

★ Topics Covered :

- ① Supervised vs Unsupervised learning
- ② Linear regression : Model & Cost Function
- ③ Gradient Descent Algorithm
- ④ Implementing Squared Error Gradient Descent for Linear regression.
- ⑤ Classification using Logistic regression
- ⑥ Cost function for log. regression
- ⑦ Implementation with Gradient descent

Notes posted in the classroom are the primary material.
 This is just a supplement to understand it better. Approaches in both may differ

SOURCE:

- Coursera :
 Supervised ML by
 Andrew Ng.

- I Supervised: → Algorithms that learn to map input x to o/p y .
- give learning Algo. examples to learn from.
 - learns by seeing correct pairs of i/p x & desired o/p label y and gives a reasonably accurate prediction or guess of o/p.

I/p :

e-mail

English

ad, user info

o/p
Spam (0/1)?

Spanish.

Click? (0/1)

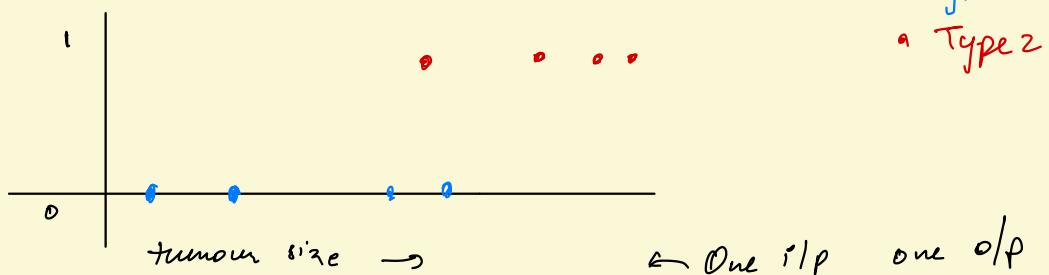
Instructor:

Andrew Ng

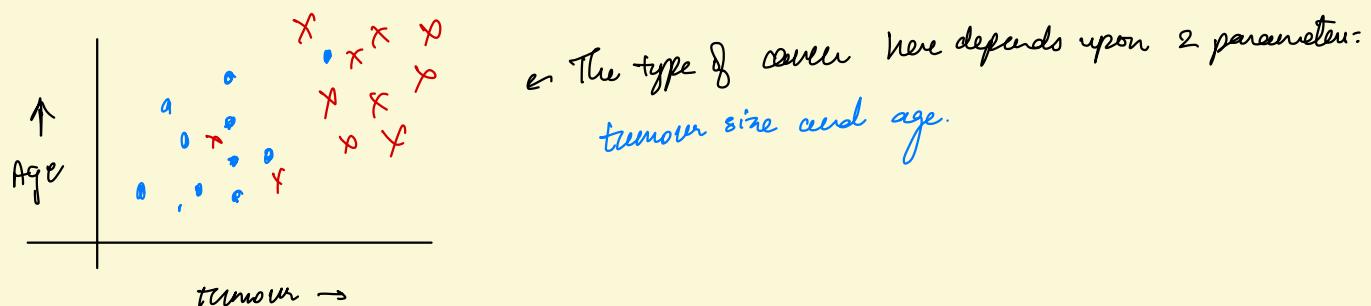
Type of Sup. learning:

- ① Regression: given set of data, predict o/p for a new data
- aim: how to get an algo. to systematically choose the most appropriate line / curve / other thing to fit this data
- Eg Given the house area vs the cost predict the cost for a given area.

② Classification:



• Type 1 cancer
• Type 2 cancer

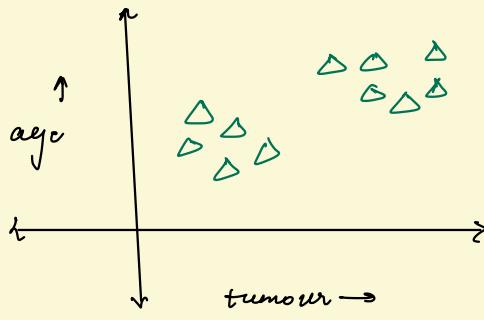


II Unsupervised learning:

We are given data that is not associated with any o/p label y .

It's unsupervised because we're not trying to supervise the algorithm to give some right answer for every P/p.

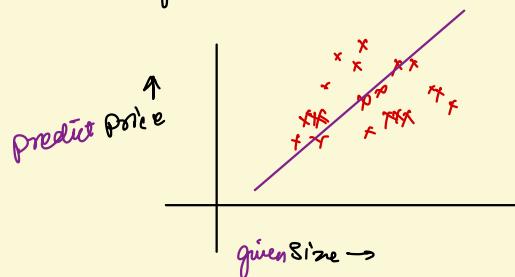
Supervised learning → data labeled coffee right answers



- ① Find something interesting in unlabeled data
like clusters or outliers
- Eg: ① Google news: groups similar stories together
② Grouping customers.

② Anomaly detection finds unusual data points

Linear Regression:



Predict the price of house based on size.

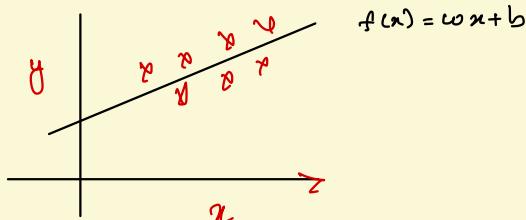
Supervised learning:
Data has "right" answers

$m = \# \text{ of training examples.}$

$(x_i, y_i) = \text{single training example.}$

$$f_{w,b}(x) = w_0 x + b$$

choose w & b to predict y .



Supervised learning

① Classification model.
predicts categories.

- ① reg. model predicts number.
② use # of possible ops

- ② small no. of possible ops

training set

learning algo

$x \rightarrow f \rightarrow \hat{y} \rightarrow \text{Estimated } y$
"model"

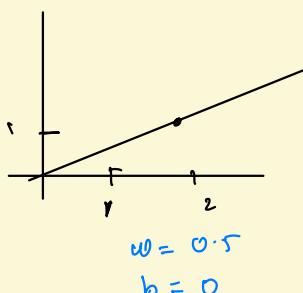
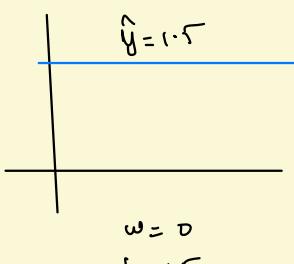
How to represent f ?

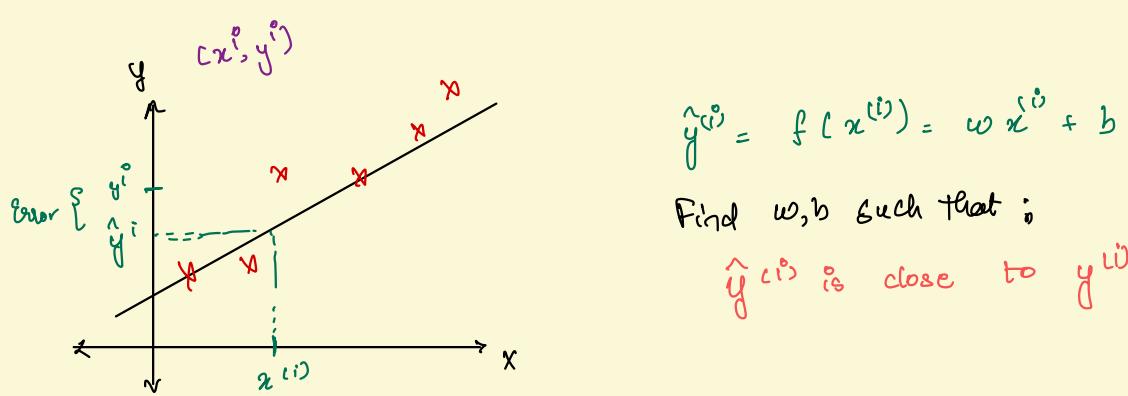
- ① Linear regression with one variable size / feature x
"Univariate" linear regression

* Cost function:

$$\text{Model: } f_{w,b}(x) = w_0 x + b$$

w, b : parameters





Cost function: $\sum (y - \hat{y})^2 \rightarrow \text{error}$ $m = \# \text{ of training examples}$

$$\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 \rightarrow \text{gets larger & larger as } m \text{ increases}$$

error

We take the average hence.

