

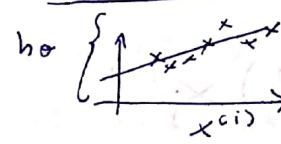
# Machine Learning

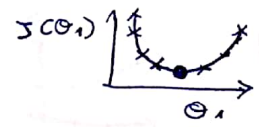
$m$  = number of training examples

$x$ s = features

$y$ s = target

## • Cost function ( $J_0$ )

$$h_0(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$


$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$


We have to minimize this

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0$$

## • Gradient descent algorithm

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1); \quad \alpha \equiv \text{learning rate}$$

Correct way

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

No. iterations

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

So que quito decir es  
que se actualiza todo a la  
...

• When we have multiple features:

$n$  = number of features

$x^{(i)}$   $\equiv$  input (features) of  $i$ th training example.

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

•  $n > 1$  (realments  $n \geq 1$ )

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

- simultaneously update  $\theta_j$  &  $\theta_0$

-  $x_0^{(i)} = 1$

### Mean normalization

We need to do this for all features

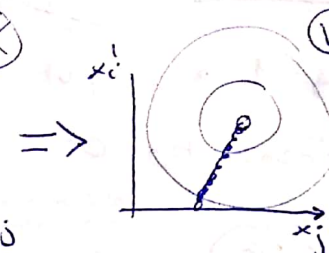
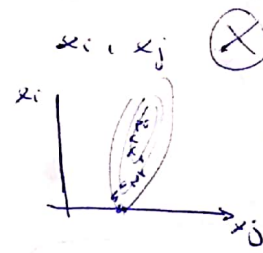
$$\forall i \text{ s.t. } i \neq 0 \Rightarrow \left\{ x_i = \frac{x_i - \mu_i}{s_i} \right\}$$

•  $i \neq 0$  because  $x_0 = 1$

•  $\mu_i$   $\equiv$  average value of  $x_i$  in training set

•  $s_i$   $\equiv$  range (max-min) or stdev of  $x_i$

### Feature scaling



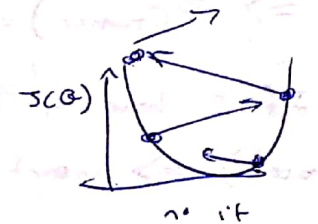
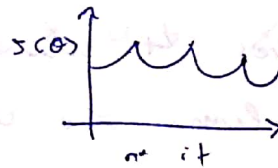
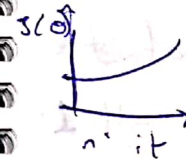
$$x_i' = \frac{x_i}{x_{i,max} - x_{i,min}}$$

$$x_j' = \frac{x_j}{x_{j,max} - x_{j,min}}$$

[This delays the convergence]

[Scaling by interval we accelerate the convergence]

• Plot  $J(\theta)$  vs number of iterations to check if it's working



This is WRONG. Try lower  $\alpha$

- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease every iteration

- If  $\alpha$  too small, reaching the convergence takes time.

## • Hypothesis representation



$h_{\theta}(x) \equiv \mathbb{P}$  that  $y=1$  on input  $x$   
parameterized by  $\theta$

$$h_{\theta}(x) = \mathbb{P}(y=1 | x; \theta)$$

$$\mathbb{P}(y=\square_1 | x; \theta) + \mathbb{P}(y=\square_2 | x; \theta) + \dots = 1$$

Ex if  $x = \begin{pmatrix} x_0 \\ x_1 \end{pmatrix} = \begin{pmatrix} 1 \\ \text{tumor size} \end{pmatrix}$  and  $h_{\theta}(x) = 0.7$   
 $\Rightarrow \mathbb{P}(\text{tumor}) = 70\%$

Ex If I have 4 features (4 columns of the dataset)  $\Rightarrow$  insert a column  $x_0$  with all 1:

dataset  $\rightarrow$

$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2	3	4	40
0	100	50	30	100

$\Rightarrow$

$x_0$	$x_1$	$x_2$	$x_3$	...
1	1	2	3	
1	0	100	50	

## • Normal equation

$$X^0 = \begin{pmatrix} 1 & 1 & 4 & 3 & 8 \\ 1 & 0 & 2 & 3 & 9 \end{pmatrix} \quad y = \begin{pmatrix} 40 \\ 100 \end{pmatrix}$$

