in some cases (~~large~~ high dim) the GD may be more appropriate.

② Longer ~~and~~ number of features. For instance

size, # bedrooms, # floors, # age of home ⟶ y (price).

$x_1$     $x_2$     $x_3$     $x_4$

Notation

feature vector

$$X = \begin{bmatrix} \text{size} & \text{\# bed} & X_3 & X_4 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix} \qquad Y = \begin{bmatrix} \text{price} \\ \vdots \end{bmatrix}$$

a feature.

# Increasing the number of features

we have considered so far linear functions of several observed ~~variables~~ features

$$x_1, x_2, \ldots$$

Suppose now that we have only one feature $\boxed{x_1}$ but we would like our predicter to be a non linear function of $\otimes$ $\hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \ldots$

We can simply define new features $x_2 = x^2$, $x_3 = x^3$ just like we did with
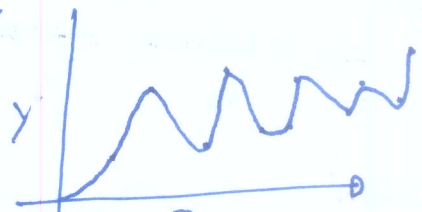
$$x_0 = 1 = x^0$$

⟶ The predicter then becomes $\hat{y}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots$

<u>It is</u> : it still fits the linear regression model but in a <u>new</u> <u>feature space</u> with additional features that are deterministic functions of our observations
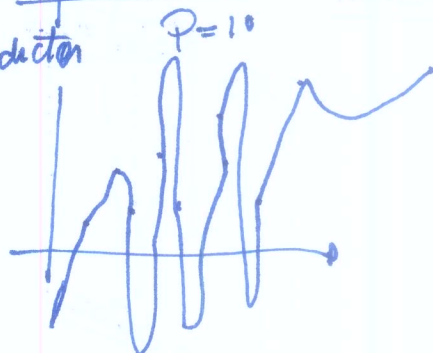
Applying the mse estimator $\hat{y}(x) = \sum_{i=1}^{10} \theta_i x_i^i$

---

## Overfitting

In our polynomial fits for 11 data points, a more complex model $\hat{y}(x)$ a $p(x)$ with degree 10 fits the data perfectly
and 11
coef.
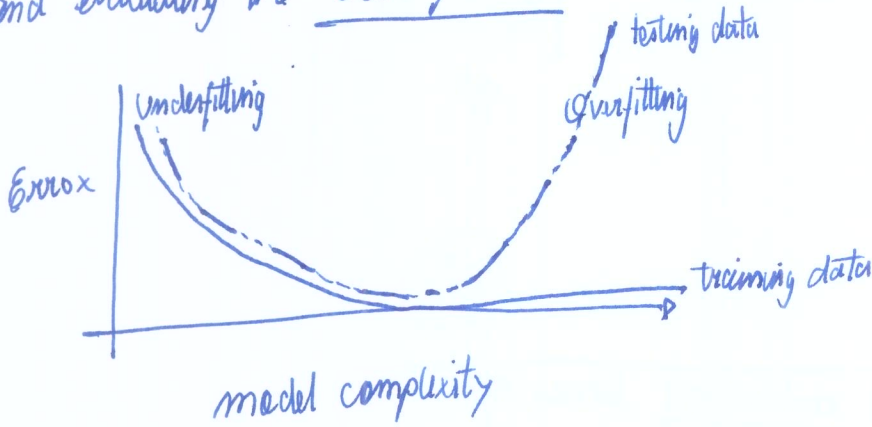
but it does not look like a good predicter

$P = 10$

we fact an <u>overfit</u> to the data

We can see the $\boxed{\text{generalization}}$ or test error simply by gathering more data and evaluating the __cost function__



# Continuing Regression with multiple variables

$$\overset{x_1}{\text{size house}}, \overset{x_2}{\#\text{bedrooms}}, \overset{x_3}{\text{floors}}, \overset{x_4}{\text{years}}, \overset{y}{\text{Price}}$$

$x^{(i)}$ training example.

$x_j^{(i)}$ value of feature $(j)$ in $(i\text{th})$ training example.

$\leadsto$ Previously we had a model / hypothesis: $\hat{y}(x) = \theta_0 + \theta_1 x$

$\leadsto$ Let's update or upgrade our model

$$\hat{y}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots \theta_4 x_4.$$

__Example__  $\hat{y}(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 \boxed{2 x_4}$

each year decreases the price.

for convenience let's define $x_0 = 1$.

zeroth feature always with the value ①.

Then $X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$   $\theta = \Theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$

0-index

0-index vector

$h_\theta(x) = \hat{y}(x)$

$\hat{y}(x) = \theta_0 \overset{1}{\cancel{x_0}} + \theta_1 x_1 \dots \theta_n x_n$

$= \theta^T x$

$\begin{bmatrix} \theta_0 & \dots & \theta_n \end{bmatrix} \quad \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix}$

$\theta^T \qquad\qquad x.$

$\boxed{\theta^T x} .$

$\boxed{\text{This is called } \underline{\text{multivariate}} \text{ linear regression}}$

multiple variables/features

---
x
---

How do we do the ⓖⓓ with multiple variables/features?

Hypothesis   $\hat{y}(x) = \theta^T x = \theta_0 \overset{1}{\cancel{x_0}} + \theta_1 x_1 \dots \theta_n x_n.$

Params   $\boxed{\theta_0, \dots, \theta_n}$ . let's think it as a vector $\overset{\text{vector}}{\theta} \in \mathbb{R}^{n+1}$

Cost function   $J(\theta_0, \dots, \theta_n) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} \left( \hat{y}(x^{(i)}) - y^{(i)} \right)^2.$

$J(\theta).$

Gradient Descent   Repeat

$\theta_j \leftarrow \theta_j - \alpha \dfrac{\partial}{\partial \theta_j} J(\vec{\theta}).$

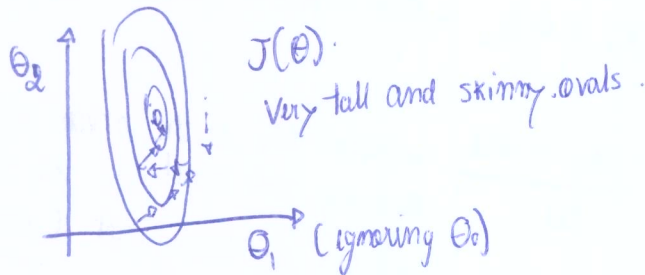update at the same time for all $j \in \{1, \dots, n\}$.

How to perform the partial derivative?

GD for $n \geq 1$. Repeat {

$\theta_j \leftarrow \theta_j - \alpha \cdot \dfrac{1}{m} \sum\limits_{i=1}^{m} \left( \hat{y}(x^{(i)} - y^{(i)} \right) \boxed{x_j^{(i)}}$   only thing that changes

$x_0^{(i)} = 1$ by definition

}

# Gradient Descent in practice 1: feature scaling

↪ what do we do if we have some features dominating the others?
↪ we need to scale the features in order $GD$ can converge more quickly.
we need to make sure the features are in the same range of values

Example $\quad x_1 = $ Size $(0 - 2000 m^2)$
$x_2 = $ nbr of bedrooms $(1-5)$



$J(\theta)$.
Very tall and skinny ovals.

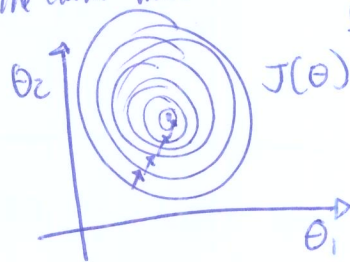$\theta_1$ (ignoring $\theta_0$)

if we run $GD$ in here, it may take a while before
converging

A nice workaround here would be to normalize the data then the contours would look like more circles

$$x_1' = \frac{\text{size } (m^2)}{2000} \qquad x_2' = \frac{\# \text{bedrooms}}{5}$$



$J(\theta)$

$\begin{cases} 0 \le x_1 \le 1 \\ 0 \le x_2 \le 1 \end{cases}$

Feature scaling ↪ Get every feature into approximately a $\underbrace{-1 \le x_i \le 1}$ range
desired but not firmly fixed

↪ $x_0$ is already in the range $(x_0 = 1)$

↪ $0 \le x_1 \le 3$ (is OK)
$-2 \le x_2 \le 4$ (is OK)
but
$-20 \le x_3 \le 40$ ✗
$-0.0001 \le x_4 \le 0.0001$ ✗.

↪ How to do it? training data.

In addition to <u>scale normalization</u>, sometimes we perform a step further and do <u>mean normalization</u>

↳ Replace $x_i$ with $x_i - \mu_i$ for zero meaning the features
↳ Do not change $x_0 = 1$

<u>Example</u> $\quad x_1 = \dfrac{size - 1000}{2000}$ $\qquad x_2 = \dfrac{\# bedrooms - 2}{5}$ $\quad$ if we do that, we will end up with values close enough to

$$-0.5 < x_1 < 0.5, \quad -0.5 < x_2 < 0.5$$

<u>General rule</u> $\quad x_i \leftarrow \dfrac{x_i - \mu_i}{S_i}$ where $\mu_i$ is the mean of $x_i$ in training set and $S_i$ is the range of values in training set $(max_i - min_i)$.
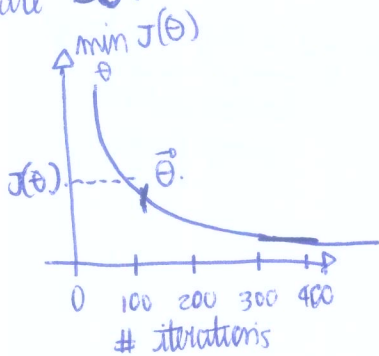
sometimes it will interesting to use the <u>std</u> of feature

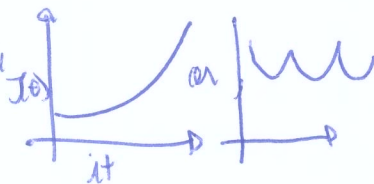① $x_i \leftarrow \dfrac{x_i - \mu_i}{\sigma_i}$ $\Big\}$ z-norm / z-score $\boxed{\text{Data standardization}}$

— ✳ —

<u>How do we know GD is working correctly and how to choose the learning rate $\alpha$.</u>


min $J(\theta)$ / $J(\theta)$ ... / $\hat{\theta}$ / what you expect.
0  100  200  300  400
# iterations
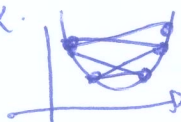
$J(\theta)$ should decrease after <u>iteration</u>
# of iterations vary according to the application

An automatic convergence would be possible, for instance if $J(\theta)$ in one iteration is only slightly lower ($10^{-3}$ for instance) than $J(\theta)$ in the previous iteration. But this is <u>difficult</u>.

if you have something like  or  then choose a smaller $\alpha$.

 Reason of increase: big $\alpha$ overshoots the minimum.

Two general observations:

① for small enough $\alpha$, $J(\theta)$ should decrease on every iteration...

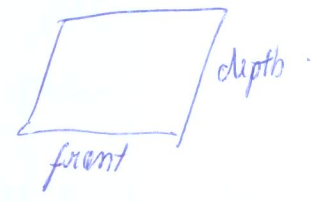② for small enough $\alpha$ or too small $\alpha$, GD can be slow to converge.