* Cost function: "Squared Error Cost Function"

$$\begin{aligned} J(w, b) &= \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 \\ \therefore J(w, b) &= \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 \end{aligned}$$

Goal: $\min_{w, b} J(w, b)$

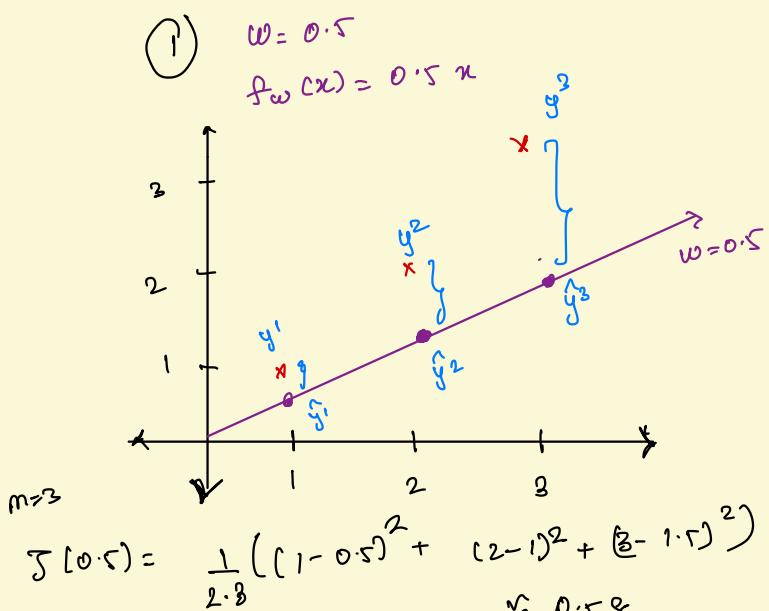
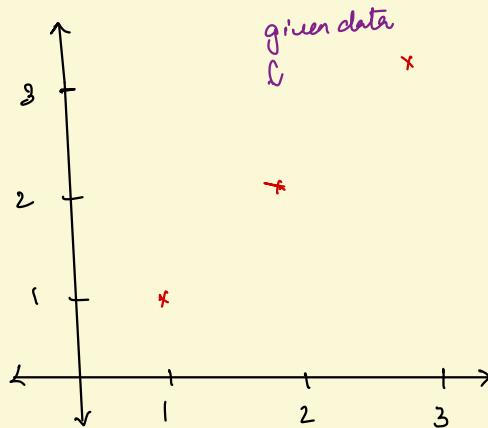
↳ How are the cost functions and the model related? $\text{let } b=0$

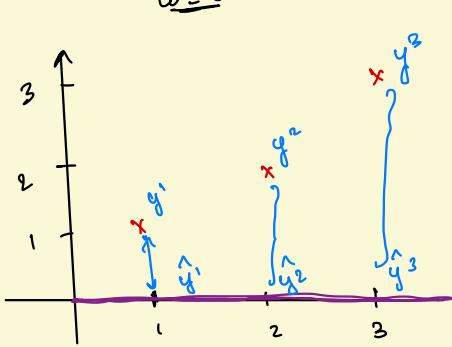
$$f_w(x) = w x$$

\Rightarrow depends upon w if x

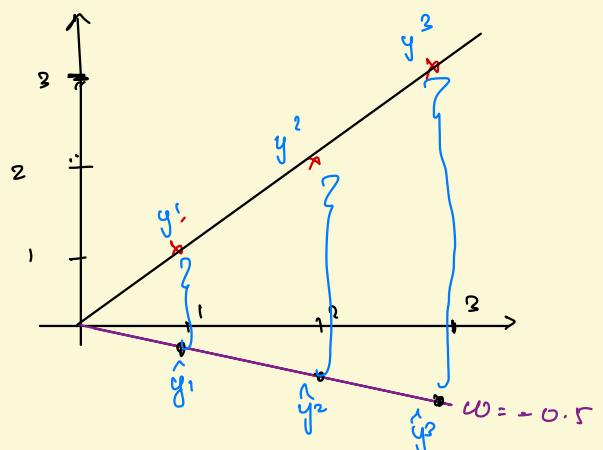
$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

\Rightarrow Depends upon w parameter





$$J(w) = \frac{1}{2m} [1^2 + 2^2 + 3^2] = 2 \cdot 3$$



$\text{eg } w=1 \Rightarrow \hat{y} = w\hat{x} = \hat{x}$

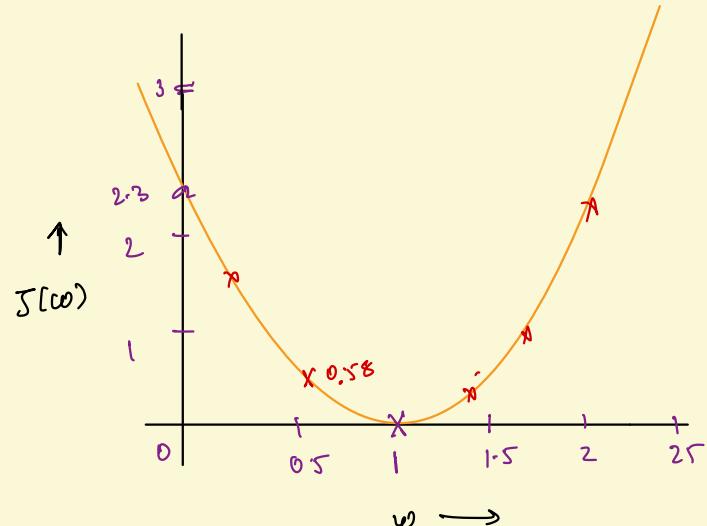
$$J(w) = \frac{1}{2m} \sum_{i=1}^m (w x^i - y^i)^2 = \frac{1}{2m} \sum_{i=1}^m (x^i - y^i)^2 = 0$$

$$w=1 \Rightarrow J(w)=0$$

$\text{eg } w=0.5 \quad \hat{y} = 0.5\hat{x}$

$$J(w) = \frac{1}{2m} \sum_{i=1}^m (0.5 x^i - y^i)^2$$

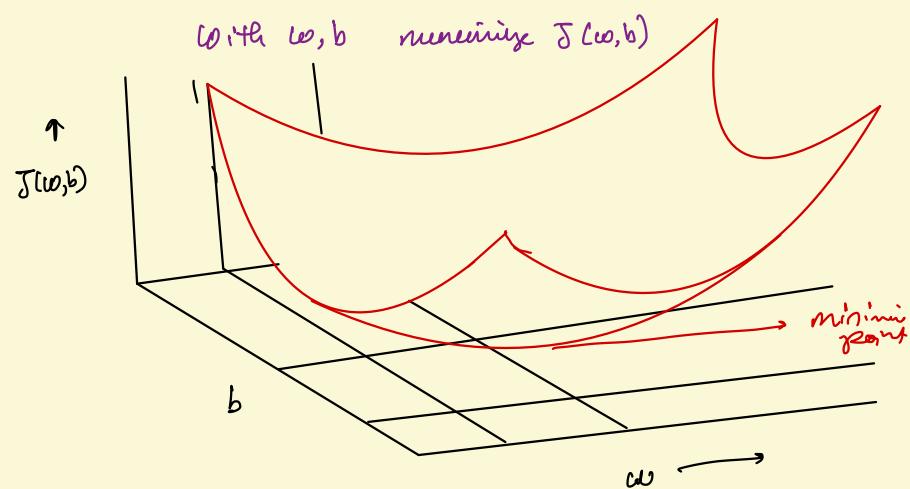
\therefore we want to minimize $J(w)$
find w



General case: Model: $f_{w,b}(x) = w x + b$

$J(w,b)$

Need an algorithm for finding parameters w & b that give the best fit for $J(w,b)$



* Gradient Descent.