$$\theta = \underbrace{(X^T X)^{-1}}_{\text{inversa de } X^T X} X^T \cdot y$$

$m \times (n+1)$   
 $m$ -dimensional vector  
When  $J(\theta)$  minimized

Esto ya está puesto como vector-columna.  
Los vectores fila están transpuestos

• En ~~del~~ Matlab  $X^+ = X^1$

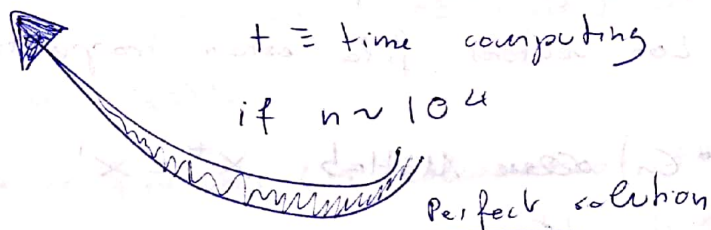


## Gradient descent (Numerical solution)

## Normal equation (Analytical solution)

- Feature scaling
- Needs to choose  $\alpha$
- Needs many iterations
- Works good when  $n$  large

- No need of feature scaling
- No need of  $\alpha$
- No need of iteration
- Slow if  $n$  very large  
 $t \sim n^3$   
 $t \equiv$  time computing  
if  $n \sim 10^4$



$$\Theta = (X^T X)^{-1} X^T Y$$

• If  $X^T X$  not invertible use pinv

$$\text{pinv}(X^T X) \cdot X^T \cdot Y$$

pinv  $\rightarrow$  pseudo invertible } use always  
inv  $\rightarrow$  invertible } pinv and ok

• If there is an invertibility problem check LD features.

there must be LD features, or if there are too many features, use regularization.

• Gradient descent  $\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$

If  $\alpha \ll \Rightarrow$  ~~gradient~~ ~~descent~~  $\nabla$  descent  
slow

If  $\alpha \gg \Rightarrow \nabla$  descent it may  
fail to converge

$$\text{If } \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$m \equiv$  batcher.

## MATLAB

• save ~~xxx~~ `□.mat` `f`; → guarda la función `f` en `□.mat`  
`.txt`

• load `□.mat` lo importa  
`.txt`

• Matrices:  $A \begin{smallmatrix} 3 \times 2 \end{smallmatrix} \Rightarrow A(3,2)$  es el elemento de esa posición

# Vector indexes in Matlab starts in 1, not in 0

Logistic regression: classification problem

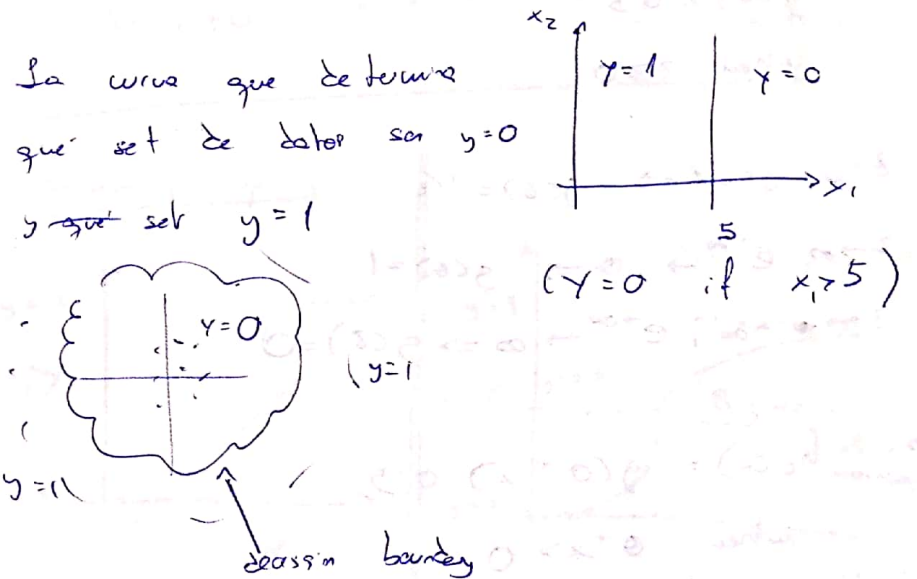
• Logistic regression  $\rightarrow h(\omega) = \sigma(\theta^T x)$