But... how to choose $\alpha$?

try... 0.001, 0.01, 0.1, 1, ... always plotting $J(\theta) \propto$ # iterations

or

0.001, 3×0.001, 0.01, 3×0.01, ···, ···

or

make it linked to the nbr of iteration $\frac{\alpha}{\text{iteration}}$.
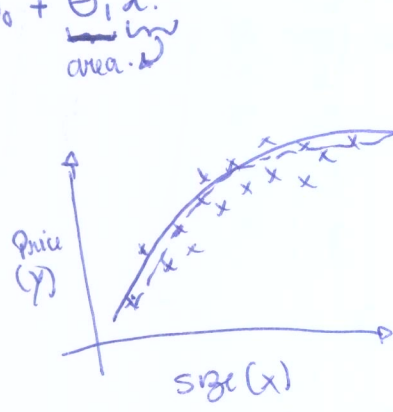
---

## Features and Polynomial Regression

Housing prices $\hat{y}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$

we can define a new feature $x = \text{front} \times \text{depth}$ (area)

then $\hat{y}(x) = \theta_0 + \underbrace{\theta_1 \underset{\text{in}}{x.}}_{\text{area.}}$

## Polynomial Regression



$$\hat{y_2}(x) = \theta_0 + \theta_1 x + \theta_2 x^2.$$

or

$$\hat{y_3}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

for using our multivariate linear regression model, we just define new features

$$\hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$\Downarrow$$

$$\hat{y}(x) = \theta_0 + \theta_1 x + \theta_2 x + \theta_3 x \quad \text{where}$$

$$x_1 = \text{size}, \quad x_2 = \text{size}^2 \quad \text{and} \quad x_3 = \text{size}^3.$$

[Note] : if we do this, feature scaling is [paramount]

$$\begin{cases} \text{size} : 1 - 10^3 \\ \text{size}^2 : 1 - 10^6 \\ \text{size}^3 : 1 - 10^9 \end{cases}$$

→ How to choose the features ? There are some algorithms for <u>that</u>.

## The Normal Equations

Gradient Descent



The <u>Normal Equation</u> is a method for solving for $\theta$ <u>analytically</u> (one step).

$$\theta \in \mathbb{R}^{n+1} \qquad J(\theta_0 \dots \theta_m) = \frac{1}{2m} \sum_{i=1}^{m} \hat{y}(x^{(i)}) - y^{(i)})^2$$

<u>minimize</u>

as we want to minimize $J(\theta)$, we solve for

$$\frac{\partial J(\theta)}{\partial \theta_j} = \dots \text{ set to } 0 \quad (\text{derive and equal to zero}) \\ \text{for every } ① .$$

solving for $\theta_0, \theta_1, \dots, \theta_n$.

<u>For example</u> : for a given problem, put the feature vectors in a matrix ⓧ and the outcomes (target) in a vector ⓨ .

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1400 & 2 & 3 & 40 \\ 1 & a & b & c & d \\ 1 & e & f & g & h \end{bmatrix}$$

area   bedrooms   floors   age

$$\mathbb{R}^{M \times (n+1)}$$

$$y = \begin{bmatrix} 460 \\ 230 \\ 500 \\ 200 \end{bmatrix}$$

price (K$)

$$\mathbb{R}^{m}$$

$$\hat{\theta} = (X^T X)^{-1} X^T y \quad \text{gives the } \underline{\text{min}} .$$

Ⓜ examples $\left(x^{(1)}, y^{(1)}\right), \ldots, \left(x^{(m)}, y^{(m)}\right)$, Ⓝ features

$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$  ∿➛ we will design what we call a design matrix  $X = \begin{bmatrix} -(x^{(1)})^T- \\ -(x^{(2)})^T- \\ \vdots \\ -(x^{(m)})^T- \end{bmatrix}$

$\mathbb{R}^{M \times (n+1)}$

what does it mean $\vec{\theta} = \left(X^T X\right)^{-1} X^T y$.

① $\left(X^T X\right)^{-1}$ is inverse of matrix $X^T X$.

∿ when using Normal Equations, $\boxed{\text{feature scaling}}$ is not necessary!

So, what should we use, Ⓖⓓ or Normal Equations?

|  Ⓖⓓ  | Normal Equation |
|---|---|
| ① Need to choose ⓐ. ↓ | ① No ⍺. ↑ |
| ② Needs many iterations ↓ | ② No iterations ↑ |
| ③ works well with a high ↑ number of features Ⓜ | ③ Needs to compute ↓ $\left(X^T X\right)^{-1}$ |
|  | ④ slow if Ⓜ is very large. |

** $X^T X \in \mathbb{R}^{n \times n}$, therefore normally it costs $O(n^3)$ for inverting on most inverse implementations

$\boxed{\text{what is large?}}$ $\begin{cases} \text{Ⓜ in hundreds, go with Normal Equation} \\ \text{Ⓜ in thousands } (< 5k), \text{ OK with Normal Eq.} \\ \text{Ⓜ} > 10k \text{ definitely go with Ⓖⓓ or some alternatives.} \end{cases}$

finally the normal equation will not work for more complex problems such as classification, that's why Ⓖⓓ is important and should be always thought as one option.

## Normal Equations and non invertibility.

~> when computing $\vec{\Theta} = \left(x^T x\right)^{-1} x^T y$, what if $\left(x^T x\right)^{-1}$ is non-invertible/singular/degenerate?

~> R and Octave and Matlab have workarounds (robust) inverse functions called pseudo-inve[rse]
that does the right thing.

~> normally there are two main causes for degeneration:

    ① Redundant features (linearly dependent)
        ↳ solution: _dim. reduction
             - feature deletion/selection

        Example  $x_1$ = size in $m^2$.      1 Km = 1000 meters
                $x_2$ = size in $Km^2$.

           So  $x_1 = \left(1000\right)^2 \cdot x_2$

    ② Too many features (# examples ⓜ ≪ # features ⓝ )
        ↳ solution: - dim reduction
             - dim/feature selection/deletion
             - Regularization

        Example $\Big\{$ $m$ = 10 examples
                    $n$ = 100 dimensions  $\overset{params}{\Theta} \in \mathbb{R}^{101}$

———————————— $x$ ————————————

## Organization of the class on August 14th, 2013

slide 9, 11, 12, 13, 14, 15, 16, 17, 18, ⑩.

# Logistic Regression and Classification

① Our outcome now is discrete-valued.

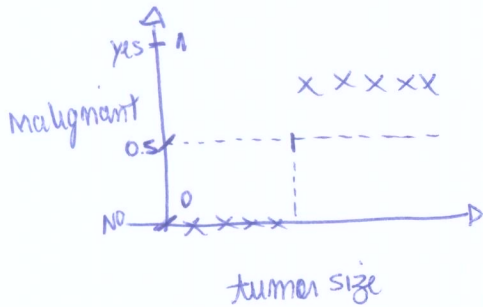② (LR) (LogR) is one of the most used learning algorithms.

Some classification problems:   spam × non spam (e-mail)
     fraudulent × non-fraudulent (transactions)
     tumor malign × benign
     going blind × not going blind (DR retinal).
         diab. retinopathy analysis

$$y \in \{0, 1\}$$

↱ positive class.

↳ neg. class.
      absence of something

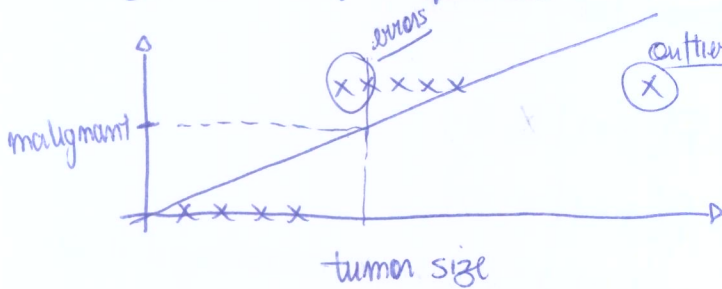binary classif. problem.



what if we approach the problem as a linear regression one?
$$\hat{y}(x) = \theta^T x.$$



Then we could threshold classifier's output at $\hat{y}(x)$ @ 0.5:

$$\begin{cases} \text{if } \hat{y}(x) \geqslant \tfrac{1}{2} \text{ predict } ① \\ \text{otherwise predict } ⓪. \end{cases}$$

But let's change the problem:



outlier.  Often linear regression isn't a good idea for classification.

↝ Another problem with (LR) (LogR) linear regression is that even with training examples with $y \in \{0,1\}$ it can output $\hat{y}(x) > 1$ or $\hat{y}(x) < 0$.

(LogR) solves this problem such that $0 \leqslant \hat{y}(x) \leqslant 1$ always.
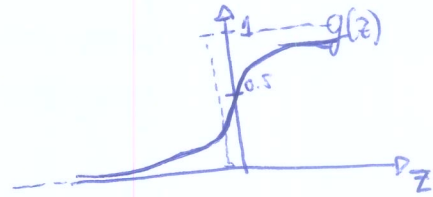
so we will use it as a classif. algorithm. The name is historical.

# Hypothesis Representation

Log R model  want  $0 \leq \hat{y}(x) \leq 1$

when we had (LR), $\hat{y}(x) = \overset{\circ}{\Theta}^T X$ , for (LogR), we change it to $\hat{y}(x) = g(\overset{\circ}{\Theta}^T x)$,

linear
regression

where $g(z) = \dfrac{1}{1+ e^{-z}}$

sigmoid or logistic function

synonyms.



By doing that, my new predictor will be

$$\hat{y}(x) = \dfrac{1}{1+ e^{-(\Theta^T x)}}$$

Now we need to (fit) parameters for $\Theta$.

## Interpretation of hypothesis Output

$\hat{y}(x) = $ estimated prob that $y=1$ on input $x$.

Example  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumon size} \end{bmatrix}$

Suppose my outcome is
$y(x) = 0.7$

Tell patient that 70% chance of tumor being malignant.

Mathmatically, we have $\hat{y}(x) = P(y=1 | x; \Theta)$
$= P(y=1 | x; \Theta)$.

prob. of $y=1$ given $x$, parameterized by $\Theta$.

Properties $\begin{cases} P(y=0 | x; \Theta) + P(y=1 | x; \Theta) = 1. \\ P(y=0 | x; \Theta) = 1 - P(y=1 | x; \Theta). \end{cases}$

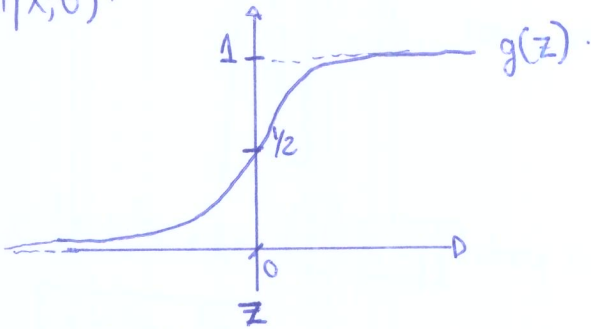*$Y$ must take values in $\{0, 1\}$ binary value only.

# Decision Boundary

$P(y=1|x;\theta)$

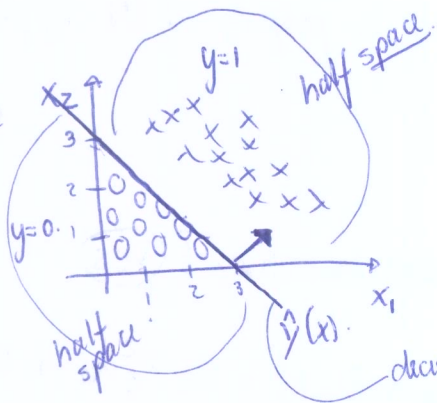$\boxed{\text{Log R}}$  $\hat{y}(x) = g(\vec{\theta}^T x)$

$g(z) = \dfrac{1}{1+e^{-z}}$



Suppose predict $y=1$ if $\hat{y}(x) \geqslant 0.5$

$y=0$, otherwise.

if we look at the sigmoid plot, we see $g(z) \geqslant 0.5$ when $z \geqslant 0$.

then given our prediction $\hat{y}(x) = g(\vec{\theta}^T x)$, it will be $\geqslant 0.5$ when $\underbrace{\vec{\theta}^T x}_{z} \geqslant 0$.

Suppose we have training set



half space

$y=1$

decision boundary

$\hat{y}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$.