Have some function $J(w, b)$

want $\min_{w, b} J(w, b)$

$\min_{w_1, w_2, \dots, w_n, b} J(w_1, w_2, \dots, w_n, b)$

Cost function for linear regression: squared error.

Cost func for neural network: may be differ

Cost func:

Error function.

* GRADIENT DESCENT ALGORITHM:

Repeat until convergence: \rightarrow until u reach a local minimum

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b) \quad \text{old } w \text{ only.}$$

α = learning rate

simultaneously update w & b

Correct:

Simultaneous update

$$\text{tmp_}w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$\text{tmp_}b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$w = \text{tmp_}w$$

$$b = \text{tmp_}b$$

Incorrect:

Simultaneous Update

$$\text{tmp_}w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

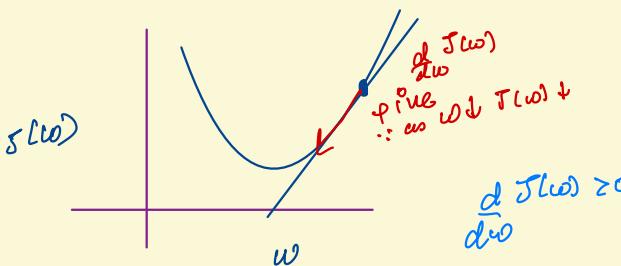
$$w = \text{tmp_}w$$

$$\text{tmp_}b = b - \alpha \frac{\partial}{\partial b} J(w, b) \quad \begin{matrix} \text{new } w \\ \times \end{matrix}$$

$$b = \text{tmp_}b \quad \text{wrong}$$

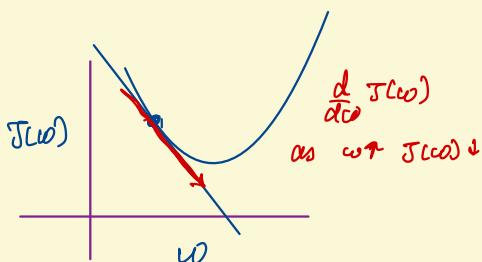
α = learning rate how big a step u take in updating w & b

Intuition when $b=0$:



$$w = w - \alpha \frac{d}{dw} J(w)$$

$$d > 0 \quad \underbrace{w = w - \alpha}_{w \downarrow} \quad \text{(positive number)}$$



$$w = w - \alpha \frac{d}{dw} J(w)$$

$$w = w - \alpha \quad [-\text{ive no.}]$$

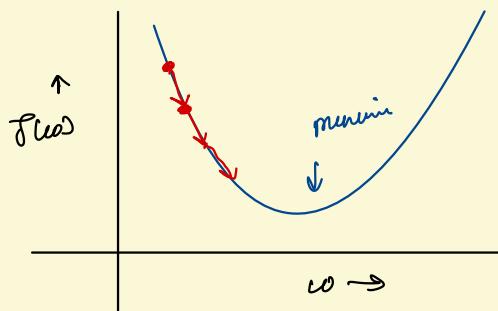
$$w = w + \alpha \quad [\text{+ive no.}] \quad w \uparrow$$

↳ Learning rate:

$$w = w - \alpha \frac{d}{dw} J(w)$$

If α is too small ...

Gradient descent may be slow



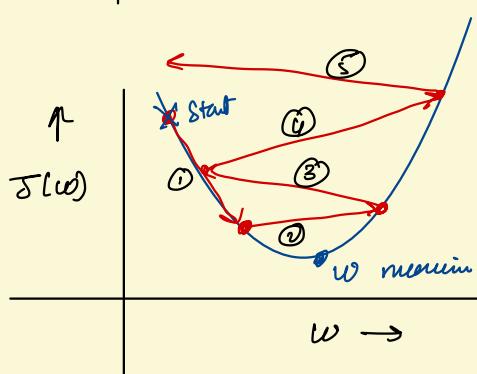
$$w = w + \alpha \frac{d}{dw} J(w)$$

If α is too large ...

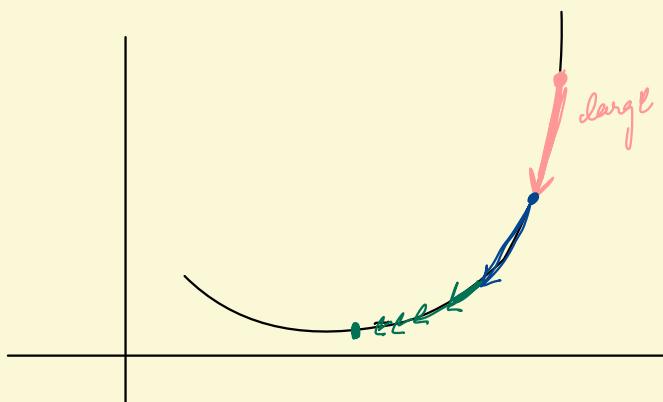
Gradient descent may,

overshoot, never reach min.

→ may diverge



At absolute minimum w , the algo does nothing



↳ Linear regression model with Gradient Descent algo:

linear reg. model

$$f(x) = wx + b$$

Cost function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^i) - y^i)^2$$

Gradient descent algorithm:

repeat until convergence Σ

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$J = \frac{1}{2m} \sum_{i=1}^m (w^T x^i + b - y^i)^2$$

?

$$\frac{\partial}{\partial w} J(w, b) = \frac{1}{2m} \sum_{i=1}^m 2 [w^T x^i + b - y^i] x^i = \sum_{i=1}^m \frac{1}{m} (w^T x^i + b - y^i) x^i$$

$$\frac{\partial}{\partial b} J(w, b) = \frac{1}{2m} \sum_{i=1}^m 2 [w^T x^i + b - y^i] = \sum_{i=1}^m \frac{1}{m} (w^T x^i + b - y^i)$$

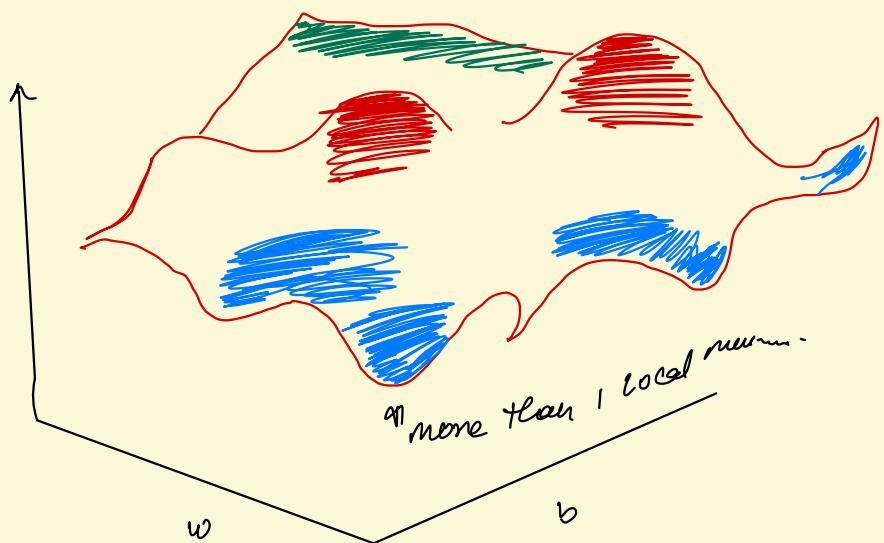
\therefore Gradient descent algorithm:

$$w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f(x^i) - y^i) x^i$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f(x^i) - y^i)$$

What if the cost func. has more than one local minima?

when squared error
cost func. with linear
reg. cost function $J(w, b)$
doesn't & will not
have multiple local minima
It's a convex function
has single global minimum



\therefore As long as α is chosen appropriately
grad. desc. always converges.

Batch Gradient descent

"Batch": Each step of gradient descent uses all training e.g. (x^p, y^p)

"Stochastic G.D": Each step of grad. descent uses a sample.