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

$$0 \leq h(\omega) \leq 1$$

• Decision boundary  $g(z)$ :  $z \equiv \theta^T x$

- If  $\theta_0 = 5, \theta_1 = -5, \theta_2 = 0, h(\omega) = g(5 - x_1)$



To discrete 0 or 1 classification

$$h_{\theta}(x) \geq 0.5 \rightarrow y=1$$

$$h_{\theta}(x) < 0.5 \rightarrow y=0$$

The logistic function works like

$$g(z) \geq 0.5$$

when  $z \geq 0$

$$z=0, e^0=1 \Rightarrow g(z)=1/2$$

$$z \rightarrow \infty, e^{\infty} \rightarrow \infty \Rightarrow g(z)=1$$

$$z \rightarrow -\infty, e^{-\infty} \rightarrow 0 \Rightarrow g(z)=0$$

$$\Rightarrow h_{\theta}(x) = g(\theta^T x) = 0.5$$

when  $\theta^T x \geq 0$

$$\Rightarrow \theta^T x \geq 0 \Rightarrow y=1$$

$$\theta^T x < 0 \Rightarrow y=0$$

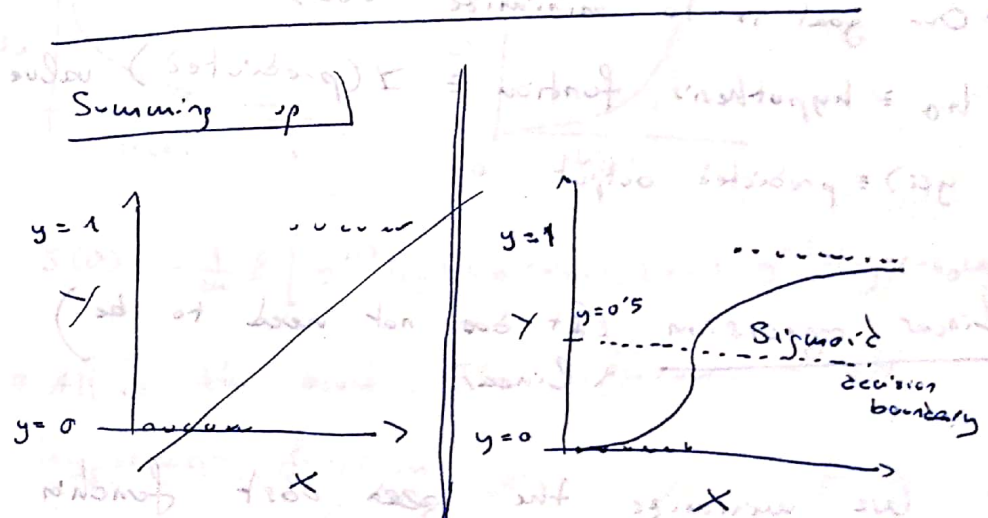
the decision boundary separates the area where  $y=0$  and where  $y=1$

$$\theta = \begin{pmatrix} 5 \\ -1 \\ 0 \end{pmatrix}$$

$$y=1 \Leftrightarrow 5 + (-1)x_1 + 0x_2 \geq 0$$

$$+5 - x_1 \geq 0 \quad ; \quad -x_1 \geq -5 \quad ; \quad x_1 \leq 5$$

boundary decision



Predicted  $y$  can exceed 0 and 1

Predicted  $y$  lies within 0 and 1

$h_{\theta} \equiv$  hypothesis function

$$h_{\theta}(x) = \beta_0 + \beta_1 x$$

$$0 \leq h_{\theta}(x) \leq 1$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$z(x) = \theta^T x$



## Cost function

Cost function  $\equiv$  error function

$$\text{Error} = \sum (\text{actual output} - \text{predicted output})^2$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

- Our goal is to minimize  $J(\theta)$
- $h_{\theta}$   $\equiv$  hypothesis function  $\equiv$   $\hat{y}$  (predicted) value
- $y^{(i)}$   $\equiv$  predicted output

## Linear regression (It does not need to be linear)

We minimize the ~~cost~~ cost function using gradient descent

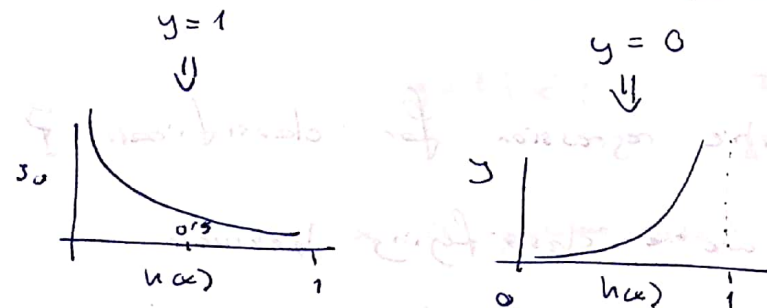
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$\alpha$  is learning rate. Be careful

Not too big (no converge), not too small (it takes a lot)

## Logistic regression (classification)

$$J(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y=0 \end{cases}$$



$$J(\theta) = -\frac{1}{m} \sum \left[ y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