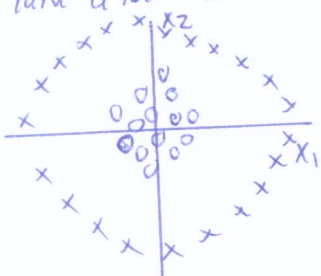We still don't know how to choose $\theta$s but

suppose $\begin{cases} \theta_0 = -3 \\ \theta_1 = 1 \\ \theta_2 = 1 \end{cases}$  $\vec{\theta} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$

$\leadsto$ $\boxed{\text{Predict}}$ $y=1$ if $\underbrace{-3 + x_1 + x_2}_{\vec{\theta}^T x} \geqslant 0$  rewriting $\boxed{x_1 + x_2 \geqslant 3}$

$\hookrightarrow y(x) = 0.5$ exactly.

$\leadsto$ The decision boundary is a property of the hypothesis including the params $\theta_0 \ldots \theta_2$. and not of the data set.

---- ✗ ----

Let's take a look at a more complex example



$\leadsto$ how can we fit $\boxed{\text{Log R}}$ params to solve this problem?

$\leadsto$ we can transform the feature space to higher order just like we did with polynomials.

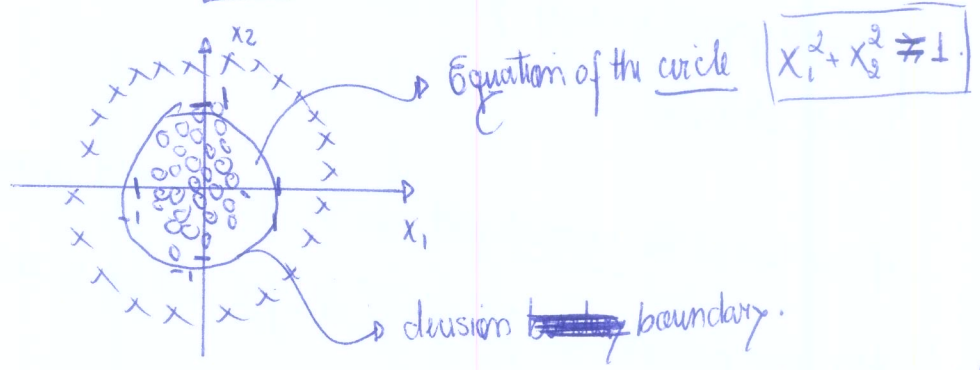$\hat{y}(x) = g\left(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 \boxed{x_1^2} + \theta_4 \boxed{x_2^2}\right)$

added features.

Suppose we have found $\vec{\Theta} = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$.

My predictor will predict $\boxed{y=1}$ (if) $-1 + x_1^2 + x_2^2 \geqslant 0$

$$\boxed{x_1^2 + x_2^2 \geqslant 1}$$



→ Equation of the circle $\boxed{x_1^2 + x_2^2 \not\geqslant 1}$

→ decision ~~boundary~~ boundary.
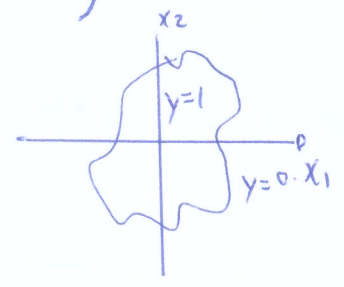
⟿ Once again, the decision boundary is a property of the hypothesis/model and not of the <u>dataset</u>

⟿ Can we have even more complex decision boundaries? (YES).

$$\hat{y}(x) = g\left( \Theta_0 + \Theta_1 x_1 + \Theta_2 x_2 + \Theta_3 x_1^2 + \Theta_4 x_1^2 x_2 + \Theta_5 x_1^2 \cdot x_2^2 + \Theta_6 x_1^3 x_2 \cdots \right)$$



—————— × ⋅——————

But, how do we [automatically] fit the params for a model/predictor?

$\underbrace{\phantom{model/predictor}}_{\text{for } \boxed{\log R}}$

training set $\left\{ \left( x^{(1)}, y^{(1)} \right), \cdots, \left( x^{(m)}, y^{(m)} \right) \right\}$

(m) examples $x \in \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix}_{R^{n+1}}$, $x_0 = 1$, $\underbrace{y \in \{0, 1\}}_{\text{classif. problem.}}$

$$\boxed{\hat{y}(x) = \frac{1}{1 + e^{-\theta^T x}}}$$
$\underbrace{\phantom{xxxxxxxxxxxx}}_{\text{hypothesis}}$

Question: how to choose parameters $\theta$?

Recall that in the linear Regression hypothesis, we had $J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} \left( \hat{y}(x^{(i)}) - y^{(i)} \right)^2$
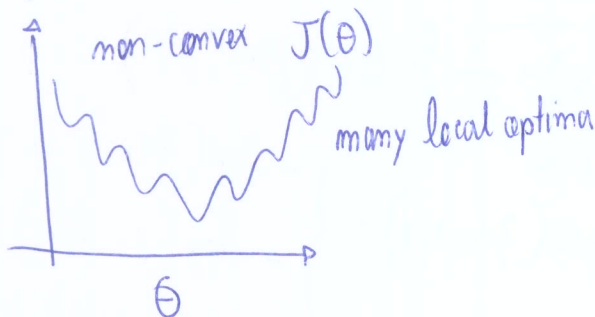
for $\boxed{\log R}$, let's call $\frac{1}{2} \left( \hat{y}(x^{(i)}) - y^{(i)} \right)^2$ as $\boxed{\text{Cost} \left( \hat{y}(x^{(i)}), y \right)}$

so $\text{Cost} \left( \hat{y}(x^{(i)}), y^{(i)} \right) = \frac{1}{2} \cdot \left( \hat{y}(x^{(i)}) - y^{(i)} \right)^2$

simplifying,

$$\text{Cost}(\hat{y}(x), y) = \frac{1}{2} \cdot \left( \hat{y}(x) - y \right)^2$$

for Log R has non linearity $\frac{1}{1 - e^{-\theta^T x}}$ complicated non linear function

this cost function works fine for $\boxed{LR}$ but we are focused on $\boxed{\text{Log R}}$ and if we minimize it for $\boxed{\text{LogR}}$, it will be clear it isn't a __convex function__



non-convex $J(\theta)$

many local optima __versus__

we would like $J(\theta)$

global min
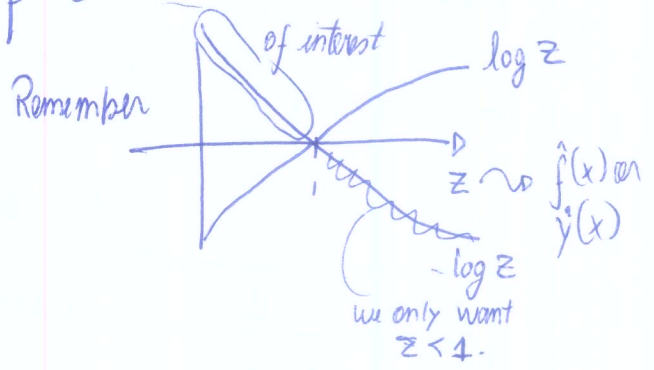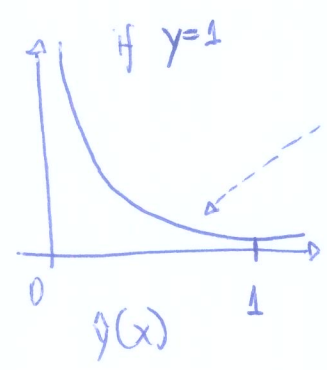
So we need to come up with a way of designing a different cost function that is __convex__.

prime

# Logistic Regression Cost function

$$Cost\left(\hat{f}(x), y\right) = Cost\left(\hat{y}(x), y\right) = \begin{cases} -\log\left(\hat{y}(x)\right) & \text{if } y=1 \\ -\log\left(1 - \hat{y}(x)\right) & \text{otherwise } (y=0) \end{cases}$$

if $y=1$



$\hat{y}(x)$

Remember

of interest

$\log z$

$z \sim \hat{f}(x)$ or $\hat{y}(x)$

$-\log z$

we only want $z < 1$.

and prediction is

$$\sim \begin{cases} Cost = 0 \text{ if } y=1, \hat{y}(x)=1 \\ \text{But} \\ \text{as } \hat{y}(x) \sim 0, \; Cost \sim \infty. \end{cases}$$

we penalize the learning alg. by a __very__ large cost as it commits mistakes (go further away from the corrects value/output).

if $y=0$.  $+\infty$



$\hat{y}(x)$

putting together



<u>Decision Boundary</u>.

---×---

# simplified Cost function and Gradient Descent

(LReg) cost function $\quad J(\theta) = \dfrac{1}{m}\sum_{i=1}^{m} Cost\left(\hat{y}(x^{(i)}, y^{(i)}\right)$

where

$$Cost\left(\hat{y}(x), y\right) = \begin{cases} -\log\left(\hat{y}(x)\right) & \text{if } y=1 \\ -\log\left(1 - \hat{y}(x)\right) & \text{otherwise} \end{cases}$$

$y \in \{0,1\}$ __binary problem__

Let's compress the two cases in one.

$$Cost\left(\hat{y}(x), y\right) = -y \log\left(\hat{y}(x)\right) - (1-y) \log\left(1 - \hat{y}(x)\right).$$

So now, we can design our cost function

as

$$\Downarrow$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} Cost\left(\hat{y}(x^{(i)}, y^{(i)})\right)$$

$$= \boxed{-\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)} \cdot \log\left(\hat{y}(x^{(i)})\right) + (1-y^{(i)}) \cdot \log\left(1 - \hat{y}(x^{(i)})\right)\right]}$$

We choose it because it can derived from statists using principle of MLE max. likelihood estim. and it is convex.

↳ To fit params $\vec{\theta}$: $\min_{\vec{\theta}} J(\vec{\theta})$.

↳ To make a prediction for a new ⓧ, Output $\hat{y}(x) = \frac{1}{1 + e^{-\vec{\theta}^T x}}$ which means $P(y=1 \mid x; \vec{\theta})$ learned

———————— x ————————

Using ⓖⓓ for minimizing $J(\vec{\theta})$.

want $\min_{\theta} J(\vec{\theta})$  Repeat {

$$\theta_j \leftarrow \theta_j - \alpha \cdot \underbrace{\frac{\partial}{\partial \theta_j} J(\theta)}_{**\alpha}$$

$$** \quad \frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \cdot \sum_{i=1}^{m} \left(\hat{y}(x^{(i)}) - y^{(i)}\right) \cdot x_j^{(i)}$$

what?  it look ✗✗ exactly the
same as before for ⓛⓡ.
But th
✗✗ ~~only~~ change is $\hat{y}(x)$
before $\hat{y}(x) = \theta^T x$.

now

now $\hat{y}(x) = \frac{1}{1 + e^{-\theta^T x}}$

In a _vectorized form_ :

$$
\begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix}_{m \times 1} \leftarrow \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix}_{m \times 1} - \alpha \left( \begin{bmatrix} Z(X) \end{bmatrix}_{m \times m} - \begin{bmatrix} Y_1 \\ \vdots \\ Y_m \end{bmatrix}_{m \times 1} * \begin{bmatrix} X \end{bmatrix}^{T}_{m \times m} \right)
$$

$$m \times 1$$

$$m \times m$$

[No] some adjustments (are) _necessary_ .

$$
\begin{bmatrix} f(A) \end{bmatrix}_{m \times m} \begin{bmatrix} B \end{bmatrix}_{m \times 1} = \begin{bmatrix} C \end{bmatrix}_{m \times 1}
$$

$$
\begin{bmatrix} f(A) \end{bmatrix}_{m \times 1} - \begin{bmatrix} Y \end{bmatrix}_{m \times 1} = \begin{bmatrix} C \end{bmatrix}_{m \times 1}
$$

## Logistic Regression - advanced optimization

So far, we have seen the opt. alg Gradient Descent and with it, we have a

cost function $J(\vec{\theta})$ and want $\min_\theta (J(\theta))$.

Given $\vec{\theta}$, we compute ① $J(\theta)$

② $\dfrac{\partial}{\partial \theta_j} J(\theta)$   $\forall j = 0, \ldots, n_{features}$.

and ⑤Ⓓ does: Repeat:

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta).$$

but this is one _alternative_. If we show how to compute $J(\theta)$ and $\dfrac{\partial}{\partial \theta_j} J(\theta)$, we

can use more sophisticated solutions such as $\begin{cases} \text{Conjugate Gradient} \\ BFGS \\ L\text{-}BFGS. \end{cases}$

Normally, such alternatives ① Do not need $\alpha$ (learning rate) manually selected.