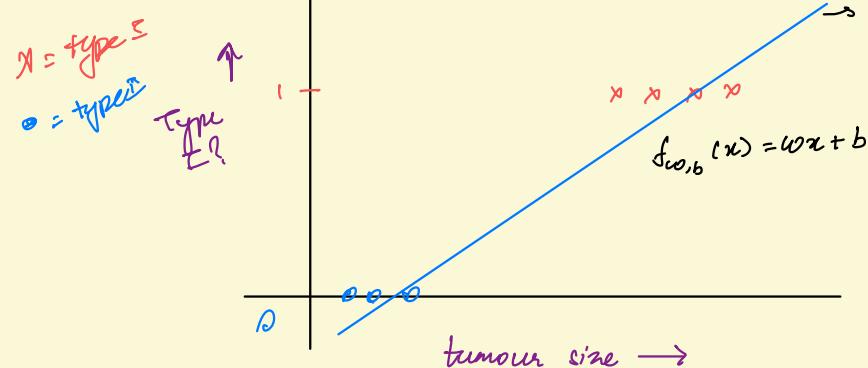
→ Classification: dep variable y can take up only a small handful values instead of range of numbers

Linear regression is not a good algorithm for classification.

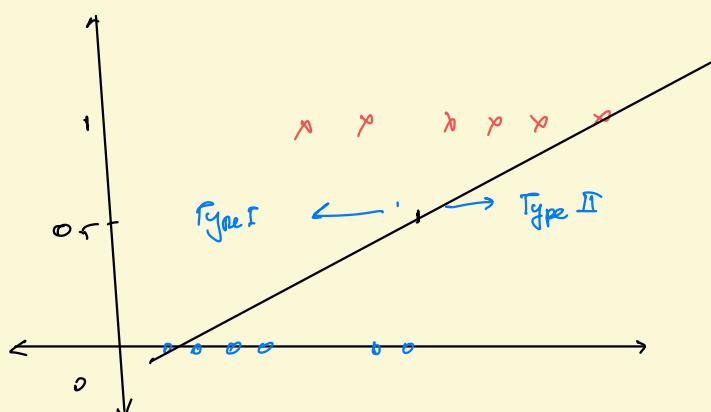
Is this email spam?
Is it a virus? Found?
Type I / Type II cancer?

no / yes
no / yes
no / yes

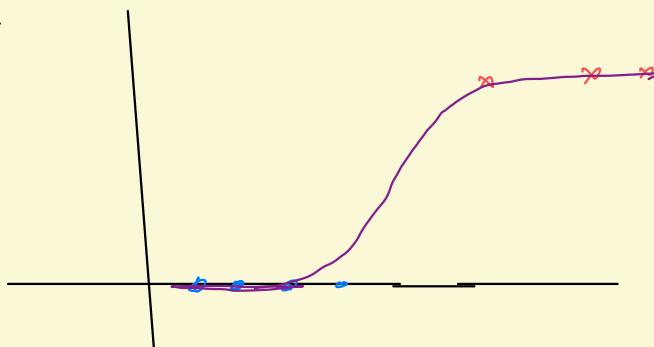
y can only be one of two values
"binary" classification



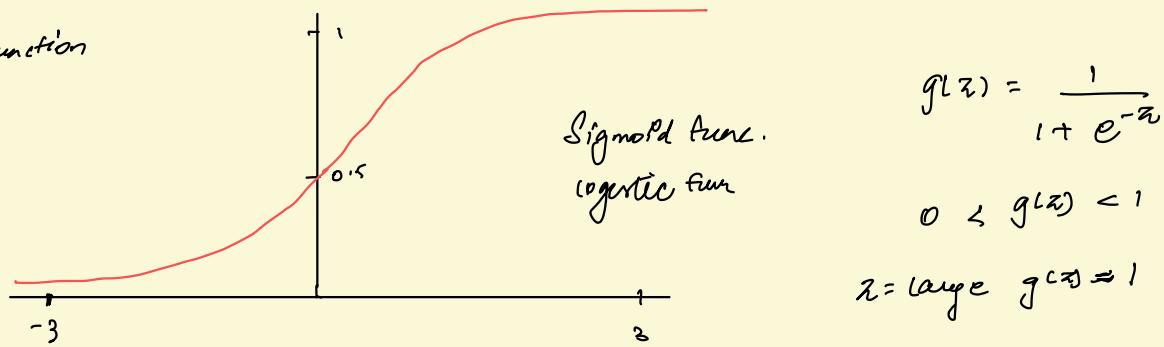
Prob.



logistic regression



Sigmoid Function



when $x = -\vec{w} \cdot \vec{x}$, $g(x) = \frac{1}{1 + e^{-(-\vec{w} \cdot \vec{x})}} \approx 0$
large no.

$$f_{w,b}(x) = \vec{w} \cdot \vec{x} + b$$

$$\pi = \vec{w} \cdot \vec{x} + b$$

$$g(x) = \frac{1}{1 + e^{-x}}$$

∴ "Logistic regression" model:

$$f_{w,b}(x) = g(\vec{w} \cdot \vec{x} + b) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

$$f_{w,b}(x) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}} = \text{"prob." that class is 1}$$

Example: \vec{x} = tumour size

y = type I or not

If $f_{w,b}(\vec{x}) = 0.7$ Model thinks that there is 70% chance that the true label $y = 1$

↳ Decision Boundary:

Logistic regression is: Take linear reg. model & sigmoid func to it.

$$f_{w,b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = \pi$$

$$\begin{matrix} \downarrow \\ z \\ \downarrow \end{matrix}$$

$$g(x) = \frac{1}{1 + e^{-x}}$$

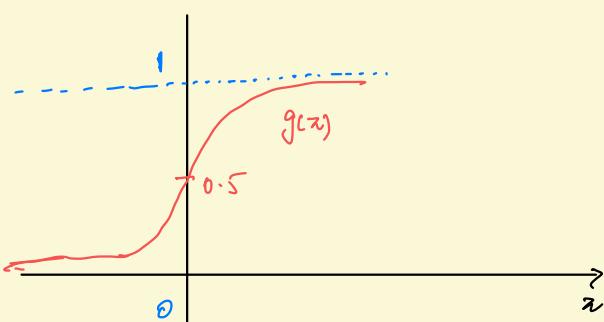
$$\therefore \text{Model P.S.: } f_{w,b}(x) = g(\vec{w} \cdot \vec{x} + b) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

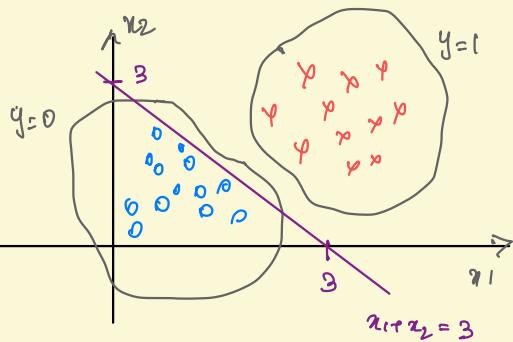
The model predicts $\hat{y} = 1$ whenever:

- $g(x) \geq 0.5$ if $x \geq 0$

- $x \geq 0$ if $\vec{w} \cdot \vec{x} + b \geq 0$

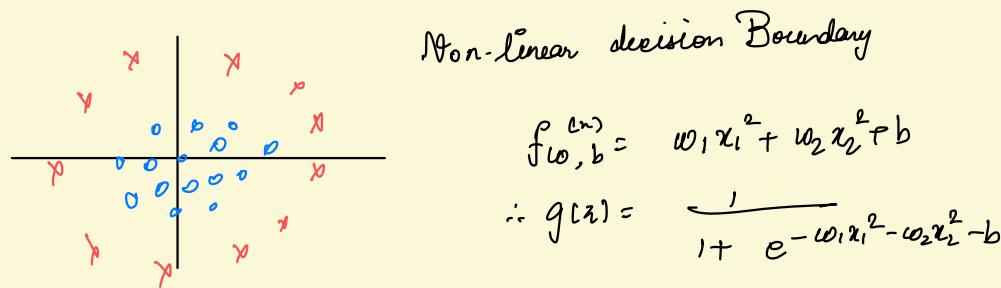
$\therefore g(x) \geq 0$ whenever $\vec{w} \cdot \vec{x} + b \geq 0$.