$\neq$  All is the same. The difference is the hypothesis function

$$J(h_{\theta}(x), y) = \begin{cases} 0 & \text{if } h_{\theta}(x) = y \\ \infty & \text{if } y=0 \text{ and } h_{\theta}(x) \rightarrow 1 \\ \infty & \text{if } y=1 \text{ and } h_{\theta}(x) \rightarrow 0 \end{cases}$$

• Using gradient descent for classification:

$$\frac{d}{dx} \sigma(x) = \frac{d}{dx} \left( \frac{1}{1+e^{-x}} \right) = \frac{-(1+e^{-x})^{-1}}{(1+e^{-x})^2} = \dots = \sigma(x)(1-\sigma(x))$$

# Why logistic regression for classification?

Suppose we're classifying between spam or not spam. In linear regression  $h(x)$  can be  $>1$  or  $<1$ , which is out of context of probabilistic probability.

With logistic regression,  $h(x) \in (0, 1)$

and we can assign a probability

Using the Carsen nomenclature:

$$h_{\theta}(x) = g(\theta^T x) \quad g(z) = \begin{cases} \rightarrow 1 & \text{as } z \rightarrow \infty \\ \rightarrow 0 & \text{as } z \rightarrow -\infty \\ \in (0, 1) \end{cases}$$

$$z = \theta^T x$$

$$g(z) = \frac{1}{1+e^{-z}} \equiv \text{sigmoid function}$$

$h_{\theta}(x) = \text{IP}(y=1 | x; \theta)$  The probability of  $y=1$  when  $x$  is parameterized to  $\theta$ .

$$\text{If } h_{\theta}(x) \geq 0.5 \Rightarrow y=1$$

$$h_{\theta}(x) < 0.5 \Rightarrow y=0$$

$$g(z) \geq 0.5 \Rightarrow \underbrace{\theta^T x}_{\geq 0.5} \Rightarrow \underbrace{z}_{\geq 0.5}$$

This is the decision boundary

The decision boundary is the line/curve that distinguishes the area where  $y=0$  and  $y=1$



## # Understanding logistic regression:

The setting of the threshold is a very important aspect of the Logistic Regression, and is dependent on the classification problem itself.

The decision for the value of the threshold is majority affected by precision and recall. Ideally, we want a precision and recall to be 1.

### • Low precision / high recall

When we want to reduce the number of false negatives without reducing the number of false positives

Ex Cancer diagnosis. We don't want any affected patient to be classified as not affected. Without giving much heed to if the patient is being wrongfully

diagnosed with cancer. This is because the absence of cancer can be detected in further tests but no further tests cancer y ser diagnosed since (perque no harm mai praveba y te monitoris).

### • High precision / low recall

When we want to reduce the number of false positives without reducing the number of false negatives.

Ex Classifying customers whether they will react positively or negatively to a personalized ad. We want to make sure they'll react positively, otherwise we can miss loss a customer.

Then, logistic regression classification:

- binomial
- multinomial

- ordinal

very poor	→ 1
poor	→ 2
good	→ 3
very good	→ 4

• Multiple features,  $n \equiv \text{features}$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$x_0^{(i)} \equiv 1$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\Rightarrow h_{\theta}(x) = \theta^T X$$

for  $n=1$

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

update  
simultaneously

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

For  $n > 1$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Simultaneously update  $\theta_j$  &  $\theta_0$

• Making sure gradient descent is working correctly

- For sufficient small  $\alpha$ ,  $J(\theta)$  decrease over every iteration

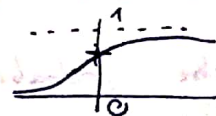
- It is not slow converging

• Octave tutorial → Lecture 5

Again, Logistic Regression is for classification,

where  $0 \leq h_{\theta} \leq 1$  thanks to

$$h_{\theta}(x) \equiv \frac{1}{1 + e^{-\theta^T x}}$$



tangente hiperbólica de tan la vida cabese (con altura)

• Logistic regression cost function

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$$\Rightarrow J(\theta) = \frac{1}{n} \sum_{i=1}^n \text{cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{n} \left[ \sum_{i=1}^n y^{(i)} \log(h_{\theta}(x^{(i)})) + \sum_{i=1}^n (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right]$$

• Gradient descent

$$\theta_{j,i} = \theta_{j,i} - \alpha \frac{\partial}{\partial \theta_{j,i}} J(\theta)$$

• Multiclass

• One-vs-all

Train a logistic regression classifier  $h_{\theta}^{(i)}(x)$  for each class  $i$  to predict the probability that  $y=i$

On a new input  $x$  to make prediction, pick the class that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

• The problem of overfitting

- too many features  
↳ model selection algorithm

- regularization

↳ keep all features, but reduce magnitude/values of parameters  $\theta_j$   
↳ works well when having a lot of features