② often faster than GD.

however, they are more _complex_.

## Example on how to use them.

$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$   $J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$   of course $\theta = \begin{bmatrix} 5 \\ 5 \end{bmatrix}$ minimizes it and that's what

we want to _find_.
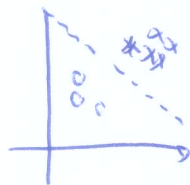
$\dfrac{\partial J(\theta)}{\partial \theta_1} = 2(\theta_1 - 5)$

$\dfrac{\partial J(\theta)}{\partial \theta_2} = 2(\theta_2 - 5)$

Normally, _all_ we need to _do_ is to write code for computing $J(\theta)$ and its

_derivatives_.

Example: digits, faces { profile left
                          profile right , age-group estimation, foldering/tagging photos, emails,
                          frontal

foldering { family
            business      $y \in \{1, \ldots, k\}$
            gym                            classes
            ⋮

Binary Classif:

Multi-class

Let's do it on-vs-all (OVA): first we transform our problem on 3 problems

$f_\theta^{(1)}(x)$                $f_\theta^{(2)}(x)$                $f_\theta^{(3)}(x)$

what we did was    $f_\theta^{(i)}(x) = P(y=i \mid x; \theta)$   $i \in \{1, \ldots, 3\}$

Basically, we train a (LogR) classifier $f_\theta^{(i)}(x)$ for each class (i) to predict

$P(y=i)$. For a new input (x), for predicting, we select the class that

maximizes $\boxed{\max_i f_\theta^{(i)}(x)}$

# Regularization and the problem of overfitting

Regularization will allow us to diminish overfitting problems.

Lets return to the problem of __housing__



size

$\theta_0 + \theta_1 x$.

$\rightsquigarrow$ underfit, high __bias__ $\Big\}$ algorithm has a very strong bias or pre-conception that the housing prices change linearly with the size even the data shows the contrary

both mean the model is not fitting the data very well.

alternatively, we could have



size

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

Overfit, "high variance".

the model fitted the data perfectly and it is to __just__ to it and the hypothesis is too variable since there is not enough data to prove it
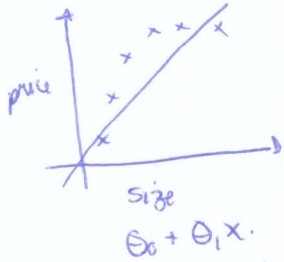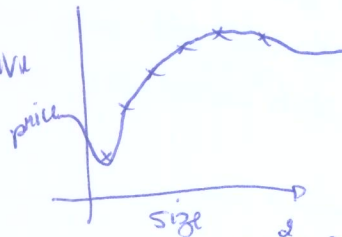
__Overfitting__ : if we have too many features, the learned hypothesis may fit ~~the~~ the training set very well $\left( J(\theta) = \dfrac{1}{2m} \sum_{i=1}^{m} f_\theta(x^{(i)}) - y^{(i)})^2 \approx 0 \right)$ but $\boxed{\text{fail}}$ to generalize to new examples (predict prices on new data).

The same thing can happen with logistic regression as well.



underfit

$f_\theta(x) = g\left( \theta_0 + \theta_1 x_1 + \theta_2 x_2 \right)$

g = sigmoid function

Just.

$f_\theta(x) = g\Big( \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 \Big)$

Overfit / high variance.

$f_\theta(x) = g\Big( \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 x_2^2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 \dots \Big)$.

# Addressing Overfitting



→ Plotting is one way but it wouldn't work for high-dimensional data

(Two) main options to deal with it

① Reduce number of features ⟨ manually which features to keep
model selection algorithm.
automatic alg. that select which features to
keep/throw away.

* however, sometimes this means you are throwing away important
info about the problem.

② Regularization

↳ Keep ALL the features but reduce magnitude/values/importance of
parameters $\theta_j$. works well when we have lots
of features, each of which contributing a bit for predicting $(y)$

---

## Regularization - Cost function

Intuition



$\theta_0 + \theta_1 x + \theta_2 x^2$



$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

Suppose we penalize and make $\theta_3$ and $\theta_4$ really small

$$\min_\theta \frac{1}{2m} \sum_{i=1}^{m} \left(f_\theta(x^{(i)}) - y^{(i)}\right)^2 + \boxed{1000\,\theta_3^2 + 1000\,\theta_4^2}$$

high penalization

this would give $\theta_3 \approx 0$
$\theta_4 \approx 0$

# Regularization Cost function

Besides the intuition, here's the general ideal:

$\sim$ small values for parameters $\theta_0, \theta_1, \ldots \theta_m$ } - Simpler hypothesis
- Less prone to Overfitting.

## Lets see the housing example

- features: $x_1, x_2, \ldots, x_{100}$
- params: $\theta_0, \theta_1, \ldots, \theta_{100} \in \mathbb{R}^{101}$

it is difficult to pick one or some to penalize beforehand. If it was easy, it would be
the same as selecting the less important features.

So, here's what we do:

$$J(\theta) = \frac{1}{2m} \left[ \overbrace{\sum_{i=1}^{m} \left( f_\theta(x^{(i)}) - y^{(i)} \right)^2}^{\text{good fitting}} + \overbrace{\lambda \sum_{i=1}^{m} \theta_j^2}^{\text{Regularization}} \right]$$

keep the params small.

→ not quadratic but smoother

→ previously.

→ do not penalize $\theta_0$

→ Regularization parameter

$\min_{\theta} J(\theta)$



price — size

$\lambda$ controls the tradeoff between good fitting to training data

(vs)

Keep the parameters small.

and with both we want to keep the hypothesis simple and avoid overfitting.

Question : what if $\lambda$ is too high (eg. $\lambda = 10^{10}$) in a regularized linear regression?

$\hookrightarrow$ all $\theta$s except $\theta_0$ will tend to (zero)
and the model will be biased



$\theta_0 + \boxed{\theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4}$ → 0 zero.

therefore $f_\theta(x) = \theta_0$.

} underfitting

too strong bias / pre-conception
the a straight line in $\theta_0$ will
fit well in spite of data on
the contrary.

↝ Be _carefull_ when choosing $\textcircled{\lambda}$.

——————— ✕ ———————

## Regularized Linear Regression

So far, for linear regression, we have seen <u>GD</u> and <u>Normal Equation</u> for fitting the parameters.

Here's the optimization function for regularized linear regression

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} \left( f_\theta(x^{(i)}) - y^{(i)} \right) + \lambda \sum_{i=1}^{m} \theta_j^2 \right]$$

⇓

$$\min_\theta J(\theta).$$

Previously, the $\textcircled{GD}$ for minimizing this was

$$\boxed{\begin{array}{l} \text{Repeat } \{ \\ \quad \theta_j \leftarrow \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^{m} \left( f(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \\ \qquad\qquad j = 0, \cdots, m. \\ \} \end{array}}$$

if we separate $\theta_0$:

$$\boxed{\begin{array}{l} \text{Repeat } \{ \\ \quad \theta_0 \leftarrow \theta_0 - \frac{\alpha}{m} \cdot \sum_{i=1}^{m} \left( f(x^{(i)}) - y^{(i)} \right) x_0^{(i)} \\ \quad \theta_j \leftarrow \theta_j - \frac{\alpha}{m} \cdot \sum_{i=1}^{m} \left( f(x^{(i)}) - y^{(i)} \right) x_j^{(i)}. \\ \} \end{array}}$$

⇓

$$\begin{array}{l} \text{Repeat } \{ \\ \quad \theta_0 \leftarrow \theta_0 - \alpha \cdot \underbrace{\frac{1}{m} \sum_{i=1}^{m} \left( f(x^{(i)}) - y^{(i)} \right) \cdot x_0^{(i)}}_{\frac{\partial J(\theta)}{\partial \theta_0}} \\ \quad \theta_j \leftarrow \theta_j - \alpha \left[ \underbrace{\frac{1}{m} \sum_{i=1}^{m} \left( f(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)} + \underbrace{\frac{\lambda}{m} \cdot \theta_j}_{\text{regularization}}}_{\frac{\partial J(\theta)}{\partial \theta_j} \text{ regularized}} \right] \\ \} \end{array}$$

We can regroup things and then

$$\Theta_j \leftarrow \Theta_j - \alpha \cdot \left[ \frac{1}{m} \sum_{i=1}^{m} \left( f(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)} + \frac{\lambda}{m} \cdot \Theta_j \right]$$

$$\Theta_j \leftarrow \Theta_j \underbrace{\left( 1 - \frac{\overset{small}{\alpha} \overset{small}{\lambda}}{\underset{large}{m}} \right)}_{< 1} - \frac{\alpha}{m} \cdot \sum_{i=1}^{m} \left( f(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)}$$

for instance $1 - \frac{\alpha \lambda}{m} < 1$

$1 - 0.01 = \boxed{0.99 \cdot \Theta_j}$ shrinks $\Theta_j$ a little bit.

what about **normal equations** ?

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \leftarrow \text{one training example}$$
$$M + (n+1)$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

$$\min_{\Theta} J(\Theta)$$

$$\overset{\circ}{\Theta} = \left( X^T X \right)^{-1} \cdot X^T y \quad \text{without Regularization}$$

for regularization;

$$\overset{\circ}{\Theta} = \left( X^T X + \lambda \begin{bmatrix} 0 & & & 0 \\ & 1 & & \\ & & 1 & \\ 0 & & & 1 \end{bmatrix} \right)^{-1} \cdot X^T y$$

$\underset{zeros}{\text{zeros}}$ $(n+1) \times (n+1)$

Example   $n = 2 \rightsquigarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

**what about non-invertibility now ?** $\overbrace{\text{(Advanced)}}^{\text{non-invertible / singular}}$

Suppose $\underset{\#\,exam}{m} < \underset{\#\,features}{n}$   $\overset{\circ}{\Theta} = \overbrace{\left( X^T X \right)^{-1}} \cdot X^T y$
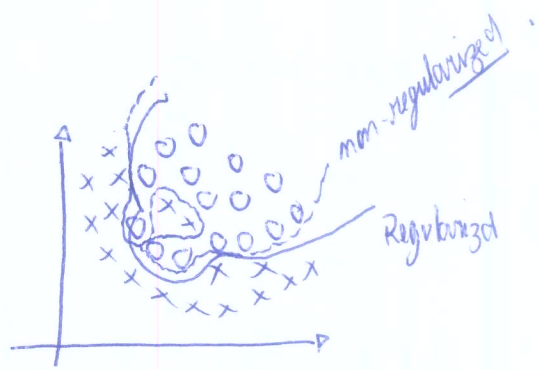
$\boxed{but}$ if $\lambda > 0$, $\overset{\circ}{\Theta} = \left( X^T X + \lambda \begin{bmatrix} 0 & & \\ & 1 & \\ & & \ddots \end{bmatrix} \right)^{-1} \cdot X^T y$

we'll have no singularity anymore. $\boxed{why}$ some features will decrease magnitude and simplify model.

# Regularized Logistic Regression

(RegL) can ~~only~~ also be prone to overfitting:



non-regularized

Regularized

$$\hat{f}_\theta(x) = \overset{sigmoid}{g}\left(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 \cdot x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots\right)$$

Cost function

$$J(\theta) = -\left[\frac{1}{m}\sum_{i=1}^{m} y^{(i)}\left(\log \hat{f}(x^{(i)})\right) + (1-y^{(i)}) \cdot \log\left(1-\hat{f}(x^{(i)})\right)\right] \overset{+}{\underset{\triangle}{}} A$$

introducing Regularization

We just add $\boxed{+ \dfrac{\lambda}{2m} \cdot \overset{m}{\underset{j=1}{\sum}} \theta_j^2 \cdot}$ for $\theta_1, \dots \theta_n$

\* dimensions

A

How to implement this?