$$\begin{aligned}
 f_{w,b}(\vec{x}) &= g(\vec{x}) = g(\vec{w} \cdot \vec{x} + b) \\
 &= g(w_1 x_1 + w_2 x_2 + b) \quad \text{let } b=3 \\
 \vec{x} &= \vec{w} \cdot \vec{x} + b = 0 \rightarrow \text{Decision Boundary} \\
 z &= x_1 + x_2 - 3 = 0 \\
 x_1 + x_2 &= 3
 \end{aligned}$$

Here, Decision Boundary is a straight line.



$$\begin{aligned}
 f_{w,b}^{(lin)} &= w_1 x_1^2 + w_2 x_2^2 + b \\
 \therefore g(z) &= \frac{1}{1 + e^{-w_1 x_1^2 - w_2 x_2^2 - b}}
 \end{aligned}$$

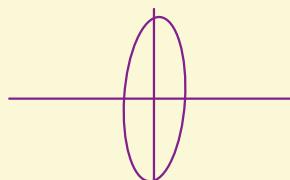
$$w_1 = w_2 = 1 \quad z = w_1 x_1^2 + w_2 x_2^2 + b, \quad w_1 = w_2 = 1 \quad b = -1$$

Type I (red) when $g(z) \geq 0.5$ when $z \geq 0$ $z = x_1^2 + x_2^2 - 1 \geq 0$ If $x_1^2 + x_2^2 \geq 1$ it's type I

$$\text{Bound: } x_1^2 + x_2^2 = 1$$

$$\text{If } f_{w,b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2 + \dots + b)$$

$$\text{If } f_{w,b}(\vec{x}) = g(w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2)$$



IMP

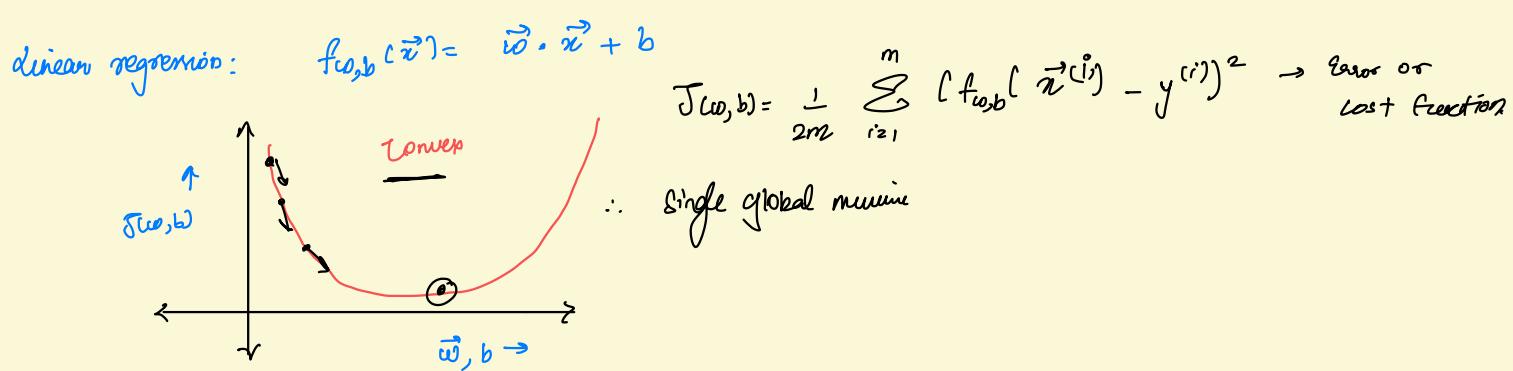
- ① For logistic regression, choose a model (linear or poly) & feed it into the sigmoid function to observe the op.
- ② In order to choose a model, we need to choose parameters w, b by minimizing a cost function.

Cost function for logistic regression model:

$$f_{w,b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

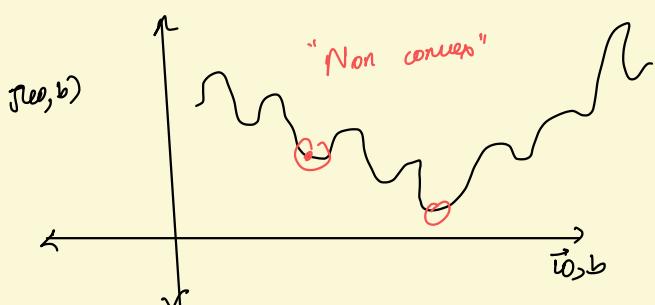
$$\therefore f_{w,b}(\vec{x}) = g(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

Choose w, b



If we use the same cost fun. for $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$, logistic regression fun.

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m \left(\frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}} - y^{(i)} \right)^2 \rightarrow \text{this cost fun. is non-convex}$$



\therefore Squared error cost function is not good enough.

\therefore We choose a different cost fun. that makes the cost func. convex again

\therefore Gradient descent can be guaranteed to converge to the global minimum

* Building a new cost function:

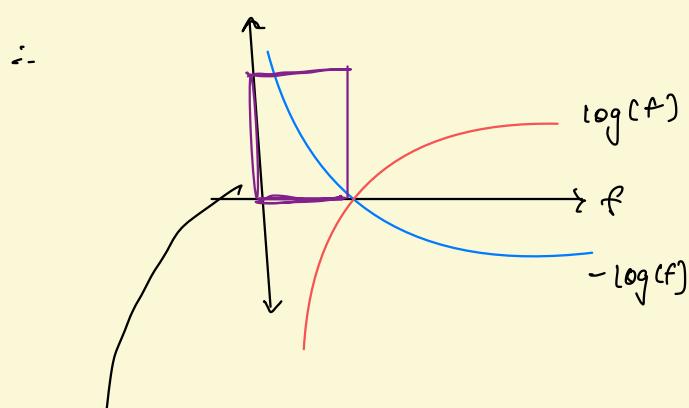
* Loss function in a cost fun.: $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$

\rightarrow Predictor o/p fun.
loss function: $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$
 \downarrow true label y .

then we give $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & , y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & , y^{(i)} = 0 \end{cases}$$

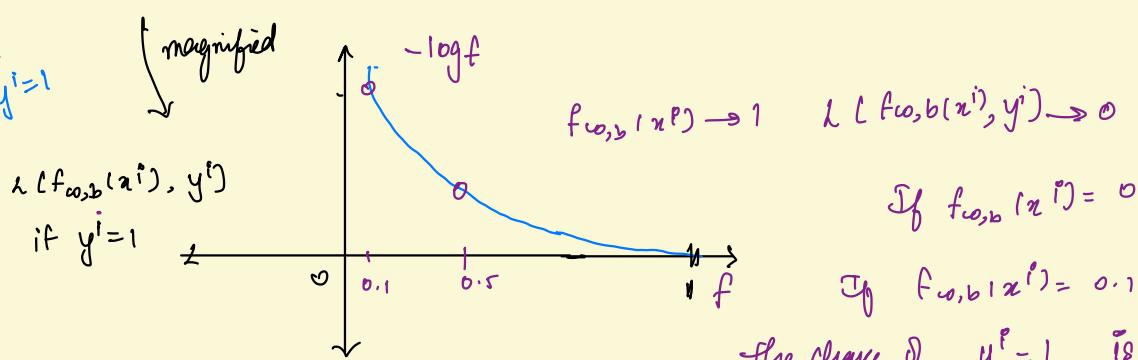
Observe $f_{\vec{w}, b}(\vec{x}^{(i)})$ is always between 0 and 1