Gradient Descent | Repeat {

$$\theta_0 \leftarrow \theta_0 - \alpha \cdot \frac{1}{m}\sum_{i=1}^{m}\left(f(x^{(i)}) - y^{(i)}\right) \cdot x_0^{(j)}$$

$$\theta_j \leftarrow \theta_j - \alpha\left[\frac{1}{m}\sum_{i=1}^{m}\left(f(x^{(i)} - y^{(i)})\right) \cdot x_j^{(i)} \oplus \boxed{+ \frac{\lambda}{m} \theta_j}\right]$$

regularization

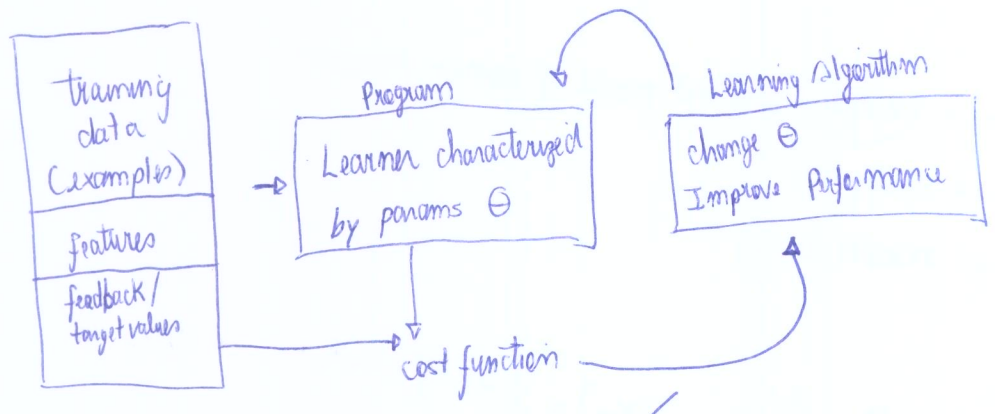$$j \in \{1, \dots, m\}.$$

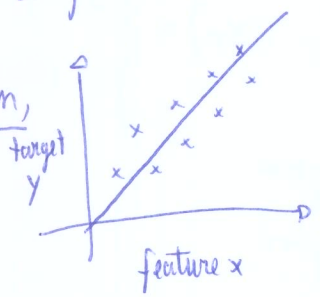$\underbrace{\qquad\qquad\qquad}_{\dfrac{\partial J(\theta)}{\partial \theta_j}}$

}

remember $f_\theta(x) = \dfrac{1}{1+e^{-\theta^T x}}$

Wrapping up what we had so far

(Notation) features X

targets Y

predictions $\hat{y}$

params $\Theta$

training data (examples)

features

feedback/ target values

→

Program
Learner characterized by params $\Theta$

Learning Algorithm
change $\Theta$
Improve Performance

cost function

In the case of <u>linear regression</u>,

target y

feature x
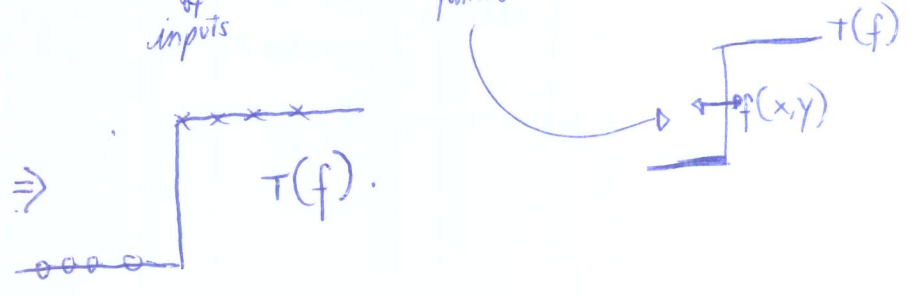
"Predictor"
Evaluate line
$r = \Theta_0 + \Theta_1 X_1$

$$\hat{y}(x) = \Theta_0 + \Theta_1 X_1$$

$$(y) \in \mathbb{R}.$$

while with classification, we seek to |predict| a discrete value/target $(y)$

Now that we already saw <u>linear regression</u> and <u>logistic regression</u> (classifier) let's
see a similar linear algorithm for classification ~~is~~ called <u>perceptron</u>.

<u>Perceptron with ② features</u>

$$\hat{f}(x) = \Theta_1 X_1 + \Theta_2 X_2 + \Theta_0.$$

$X_1$  $\Theta_1$
$X_2$  $\Theta_2$
$1$  $\Theta_0$
→ classifier
weighted sum of inputs

→ $T(f)$
threshold function

→ $\hat{C}(x)$
output of a $\{-1, +1\}$ class

$f(x)$

$x \ x \ x \ x$

$y$. $\Rightarrow$

$T(f)$.

$T(f)$

$f(x, y)$

→ Perceptron is a linear classifier

→ The $(w_s)$ will be the weights ($\theta$s earlier).
   $w_0$ is the intercept $w_0 = 1$ initially.

A perceptron calculates ② quantities
{
① a weighted sum of the input features
② a threshold on this sum by the ① function.
}

The perceptron is a $(simple)$ artificial model of human neurons

   weights → synapses
   threshold → neuron firing

Perceptron operation:

● $$\hat{f}(x_1, \ldots x_n) = \begin{cases} 1 & \text{if } \theta_1 x_1 + \ldots + \theta_m x_m > 0 \\ 0 & \text{otherwise (can also be defined as } -1) \end{cases}$$
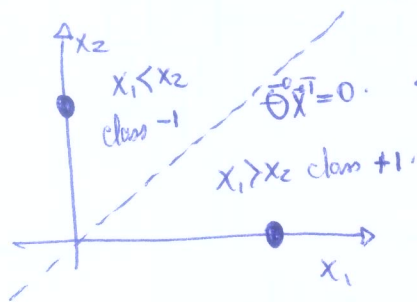
$\vec{\theta}$ → weight vector $\in \mathbb{R}^{m+1}$

$\vec{x}$ → feature vector $\in \mathbb{R}^{m+1}$.

$\theta_0 x_0 + \theta_1 x_1 + \ldots \theta_m x_m = \underbrace{\vec{\theta} \vec{x}^T}_{\text{vector inner product}}$

{
The perceptron represents a hyperplane decision surface in $n$-dimensional space
{
2d line $\downarrow$
3d plane
}

The equation of the hyperplane is $\underline{\vec{\theta}\vec{x}^T = 0}$
   This is the equation for points in x-space that are
   on the boundary.
}

Example



$\vec{\theta}\vec{x}^T = 0$. Suppose $\vec{\theta} = (\theta_1 = 1, \theta_2 = -1, \theta_0 = 0)$.

$x_1 < x_2$ class $-1$

$x_1 > x_2$ class $+1$.

$\vec{\theta}\vec{x}^T \Rightarrow 1 x_1 + (-1 . x_2) + 0(x_0) = 0$

$x_1 - x_2 = 0$

$\boxed{x_1 = x_2}$ decision boundary equation

* A data set is ~~linearly~~ [separable] by a learner (if) there is some instance of that learner that correctly predicts (all) data points

* The learner is [linearly separable] if:
  - Can separate the two classes using a straight line if $X \in \mathbb{R}^2$ or with a hyperplane if $X \in \mathbb{R}^m$.

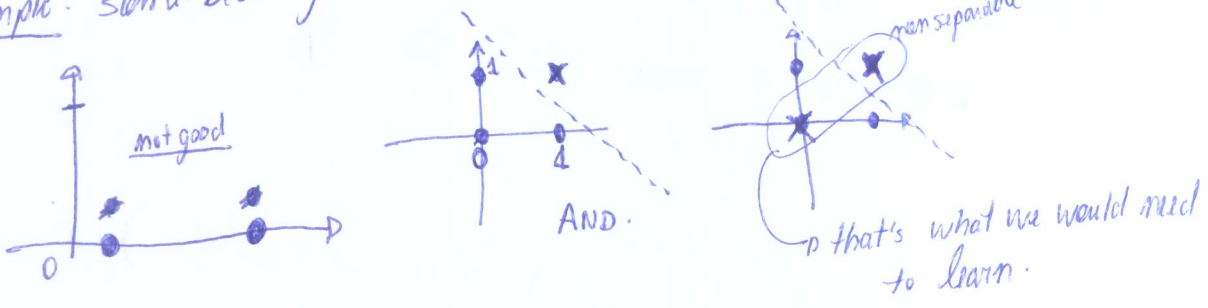Class Overlap        ~~Overlap~~ { Common in practice
                     Some observation values possible for both classes ( benign/malign cells look similar).

                     what to do? { - Increase number of features.
                                    - Increase model complexity.

# Representational Power of perceptrons

↝ what mappings a perceptron can represent perfectly?
well, it is a linear classifier, so it can represent any mapping that is [linearly separable]

Example: some boolean functions (AND) but no [XOR]



not good

AND.

non separable

↳ that's what we would need to learn.

[Remember] Effects of dimensionality

① Data are increasingly separable in high dimension — Is this a good thing?

Good { - Separation is easier in higher dimensions (for fixed (M))
       - Increase nbr of features, and even a linear classifier would do the trick!          (↳ nbr examples.

Bad { - training (vs) test error; overfitting; bias (vs) variance
      - Increasingly complex decision boundaries can eventually get all training data right but it doesn't mean necessarily good for testing data ⟹ [lacks generalization]
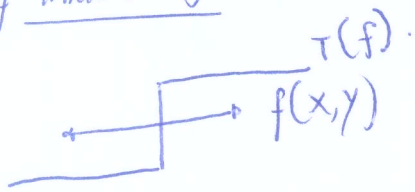
<u>How to update weights?</u>

for each data point $x^{(j)}$:

$$f(x^{(j)}) \leftarrow T(\vec{w} * x^{(j)})$$
$$\vec{w} \leftarrow w + \alpha \left( f(x^{(j)}) - y^{(j)} \right) \cdot x^{(j)} \quad \underline{\text{gradient-like step}}$$
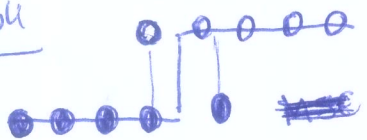
○ if $f(x^{(j)})$ correctly predicted : <u>no update</u>.

otherwise : update weights.

Another way of <u>thresholding</u> would be using smoother functions

$T(f)$.

$\rightarrow f(x,y)$

$g(f)$ sigmoid for instance

$\rightarrow f(x,y)$

$\begin{cases} \text{if we are far from decision boundary, } |f(\cdot)| \text{ is large, small error.} \\ \text{nearby the boundary : } (f(\cdot)) \text{ near } \frac{1}{2}, \text{ large error.} \end{cases}$

<u>Example</u>

$J(\theta) = \frac{2}{10}$ normal thresholding

<u>while</u>

$J(\theta) = \left( 0^2 + 1^2 + .2^2 + .25^2 + \dots \right) / 10.$

<u>Summary</u>
- linear classifier $\iff$ perceptron
- visualizing the decision boundary : useful but not <u>always</u> possible
- Measuring quality : Sum of squared errors (SSE).
- Learning the weights of a linear classifier reduces to an optimization problem. For (SSE) we can do gradient descent but there are others

# Neural Networks — Non-linear hypotheses

- Quite old idea
- Out of favor for sometime but in the hype again for several recent results such as Deep Belief Networks (Deep learning) and results such as image nets
- Why do we need yet another learning algorithm?

Consider a problem



→ we could use LogReg such as

$$\hat{f}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1 x_2^2 + \dots)$$

poly terms.

→ normally it works ok for a few features (2,3) but not as good for many features. (how to derive polynomials then?)

if we have 100 features and want to just include 2nd order polynomials, we would have

$$\hat{f}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_{100} x_{100} + \theta_{101} x_1^2 + \theta_{102} x_1 x_2, \dots, \theta_K x_1 x_{100} + \theta_{K+1} x_2 x_3 + \dots) \approx 5K \text{ features. } \theta(n^2)$$

orig. feature space. $\approx \frac{n^2}{2}$

↝ If we want 3rd order polys, we would have $\theta(n^3)$. For 100 features ≈ 170K features!