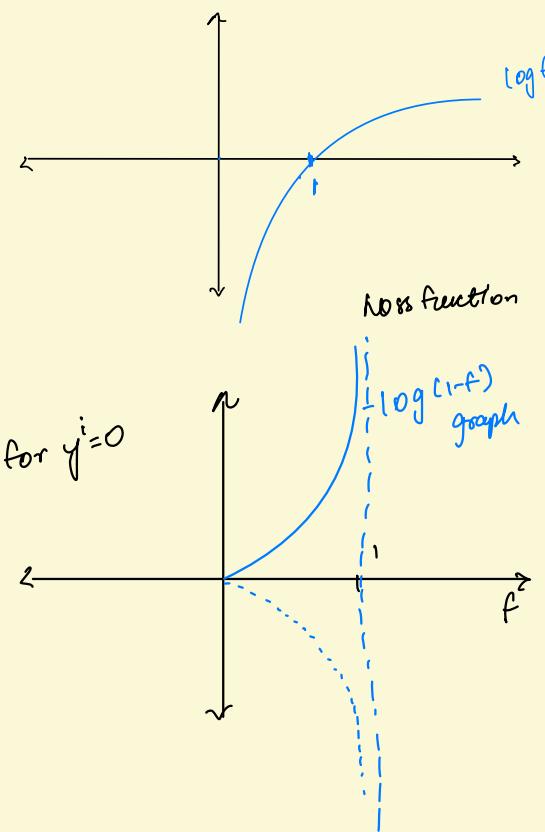
\therefore If algo. predicts $f_{\vec{w}, b}(\vec{x}^{(i)})$ close to 1 & true label is 1,

For logistic reg. o/p is close to 1

Case I
label $y^i = 1$

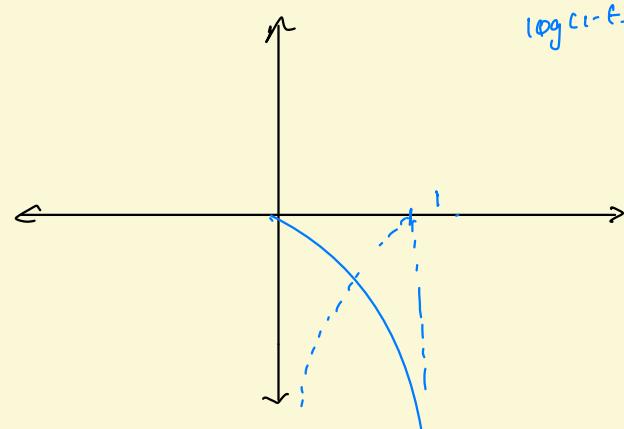


Case II
label $y^i = 0$



∴ When $f = 0$, or very close to zero, loss func $\rightarrow 0$
i.e. true label is zero & predictor label = 0.
When $f \rightarrow 1$, loss func $\rightarrow \infty$.

for $0 < f < 1$,
 $\log(1-f)$ graph looks like this



→ logistic squared error not a good choice.

$$\cdot \lambda [f_{w,b}(\vec{x}^i), y^i] = \begin{cases} -\log(f_{w,b}(\vec{x}^i)) & , y^i = 1 \\ -\log(1 - f_{w,b}(\vec{x}^i)) & , y^i = 0 \end{cases}$$

Also, this choice of loss func. is convex.

$$\therefore J(w, b) = \frac{1}{m} \sum_{i=1}^m \lambda [f_{w,b}(\vec{x}^i), y^i]$$

~~$\ln \left(\frac{1}{1 + e^{-w_i x_i}} \right)$~~
 $\ln \left(\frac{1}{1 + e^{-w_i x_i}} \right)$

⇒ Simplified cost func. for Logistic Regression:

$$\lambda [f_{w,b}(\vec{x}^i), y^i] = \begin{cases} -\log(f_{w,b}(\vec{x}^i)) & , y^i = 1 \\ -\log(1 - f_{w,b}(\vec{x}^i)) & , y^i = 0 \end{cases}$$

i.e. loss function:

$$h(f_{\omega, b}(x^i), y^i) = -y^i \log(f_{\omega, b}(x^i)) - (1-y^i) \log(1-f_{\omega, b}(x^i))$$

$$\therefore J(\vec{\omega}, b) = \frac{1}{m} \sum_{i=1}^m h(f_{\omega, b}(x^i), y^i) = \frac{1}{m} \sum_{i=1}^m (y^i \log(f_{\omega, b}(x^i)) + (1-y^i) \log(1-f_{\omega, b}(x^i)))$$

This particular cost func. is derived from statistics

using Max Likelihood estimator.

* Gradient Descent: to find the min of $J(\vec{\omega}, b)$

Cost: $J(\vec{\omega}, b) = -\frac{1}{m} \sum_{i=1}^m (y^i \log(f_{\omega, b}(x^i)) + (1-y^i) \log(1-f_{\omega, b}(x^i)))$

$$f(\vec{\omega}, b) = \frac{1}{1 + e^{-(\vec{\omega} \cdot \vec{x} + b)}}$$

repeat

$$\left\{ \begin{array}{l} \omega_j = \omega_j - \alpha \frac{\partial}{\partial \omega_j} J(\vec{\omega}, b) \quad j = 1, \dots, n \end{array} \right.$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{\omega}, b) \quad \left\{ \text{simultaneous updates} \right.$$

$$\frac{\partial}{\partial \omega_j} J(\vec{\omega}, b) = -\frac{1}{m} \sum_{i=1}^m \frac{y_i}{f_{\omega, b}(x^i)} \left(\frac{1}{1 + e^{-(\vec{\omega} \cdot \vec{x} + b)}} \right)^2 e^{-(\vec{\omega} \cdot \vec{x} + b)} \\ - \frac{(1-y^i)}{1-f_{\omega, b}(x^i)} \left(\frac{1}{1 + e^{-(\vec{\omega} \cdot \vec{x} + b)}} \right)^2 e^{-(\vec{\omega} \cdot \vec{x} + b)}$$

$$\left(\begin{array}{l} \left(\frac{y_i}{f} - \frac{1-y^i}{1-f} \right) = \frac{y_i - y_i f - 1 + f}{f(1-f)} = \frac{y_i - f}{f(1-f)} \\ f(1-f) = \left(\frac{1}{1 + e^{-(\vec{\omega} \cdot \vec{x} + b)}} \right) \left(1 - \frac{1}{1 + e^{-(\vec{\omega} \cdot \vec{x} + b)}} \right) = \left(\frac{1}{1 + e^{-(\vec{\omega} \cdot \vec{x} + b)}} \right)^2 e^{-(\vec{\omega} \cdot \vec{x} + b)} \end{array} \right)$$

$$\frac{\partial}{\partial \omega_j} J(\vec{\omega}, b) = -\frac{1}{m} \sum_{i=1}^m \left(\frac{y_i}{f} - \frac{1-y^i}{1-f} \right) \left(-y_j f (1-f) \right)$$

$$= -\frac{1}{m} \sum_{i=1}^m (y_i - f) y_j f = -\frac{1}{m} \sum_{i=1}^m (y_i - f) y_j = \frac{1}{m} \sum_{i=1}^m (f - y_i) y_j$$

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m (y^i \log(f_{w,b}(x^i)) + (1-y^i) \log(1-f_{w,b}(x^i)))$$

$$f(w,b) = \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}}$$

$$\frac{\partial}{\partial b} J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m (y^i \left[\frac{1}{f} \right] \left[\frac{e^{-(\vec{w} \cdot \vec{x} + b)}}{(1+e^{-(\vec{w} \cdot \vec{x} + b)})^2} \right] - \frac{(1-y^i)}{1-f} \left[\frac{e^{-(\vec{w} \cdot \vec{x} + b)}}{(1+e^{-(\vec{w} \cdot \vec{x} + b)})^2} \right])$$

$$\left(\frac{y^i}{f} - \frac{(1-y^i)}{1-f} \right) = \frac{y^i(1-f) - f(1-y^i)}{f(1-f)} = \frac{y^i - f}{f(1-f)}$$

$$f(1-f) = \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}} \left(1 - \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}} \right) = \frac{\frac{e^{-(\vec{w} \cdot \vec{x} + b)}}{1+e^{-(\vec{w} \cdot \vec{x} + b)}}}{\left(\frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}} \right)}$$

$$\frac{\partial}{\partial w} J(w, b) = -\frac{1}{m} \sum_{i=1}^m \left(\frac{y^i}{f} - \frac{(1-y^i)}{1-f} \right) f'(1-f) = \frac{1}{m} \sum_{i=1}^m (y^i - f_{w,b}(x^i))$$

∴

Gradient descent algorithm:

repeat

do

$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^i) - y^i) x_j^i \right]$$

$$b = b - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^i) - y^i) \right]$$

} simultaneous updates

↓
with old w only

Linear Regression using Closed form.