↝ In practice, we have problems with thousands or millions of features in the original feature space.
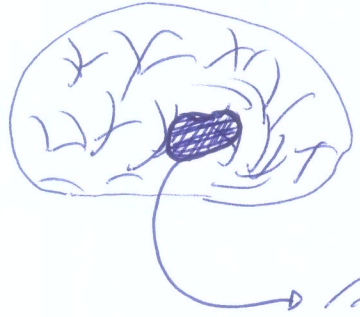
———————— + ————————

# Neurons and the brain

Origins:
{ algorithms that try to mimic the brain.
widely used in the 80's and early 90s. Popularity went down in late 90s.
Recent hype due to state-of-the-art results in many applications.
the reason is that now we have enough computational power to design and use large-scale neural networks

expensive to run

The "one learning algorithm" hypothesis

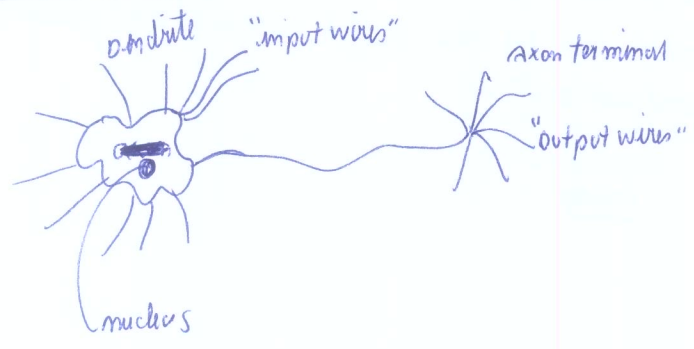→ Some people believe the brain does all of its amazing things based on one learning algorithm.

Auditory cortex learns to ~~see~~ hear/listen but researchers have shown that cutting the connection of Auditory cortex to the _ear_ and rewiring it to the _eye_ leads to a learning process in which the auditory cortex learns to see!

neuro-rewiring experiments

Examples: — Braimport experiment: Camera for low-res grayscale image on top of the head and a wire to connect each pixel to our tongue (low voltage → dark pixel and we can "see" through our tongue.                                              high voltage → bright pixel)

— Human echolocation (sonar): snap fingers. Boy that had his eyeballs removed learned to navigate using this "rewiring" technique.

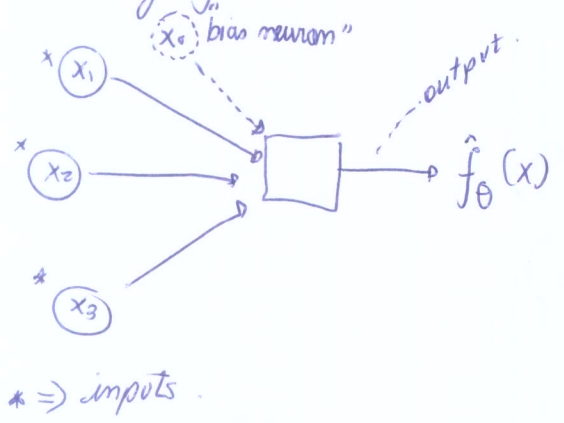— Implanting a 3rd eye in a frog will lead him to actually learn how to use it.

---- * ----

Neural Networks - Representations

Dendrite    "input wires"                    Axon terminal

                                             "output wires"

nucleus

Computatiolly speaking, we can think of a neuron as a computational unit that receives a set of inputs, performs some computation and outputs some signals through the axons to other neurons.

We are going to model this as a _logistic unit_



$$\hat{f}_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\vec{x} = \begin{bmatrix} x_0 \\ \vdots \\ x_m \end{bmatrix} \qquad \vec{\theta} = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix}$$

* ⟹ inputs.

$x_0$ ⟹ bias neuron always equal to ①.

In (NN) terminology the sigmoid/logistic activation function and the $\vec{\theta}$ parameters are known as _weights_.

represents

Generalizing, a (NN) is just a group of several neurons.



Layer ①    Layer ②    Layer ③

⇓

input layer    hidden layer.    output layer.

because it is not Ⓧ or Ⓨ (input/output)
therefore the term "hidden".

↝ We may have $\boxed{NN_s}$ with several hidden _layers_.

$a_\lambda^{(j)}$ = "activation" of unit $(i)$ in layer $(j)$

$\Theta^{(j)}$ = matrix of weights controlling function mapping from layer $(j)$ to layer $(j+1)$.

Here's the computation :

$$a_1^{(2)} = g\left( \Theta_{10}^{(i)} x_0 + \Theta_{11}^{(i)} \cdot x_1 + \Theta_{12}^{(i)} x_2 + \Theta_{13}^{(i)} x_3 \right)$$

$$a_2^{(2)} = g\left( \Theta_{20}^{(i)} x_0 + \Theta_{21}^{(i)} x_1 + \Theta_{22}^{(i)} x_2 + \Theta_{23}^{(i)} x_3 \right).$$

$$a_3^{(3)} = g\left( \Theta_{30}^{(i)} x_0 + \Theta_{31}^{(i)} x_1 + \Theta_{32}^{(i)} x_2 + \Theta_{33}^{(i)} x_3 \right).$$

↳ $f_\theta(x) = a_1^{(3)} = g\left( \Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right).$

$\Theta^{(1)}$ is a matrix mapping layer $(1)$ to $(2)$, so $\Theta^{(1)} \in \mathbb{R}^{3\times4}$.

↝ If network has $(s_j)$ units in layer $(j)$, and $(s_{j+1})$ in layer $(j+1)$, then $\Theta^{(j)}$ will

have dimensions $(s_{j+1}) \times (s_j + 1)$.

↝ with this, we have ▸ $\boxed{\hat{f}_\Theta(x)}$

our classifier.

————————————————— ✳ —————————————————

model Representation, intuition and calculations

forward propagation : vectorized implementation

let's call

$z_1^{(2)}$, therefore $a_1^{(2)} = g\left( z_1^{(2)} \right)$.

$$a_1^{(2)} = g\left( \underbrace{\Theta_{10}^{(i)} x_0 + \Theta_{11}^{(i)} x_1 + \Theta_{12}^{(i)} x_2 + \Theta_{13}^{(i)} x_3} \right)$$

$$a_2^{(2)} = g\left( \Theta_{20}^{(i)} x_0 + \Theta_{21}^{(i)} x_1 + \Theta_{22}^{(i)} x_2 + \Theta_{23}^{(i)} x_3 \right) \quad z_2^{(2)}$$

$$a_3^{(2)} = g\left( \Theta_{30}^{(i)} x_0 + \Theta_{31}^{(i)} x_1 + \Theta_{32}^{(i)} x_2 + \Theta_{33}^{(i)} x_3 \right) \quad z_3^{(2)}.$$

$\boxed{\Theta^{(1)} \cdot \vec{x}}$ or $\boxed{\overset{\to}{\Theta}{}^{(1)} \overset{\to}{x}}$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_3 \end{bmatrix} \quad ; \quad Z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad ; \quad$$

$$Z^{(2)} = \vec{\Theta}^{(1)} \vec{x}$$

$$a^{(2)} = g\left(Z^{(2)}\right).$$
$\longrightarrow \mathbb{R}^3. \qquad \longrightarrow \mathbb{R}^3$
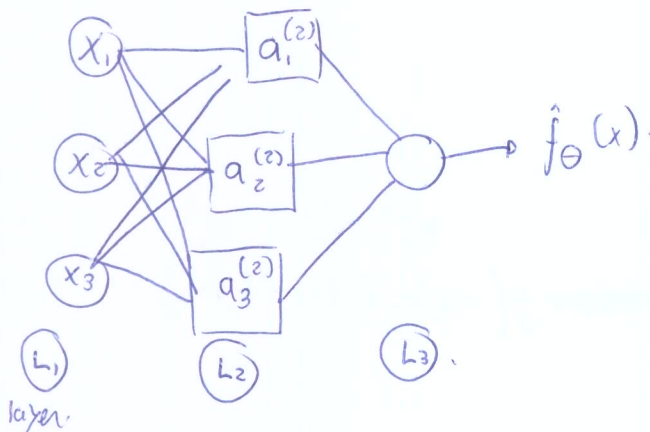
$g$ function applies __sigmoid__ function on each $z_i^{(i)}$ element-wise. Given all elements in $z^{(i)}$ (e.g. in $\mathbb{R}^3$), we apply $g$ on each individually.

Let's define $a^{(1)} = \vec{x}$, therefore $z^{(2)} = \vec{\Theta}^{(1)} \cdot \vec{x}$ now becomes $\boxed{z^{(2)} = \vec{\Theta}^{(1)} \cdot a^{(1)}.}$

Now we $\boxed{\text{add}}$ $\boxed{a_0^{(2)} = 1}$ so that $a^{(2)} \in \mathbb{R}^4$
$\underset{\text{bias unit}}{} \qquad \underline{\text{and}}$

$$z^{(3)} = \vec{\Theta}^{(2)} \cdot a^{(2)}$$

$$\hat{f}_\Theta(x) = a^{(3)} = g\left(z^{(3)}\right).$$

This algorithm is known as $\boxed{\text{forward propagation}}$.

---

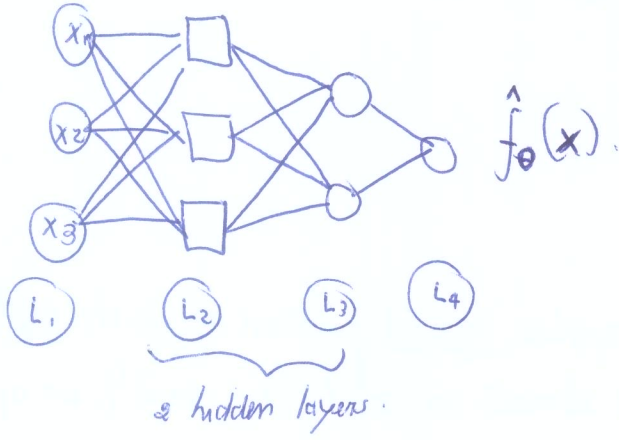Neural Network learning its own features



L₁ layer

* Suppose we don't see the left part of it. if we do that, we end up having a __logistic Regression__ like situation

$$\hat{f}_\Theta(x) = g\left( \Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)} \right).$$
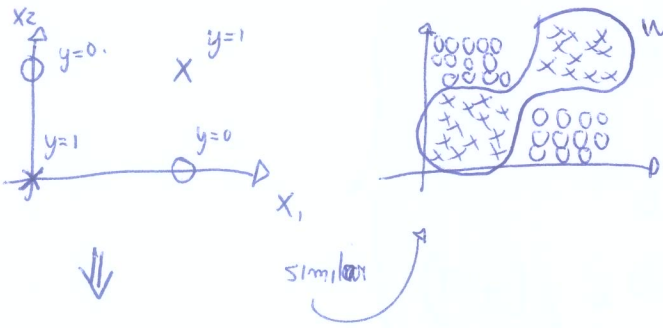
The $\boxed{\text{nice}}$ thing here is the the features $\boxed{a_\lambda}$ are __learned__ from the input.

Observation : we can have other (NN) diagrams/architectures

└─ how (NN) ~~units~~ units connect
to each other.



$\hat{f}_\theta(x)$

2 hidden layers.