CLOSED FORM:

loss func. : $J(\omega, b) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - f_{\omega, b}(x^{(i)}))^2$

$$f_{\omega, b}(x^{(i)}) = \omega_1 x_1 + \omega_2 x_2 + b \cdot 1$$

$$\vec{\omega} = \begin{bmatrix} \omega_1 \\ \omega_2 \\ b \end{bmatrix} \quad \vec{x}_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ 1 \end{bmatrix} \quad \therefore f_{\omega, b}(x^{(i)}) = [\omega_1 \ \omega_2 \ b] \begin{bmatrix} x_{1i} \\ x_{2i} \\ 1 \end{bmatrix}$$

$$f_{\omega, b}(x^{(i)}) = \vec{\omega}^T \vec{x}_i$$

$$J(\omega, b) = J(\vec{\omega})$$

$$J(\vec{\omega}) = \frac{1}{2m} \sum_{i=1}^m (y^{(i)} - \vec{\omega} \cdot \vec{x}_i)^2$$

Minimize $J(\vec{\omega})$: $\frac{\partial J}{\partial \vec{\omega}} = 0$

$$\frac{\partial J}{\partial \vec{\omega}} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \vec{\omega} \cdot \vec{x}_i) (-\vec{x}_i) = 0$$

$$\frac{\partial J}{\partial \vec{\omega}} = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \vec{x}_i + \frac{1}{m} \sum_{i=1}^m (\vec{\omega} \cdot \vec{x}_i) \vec{x}_i = 0$$

$$(\vec{\omega} \cdot \vec{x}_i) \vec{x}_i = (\vec{x}_i \vec{x}_i^\top) \vec{\omega}$$

$$\sum_{i=1}^m (\vec{\omega} \cdot \vec{x}_i) \vec{x}_i = \sum_{i=1}^m (\vec{x}_i \vec{x}_i^\top) \vec{\omega} = 0$$

$$\frac{\partial J}{\partial \vec{\omega}} = -\sum_{i=1}^m y^{(i)} \vec{x}_i \rightarrow \left(\sum_{i=1}^m (\vec{x}_i \vec{x}_i^\top) \right) \vec{\omega} = 0$$

For $m=4$:

$$\sum_{i=1}^4 y^{(i)} \vec{x}_i = y^1 \begin{bmatrix} x_{11} \\ x_{12} \\ 1 \end{bmatrix} + y^2 \begin{bmatrix} x_{21} \\ x_{22} \\ 1 \end{bmatrix} + y^3 \begin{bmatrix} x_{31} \\ x_{32} \\ 1 \end{bmatrix} + y^4 \begin{bmatrix} x_{41} \\ x_{42} \\ 1 \end{bmatrix} = b$$

$$\sum_i \alpha_i x_i^T = \begin{bmatrix} x_{11} \\ x_{12} \\ 1 \end{bmatrix} [x_{11} \ x_{12} \ 1] + \begin{bmatrix} x_{21} \\ x_{22} \\ 1 \end{bmatrix} [x_{21} \ x_{22} \ 1] + \begin{bmatrix} x_{31} \\ x_{32} \\ 1 \end{bmatrix} [x_{31} \ x_{32} \ 1] + \begin{bmatrix} x_{41} \\ x_{42} \\ 1 \end{bmatrix} [x_{41} \ x_{42} \ 1]$$

$$\therefore -\vec{b} + A\vec{w} = 0$$

$$\vec{w} = (A^T)^{-1} b$$

Disadvantage: ① A^T not always invertible
 ② For problems having more no. of features it's difficult to invert mat A

Obtain

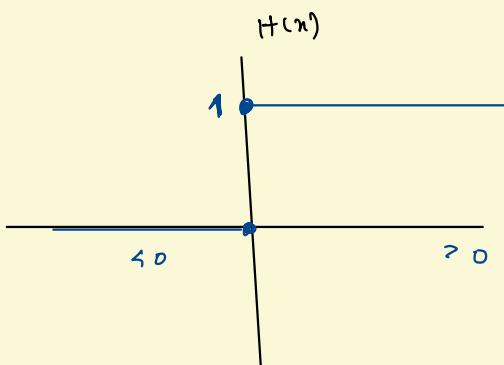
$$y = \text{sgn}(w^T x) = \begin{cases} +1 & \text{if } w^T x \geq 0 \\ -1 & \text{otherwise.} \end{cases}$$

Minimizing Classification Error:

(Aim: to find optimum \vec{w})

Idea: to minimize misclassification

$$\textcircled{1} \quad -y_i^* w^T x_i = \begin{cases} \geq 0, & \text{misclassification} \\ < 0, & \text{correct classification.} \end{cases} \quad (\text{Do we agree?})$$



$$E_0(w) = \sum_{i=1}^N H(-y_i^* w^T x_i)$$

"Hard errors"

Claim: $E_0(w)$ counts the misclassification.

if misclassification: $-y_i^* w^T x_i \geq 0 \quad \& \quad \therefore H(\geq 0) = 1$
 It assigns count 1 for each misclassification

if correct classification: $y_i^* w^T x_i < 0 \quad \& \quad H(< 0) = 0$
 It assigns count 0 for each correct classification.

So if we minimize $E_0(w)$ & find optimum w we're done.

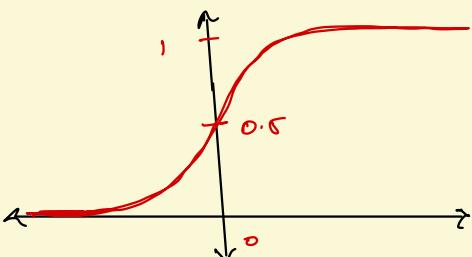
problem: The obj. func. E_0 is difficult to optimize.
 why? $\nabla H = 0$ a.e except @ origin

\therefore Minimizing $E_0(w)$ not a good idea.

2

Using the Sigmoid function.

$$l(x) = \frac{1}{1 + e^{-x}}$$



as $x \rightarrow \infty \quad l(x) \rightarrow 1 \quad \text{and} \quad x \rightarrow -\infty \quad l(x) \rightarrow 0$

loss func is probability of misclassification.

threshold = 0.5

$$-y_i \omega^T x_i = \begin{cases} \geq \text{threshold}, & \text{misclassification} \\ < \text{threshold}, & \text{correct classification.} \end{cases}$$

"Soft Error"