---

Let's return to ~~our~~ logical operators _example_ in which $x_1$ and $x_2$ are binary $(0$ or $1)$.



we want to learn a non-linear separator/classifier

similar
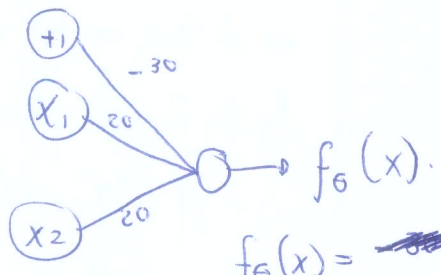
$$y = x_1 \text{ XOR } x_2$$
$$y = \text{XNOR } x_2$$
$$\boxed{\text{Not}(x_1 \text{ XOR } x_2)}.$$

Let's start with a _simple_ example for logical (AND).

$$x_1, x_2 \in \{0, 1\}$$
$$y = x_1 \text{ AND } x_2.$$



$f_\theta(x).$

$\theta_{10}^{(1)} \quad \theta_{11}^{(1)} \quad \theta_{12}^{(1)}$

$$f_\theta(x) = ~~\text{xxxx}~~ g\left(-30 x_0 + 20 x_1 + 20 x_2\right)$$

sigmoid



| $x_1$ | $x_2$ | $\hat{f}_\theta(x)$ |
|---|---|---|
| 0 | 0 | $g(-30) \approx 0$ |
| 0 | 1 | $g(-10) \approx 0$ |
| 1 | 0 | $g(-10) \approx 0$ |
| 1 | 1 | $g(10) \approx 1$ |

$$\boxed{f_\theta(x) \approx x_1 \text{ AND } x_2}$$

Example for (OR) function



| $x_1$ | $x_2$ | $\hat{f}_\theta(x)$ |
|---|---|---|
| 0 | 0 | $g(-10) \approx 0$ |
| 0 | 1 | $g(+10) \approx 1$ |
| 1 | 0 | $g(+10) \approx 1$ |
| 1 | 1 | $g(40) \approx 1$ |

$$\hat{f}_\theta(x) = g\left(-10 x_0 + 20 x_1 + 20 x_2\right).$$

Example for [NOT $x_1$]



| $x_1$ | $\hat{f}_\theta(x)$ |
|---|---|
| 0 | $g(10) \approx 1$ |
| 1 | $g(-10) \approx 0$ |

$$\hat{f}_\theta(x) = g\left(10 - 20 x_1\right)$$

Exercize: how to compute
$\rightsquigarrow$ $\left(\text{Not } \textcircled{x_1}\right) \text{AND} \left(\text{Not} \textcircled{x_2}\right)$ ?

---

$\textcircled{x_1}$ NOR $\textcircled{x_2}$



| $x_1$ | $x_2$ | $a_1^{(2)}$ | $a_2^{(2)}$ | $\hat{f}_\theta(x)$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 0 | 1 |

$\begin{cases} a_1^{(2)} \text{ is using } x_1 \text{ AND } x_2 \text{ network} \\ a_2^{(2)} \text{ is using } \neg x_1 \text{ and } \neg x_2 \text{ or } \left(\text{not } x_1\right) \text{AND} \left(\text{not } x_2\right). \end{cases}$
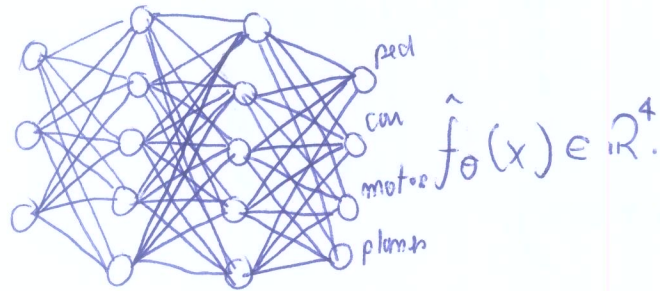
\* Each layer in the sequence can learn more complex functions ~~them~~ on top of its predecessors output layers.

---

Basically, it is an extension of the one (vs) all method.

Pedestrian    Car    motorcycles    Planes



$$\hat{f}_\theta(x) \in \mathbb{R}^4.$$

ped
car
motor
plane

We want $\hat{f}_\theta(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ when pedestrian, $\hat{f}_\theta(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ when car, $\hat{f}_\theta(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ when

motorcycles and $\hat{f}_\theta(x) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$ when plane.

In the training set, $\left(x^{(1)}, y^{(1)}\right), \left(x^{(2)}, y^{(2)}\right), \left(x^{(3)}, y^{(3)}\right), \ldots, \left(x^{(m)}, y^{(m)}\right)$

$y^{(i)}$ is going to be one of $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

                                               ped       car      motor     plane

So previously, we have represented $y \in \{1, 2, 3, 4\}$ (e.g., K-NN). Now we represent it in a binary form. So, for 4 classes,

$$\hat{f}_\theta(x) \underset{\in \mathbb{R}^4}{\approx} \underset{\in \mathbb{R}^4}{y^{(i)}}$$
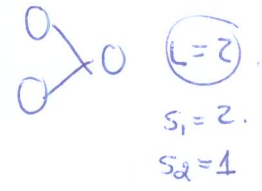
———————— x ————————

# Cost function for Neural Networks

~ Let's start with the classification problem

$$\{ (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)}) \}.$$

$L = 2$
$S_1 = 2.$
$S_2 = 1$

$(L)$ will be the number of layers in the network.

$(S_\ell)$ = number of units ( excluding the "bias" unit ) in layer $(\ell)$.

In the binary classification, $y=0$ or $y=1$, therefore we will have $(1)$ output unit. $\mathbb{R}^1$. For a multiclass classification problem ( $K$ classes), $y \in \mathbb{R}^K$, therefore we will have $(K)$ output units.

Let's now define a $\boxed{\text{cost function}}$. Recall that for logistic regression, we had :

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \log f_\theta(x^{(i)}) + (1 - y^{(i)}) \cdot \log \left( 1 - \hat{f}_\theta(x^{(i)}) \right) \right] + \frac{\lambda}{2m} \sum_{j=1}^{m} \Theta_j^2$$

Regularization
$\Theta_0 = 1$

for a _neural network_ , instead of having just one logistic regression output unit, we will have $(K)$ of _them_.

$$\hat{f}_\theta(x) \in \mathbb{R}^K \quad \left( \hat{f}_\theta(x) \right)_i = (i^{th}) \text{ output}.$$

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log(\hat{f}_\theta(x^{(i)}))_k + (1 - y_k^{(i)}) \cdot \log \left( 1 - \hat{f}_\theta(x^{(i)}) \right)_k \right] +$$

$$\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_\ell} \sum_{j=1}^{S_{\ell+1}} \left( \Theta_{ji}^{(l)} \right)^2$$

Wow! How can we optimize such cost function to actually find our weights or parameters ?
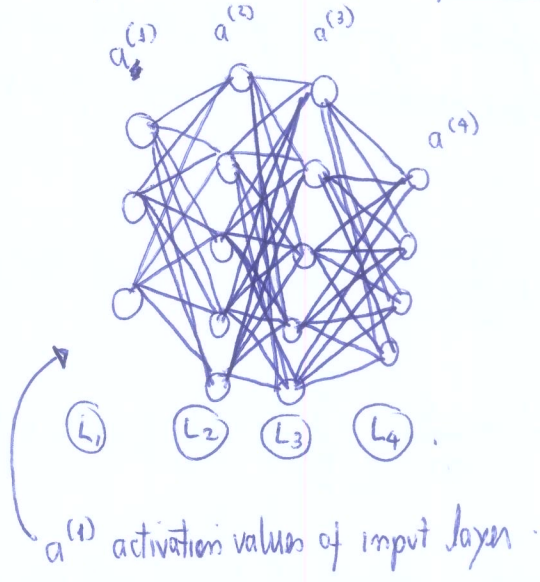
## Backpropagation Algorithm

Given our $J(\theta)$, we want to $\min_\theta J(\theta)$. For that we need to compute:

① $J(\theta)$

② $\dfrac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$ ; recall $\theta_{ij}^{(l)} \in \mathbb{R}$.

Suppose we only have one training example $(x, y)$. Using the _forward propagation_, we have

$$a^{(1)} = x.$$
$$z^{(2)} = \theta^{(1)} \cdot a^{(1)}$$
$$a^{(2)} = g(z^{(2)}) \text{ add } a_0^{(2)}$$
$$z^{(3)} = \theta^{(2)} \cdot a^{(2)}$$
$$a^{(3)} = g(z^{(3)}) \cdot \text{ add } a_0^{(3)}$$
$$z^{(4)} = \theta^{(3)} a^{(3)}$$
$$a^{(4)} = \hat{f}_\theta(x) = g(z^{(4)}).$$

$a^{(1)} \quad a^{(2)} \quad a^{(3)}$
$a_0^{(1)}$
$a^{(4)}$

$L_1 \quad L_2 \quad L_3 \quad L_4$.

$a^{(1)}$ activation values of input layer.

### Gradient computation: Backpropagation algorithm

__Intuition__ $\delta_j^{(l)} = $ "error" of node $(j)$ in layer $(l)$.

$a_j^{(l)} \Rightarrow$ activation of $(jth)$ unit in layer $(l)$.

__Concretly__, if we have the $(NN)$ above, for each output unit (layer $L = 4$),

$$\delta_j^{(4)} = \boxed{a_j^{(4)}} - y_j$$
$$\big( p(f_\theta(x)) \big);$$

if we vectorize it, we have $\vec{\delta}^{(4)} = \vec{a}^{(4)} - \vec{y}$.

$\in \mathbb{R}^K$ where $(K)$ is the number of output units in the $(NN)$.

__After__ that, we compute the $(\delta)$ terms for the _earlier layers_:

$$\delta^{(3)} = \left(\theta^{(3)}\right)^T \cdot \delta^{(4)} \circledast g'(z^{(3)}).$$
$$\delta^{(2)} = \left(\theta^{(2)}\right)^T \delta^{(3)} \circledast g'(z^{(2)}).$$

$(\circledast)$ element-wise multiplication.

The derivative $g'(z^{(3)})$ is basically $\underset{\text{vector of one}}{\underbrace{\phantom{xxx}}}$

$$g'(z^{(3)}) = \underbrace{a^{(3)}}_{\text{vector}} \cdot * \underbrace{(\vec{1} - a^{(3)})}_{\text{vector}}.$$

$$g'(z^{(2)}) = a^{(2)} \cdot * (1 - a^{(2)}).$$

↝ There is no $\delta^{(1)}$ since they are the features we observe in our training set

↝ The name **back propagation** comes from the fact we start in the output layer, calculate the error $\delta^{(4)}$ there and go back to the 3rd layer & where we compute $\delta^{(3)}$ and go back to layer ② and calculate $\delta^{(2)}$ in such a way it is propagating the error from layer $(l+1)$ to ① and from ⓛ to $(l-1)$ and so on.

**finally** $\boxed{\dfrac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)}}$ when ignoring the regularization term $(\lambda = 0)$.

Now, for a larger training set $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})\}$, set

$$\Delta_{ij}^{(l)} = 0 \quad \forall i,j,l.$$

↳ will be used as accumulators for computing $\dfrac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$.

For $\underline{i=1}$ to $\underline{m}$ {    // on each iteration, we work with one example.

    Set $a^{(1)} = x^{(i)}$

    Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \ldots, \underline{L}$.

    Using $y^{(i)}$, compute $\boxed{\delta^{(L)} = a^{(L)} - y^{(i)}}$ . // computing error for this example for the output layer.

    Compute $\delta^{(L-1)}, \delta^{(L-2)}, \ldots, \delta^{(2)}$

          used to compute the // partial derivatives in the end $\dfrac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$.

    $\boxed{\Delta_{ij}^{(l)} \leftarrow \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}}$

}

↳ if we vectorize, $\boxed{\Delta^{(l)} \leftarrow \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T}$.

$\begin{cases} D_{ij}^{(l)} \leftarrow \dfrac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \text{ if } j \neq 0 \\ D_{ij}^{(l)} \leftarrow \dfrac{1}{m} \cdot \Delta_{ij}^{(l)} \text{ if } j = 0 \cdot 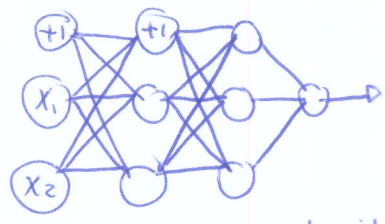\text{ (bias term) case} \end{cases}$ is exactly $\boxed{\dfrac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}}$
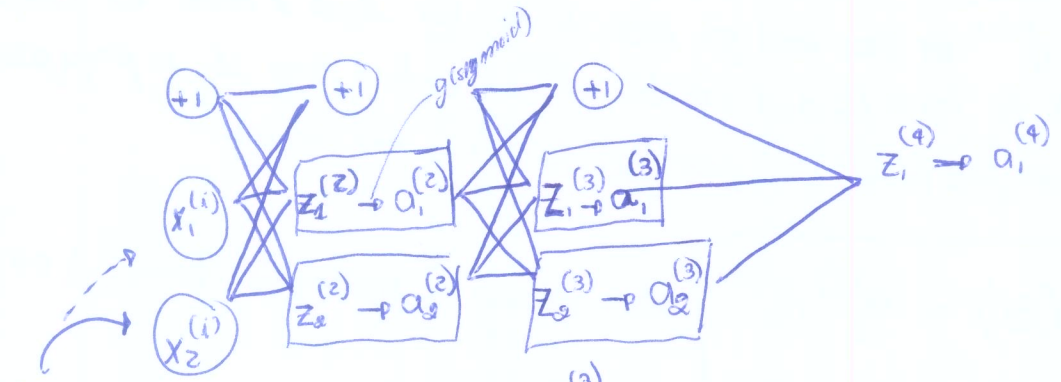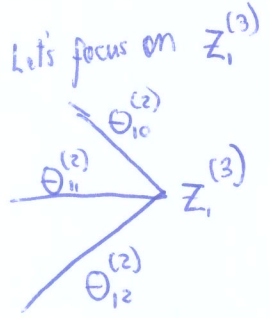
# Backpropagation Intuition

mechanical steps.

Let's see again _forward propagation_



2 units 2 un  2 un  1 unit
not counting bias



$g$ (sigmoid)

$(x^{(i)}, y^{(i)})$
single example.

Let's focus on $Z_1^{(3)}$

$\Theta_{10}^{(2)}$
$\Theta_{11}^{(2)}$ $\quad Z_1^{(3)}$
$\Theta_{12}^{(2)}$

So $\boxed{Z_1^{(3)} = \Theta_{10}^{(2)} \times 1 + \Theta_{11}^{(2)} \cdot a_1^{(2)} + \Theta_{12}^{(2)} a_1^{(2)}}$

this is forward propagation.

×

# What does Backpropagation do?

Recall its formulation

$$J(\theta) = -\frac{1}{m}\left[ \sum_{i=1}^{m} y^{(i)} \log\left(f_\theta(x^{(i)})\right) + (1-y^{(i)}) \cdot \log\left(1 - f_\theta(x^{(i)})\right) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1}\sum_{i=1}^{S_l}\sum_{j=1}^{S_{l+1}} \left(\Theta_{ji}^{(l)}\right)$$

Regularization.

focusing on a single example $(x^{(i)}, y^{(i)})$, the case of ① output unit and $\boxed{\lambda = 0}$

$$\text{Cost}(i) = y^{(i)} \log f_\theta(x^{(i)}) + (1-y^{(i)}) \log \left(1 - f_\theta(x^{(i)})\right)$$

the correct form is

think of $\text{Cost}(i) \approx \left(f_\theta(x^{(i)}) - y^{(i)}\right)^2$

How to th ⓃⓃ on example ⓘ?

\* $\log\left(f_\theta(x^{(i)})\right)$
but should be
$\log\left(1 - f_\theta(x^{(i)})\right)$

# Forward Propagation



+1   +1   +1

$x_1$

$x_2$

$\delta_j^{(l)} = $ error of cost for $a_j^{(l)}$ ( unit $j$ in layer $l$ ).

Formally, $\delta_j^{(l)} = \dfrac{\partial}{\partial z_j^{(l)}}$ Cost (i) $\left(\text{for } j \geq 0\right)$ where

change here affects *

$$\text{Cost}(i) = y^{(i)} \cdot \log \underbrace{f_\theta(x^{(i)})}_{* \text{ here}} + (1 - y^{(i)}) \cdot \log \underbrace{\left(1 - f_\theta(x^{(i)})\right)}_{\text{and } * \text{ here}}$$

Let's see it:

for $\quad \delta_1^{(4)} = y^{(i)} - a_1^{(4)} \quad$ output.

how to get $\quad \delta_2^{(2)} = ?$



activation function g

$$\begin{cases} \delta_2^{(2)} = \Theta_{12}^{(2)} \, \delta_1^{(3)} + \Theta_{22}^{(2)} \cdot \delta_2^{(3)} \\[2mm] \text{and} \\[2mm] \delta_{(2)}^{(3)} = \Theta_{12}^{(3)} \cdot \delta_1^{(4)} \end{cases}$$

———————— * ————————

# Implementation Note

Example

$s_1 = 10$ units

$s_2 = 10$

$s_3 = 1$ output unit

$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}$     $D^{(1)} \in \mathbb{R}^{10 \times 11}$   (derivatives)

$\Theta^{(2)} \in \mathbb{R}^{10 \times 11}$     $D^{(2)} \in \mathbb{R}^{10 \times 11}$

$\Theta^{(3)} \in \mathbb{R}^{1 \times 11}$     $D^{(3)} \in \mathbb{R}^{10 \times 11}$.

If we want to <u>use</u> this with more advanced techniques of optimization, we may need to deal with them as ~~matrices~~ vectors instead of as matrices.
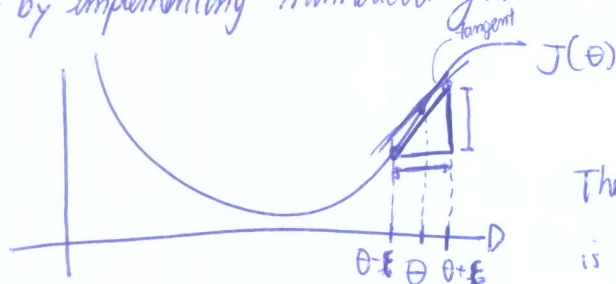
E.g., when optimizing in <u>Octave</u>.

———————— x ————————

# Gradient Checking

How can we check the Backpropagation algorithm along with its optimization function for finding the parameters (e.g., gradient descent) are actually working?

~ • We start by implementing numerical gradient checking.



$\theta \in \mathbb{R}$.

The slope of the line $\left( J(\theta+\varepsilon), (\theta+\varepsilon); J(\theta-\varepsilon), (\theta-\varepsilon) \right)$ is an _approx_ for the derivative of $J(\theta)$.

$$\text{slope} = \frac{\text{vertical height}}{\text{horizontal width}} = \frac{J(\theta+\varepsilon) - J(\theta-\varepsilon)}{2\varepsilon}$$

Therefore $\quad \dfrac{\partial}{\partial \theta} J(\theta) \approx \dfrac{J(\theta+\varepsilon) - J(\theta-\varepsilon)}{2\varepsilon} \quad \Big\}$ two sided difference.

$\varepsilon = 10^{-4}$ (the smaller the $\varepsilon$, the closer to the actual derivative.

In a more general _case_, $\theta \in \mathbb{R}^n$ (e.g., $\theta$ is "unrolled" version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

unrolled matrix $2 \times 2 \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$.

$\theta = \left[ \theta_1, \theta_2, \ldots, \theta_m \right]^T$.

we can do:

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\boxed{\theta_1 + \varepsilon}, \theta_2, \theta_3, \ldots \theta_m) - J(\boxed{\theta_1 - \varepsilon}, \theta_2, \ldots, \theta_m)}{2\varepsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \boxed{\theta_2 + \varepsilon}, \theta_3, \ldots \theta_m) - J(\theta_1, \boxed{\theta_2 - \varepsilon}, \theta_3, \ldots \theta_m)}{2\varepsilon}$$

$$\frac{\partial}{\partial \theta_m} J(\theta) \approx \frac{J(\theta_1, \theta_2, \ldots \boxed{\theta_m + \varepsilon}) - J(\theta_1, \theta_2, \ldots, \boxed{\theta_n - \varepsilon})}{2\varepsilon}$$

These equations allow us to numerically estimate the derivatives with respect to any parameter $\theta_i$ and then can be used to double check the derivatives calculated by backpropagation.
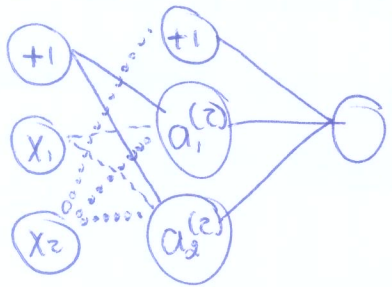
↳ Backpropagation is <u>very</u> efficient compared to gradient numerical checking. Therefore we <u>should</u> use it just in one or a few iterations to ~~double~~ double check backpropagation results.

---------- x ----------

## Random Initialization

We need to pick some initial **value** for $\Theta$.

Initializing $\Theta$ with <u>zeros</u> worked OK for logistic regression but is doesn't for (NNs).

The reason is that the network will be <u>redundant</u>.

$$a_1^{(2)} = a_2^{(2)} \quad \text{and} \quad \delta_1^{(2)} = \delta_2^{(2)}.$$

After each update, parameters corresponding to inputs going into each of two hidden units are identical



$$\boxed{a_1^{(2)} = a_2^{(2)}} \text{ if } \Theta = 0 \text{ for all elements in the beggining}$$

This is known as the problem of <u>symmetric weights</u>. For breaking the symmetry, we <u>must</u> iniatialize $\Theta_{ij}^{(l)}$ to a random value in $[-\varepsilon, \varepsilon]$ $(\text{this means } -\varepsilon \le \Theta_{ij}^{(l)} \le \varepsilon)$

This is not $\varepsilon$ of before.

---------- x ----------

## Wrapping up

The first thing we need to do when training a (NN), is to select a (NN) architecture (connectivity pattern)

How to choose this?

① input units: dimensionality of the problem $x^{(i)}$.

② output units: number of classes.

③ for the internal (hidden layers), a reasonable default is [One] hidden layer. But if we are using more than one hidden layer, it is interesting to have the same number of ~~&~~ units in all layers ( the more, the **better**, normally).
                                                                        ↳ hidden layers

Note : 
{ the more hidden layers, the better at the cost of much more computation.

{ nbr of units in hidden layer normally is comparable to that in the input layer.

## Training a neural network

① Initialize weights randomly.

② Implement forward propagation to obtain $f_\Theta(x^{(i)})$ $\forall x^{(i)}$.

③ Compute cost function $J(\Theta)$.

④ Backpropagation to compute partial derivatives $\dfrac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$.

Example     for $i=1$ to ⓜ}  ⁿᵇʳ ᵉˣᵃᵐᵖˡᵉˢ  $// x^{(1)}, y^{(1)} ; \ldots ; x^{(m)}, y^{(m)}$

|     Perform forward prop. and backprop using $x^{(i)}, y^{(i)}$.
|     ( Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l=2, \ldots, L$).

~~Batch~~ Batch

|     $\Delta^{(l)} \leftarrow \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$.  // accumulators having, in the end, the partial derivatives
}                                                            // $\dfrac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

compute $\dfrac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

⑤ Use gradient checking to compare $\dfrac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed in ④ with a numerical estimate of gradient of $J(\Theta)$.

⑥ Disable gradient checking.

⑦ Use Gradient Descent or other opt method with **backprop** to minimize $J(\Theta)$ as a function of params $\Theta$.

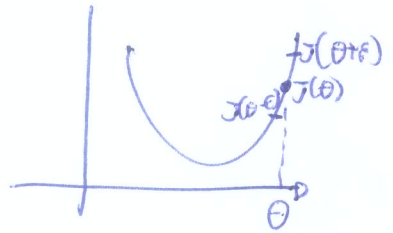{ As we have $J(\Theta)$ as a non-convex function, we may ~~have~~ converge to local minima.

Backprop → calculates the derivatives (direction of a step)
the gradient.

Gradient Descent → take little baby steps downhill.

————————————— x —————————————

## Last notes on Gradient Checking

- Always check the first results of Gradient Descent using gradient checking.



- Gradient checking is <u>too</u> slow. Do not use it for learning, just to check if it works
  in the first <u>iteration</u>.

————————————— x —————————————