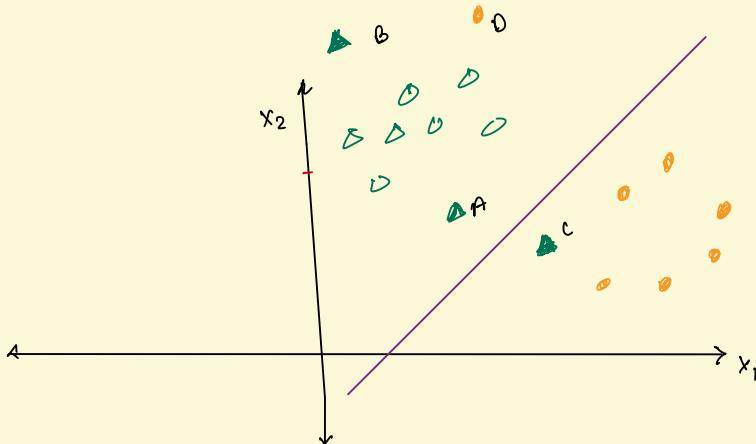
$$E_1(\omega) = \sum_{i=1}^N l(-y_i \omega^T x_i)$$

loss function = 1 when misclassified
loss func. = 0 when correctly classified

Observation: ① $l(-y_i \omega^T x_i)$ is a quantity b/w 0 and 1

$$-y_i \omega^T x_i > 0 \quad E_1(\omega) = \sum_{i=1}^N l(-y_i \omega^T x_i)$$

Just think about this

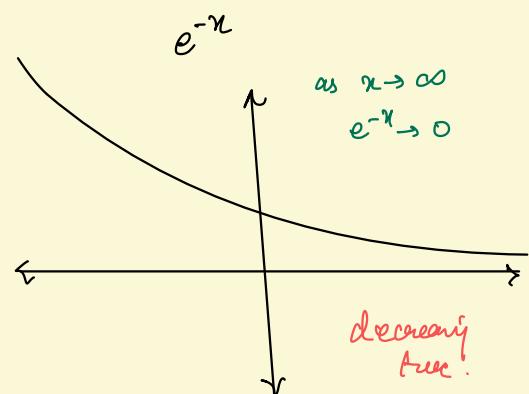


what can you say about data points value of l at points :

(a) A and B

(b) A : $l(-y_i \omega^T x_i)$ threshold

(c) B : $l(-y_i \omega^T x_i)$ threshold



some more : $x_1 < x_2$

$$e^{-x_1} \geq e^{-x_2} \Rightarrow 1 + e^{-x_1} \geq 1 + e^{-x_2}$$

$$\frac{1}{1 + e^{-x_1}} \leq \frac{1}{1 + e^{-x_2}} \Rightarrow l(x_1) \leq l(x_2)$$

Both $-y_i \omega^T x_i$ and $-y_i^* \omega^T x_i^*$ are -ive at A & B. L: they are correctly classified.

$$-y_i \omega^T x_i @ B < -y_i^* \omega^T x_i^* @ A$$

$$l(-y_i \omega^T x_i) @ B < l(-y_i^* \omega^T x_i^*) @ A$$

loss @ B \approx loss @ A \therefore B is more correctly classified than A

(B) at C & D

$$\ell(-y_i \omega^T x_i) < \text{threshold} @ C$$

$$\ell(-y_i \omega^T x_i) < \text{threshold} @ D$$

both

$-y_i \omega^T x_i @ C$ &

$-y_i \omega^T x_i @ D$ are
true. (\because misclassified)

$$-y_i \omega^T x_i @ C < -y_i \omega^T x_i @ D$$

$$\therefore \ell(-y_i \omega^T x_i) @ C < \ell(-y_i \omega^T x_i) @ D$$

(It's a misclass.
(It is a negative
value
remember))

$\hookrightarrow \text{loss } @ C < \text{loss } @ D$

i.e. C is comparatively less misclassified than D.

Minimizing E_1 :

$$E_1(\omega) = \sum_{i=1}^N \ell(-y_i \omega^T x_i) = \text{Probability of misclassification}$$

Minimizing $E_1(\omega)$

$$E_1(\omega) = \sum_{i=1}^N \frac{1}{1 + e^{y_i \omega^T x_i}}$$

Minimize ω :

$$\frac{\partial E_1}{\partial \omega} = 0$$

$$\frac{\partial E_1}{\partial \omega} = \sum_{i=1}^N \frac{\partial \ell(-y_i \omega^T x_i)}{\partial \omega}$$

$$\frac{\partial \ell(-y_i \omega^T x_i)}{\partial \omega} = \frac{\partial}{\partial \omega} \left(\frac{1}{1 + e^{y_i \omega^T x_i}} \right) = -\frac{e^{y_i \omega^T x_i} (y_i x_i)}{(1 + e^{y_i \omega^T x_i})^2}$$

$$\ell(-y_i \omega^T x_i) = \frac{1}{1 + e^{y_i \omega^T x_i}}$$

$$1 - \ell(-y_i \omega^T x_i) = -\frac{e^{y_i \omega^T x_i}}{1 + e^{y_i \omega^T x_i}}$$

$$1 - \ell(-x) = 1 - \frac{1}{1 + e^x} = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}} = \ell(x)$$

$$1 - \ell(-x) = \ell(x)$$

$$\frac{\partial E_1}{\partial w} = \sum_{i=1}^N y_i \ln(y_i w^T x_i) + (1 - \ln(y_i w^T x_i)) x_i$$

Thus by minimizing $E_1(w)$ we can get the optimal w .

We can view $\ln(1 - y_i w^T x_i)$ as prob. of making an error on training sample (x_i, y_i) using model w .

* Logistic Regression :-

Joint probability distribution:

correct classification:

prob. of misclassification: $\ln(1 - y_i w^T x_i)$

prob. of correct classification: $1 - \ln(1 - y_i w^T x_i)$ (But $1 - \ln(-x) = \ln(x)$)

Prob. of a correct classification @ (x_i, y_i) is: $\ln(y_i w^T x_i)$

\therefore Prob. of correct classification for all samples in training set assuming samples are independent & iid.

$$E_{LR}(w) = \arg \max_w L(w) = \arg \max_w \prod_{i=1}^N \ln(y_i w^T x_i) \rightarrow \text{Maximize the prob. of correct classification.}$$

$$= \arg \max_w \ln L(w) = \arg \max_w \sum_{i=1}^N \ln(\ln(y_i w^T x_i))$$

$$(1 - \ln(-x) = \ln(x))$$

Equivalent to

$$\text{minimize} \quad \sum_{i=1}^N -\ln(1 - \ln(y_i w^T x_i))$$

$$\sum_{i=1}^N -\ln(1 - \ln(y_i (\omega^T x_i + b))) =$$

$$\frac{1}{1 - \ln(-x)} = \frac{1}{1 - \ln(x)}$$

$$F(\omega) = \sum_{i=1}^N -\ln(\ell(y_i^\circ \omega^T x_i^\circ))$$

$$\ell(-x) = 1 - \ell(x)$$

No constraints.

$$\log n = \frac{1}{2}$$

$$F(\omega) = \sum_{i=1}^N -\ln\left(\frac{1}{1 + e^{-y_i \omega^T x_i^\circ}}\right)$$

$$\frac{\partial F}{\partial \omega} = \sum_{i=1}^N -\left(1 + e^{-y_i \omega^T x_i^\circ}\right)^{-1} \frac{e^{-y_i \omega^T x_i^\circ} (y_i^\circ x_i^\circ)}{(1 + e^{-y_i \omega^T x_i^\circ})^2}$$

$$= -\sum_{i=1}^N \frac{y_i e^{-y_i \omega^T x_i^\circ}}{(1 + e^{-y_i \omega^T x_i^\circ})^2}$$

$$1 - \ell(y_i^\circ \omega^T x_i^\circ) = 1 - \frac{1}{1 + e^{-y_i \omega^T x_i^\circ}} = \frac{e^{-y_i \omega^T x_i^\circ}}{1 + e^{-y_i \omega^T x_i^\circ}}$$

* Multivariable linear regression

* The model is given as:

$$f(\omega, b) = \omega_1 x_1 + \omega_2 x_2 + b$$

$$x_{\text{train}} = \begin{bmatrix} 1 & x_1 & \dots & x_n \\ \vdots & \vdots & & \vdots \\ m & x_m & & \end{bmatrix}$$

* Y-pred.

"m samples"

$$y_{\text{pred.}} = \begin{bmatrix} 1 \\ 2 \\ \vdots \\ m \end{bmatrix} \begin{bmatrix} \omega_1 x_{11} + \omega_2 x_{12} + b \\ \omega_1 x_{21} x_{22} + b \\ \vdots \\ \omega_1 x_{m1} + \omega_2 x_{m2} + b \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ \vdots \\ m \end{bmatrix} \begin{bmatrix} [x_{11}, x_{12}] \cdot [\omega_1, \omega_2] + b \\ [x_{21}, x_{22}] \cdot [\omega_1, \omega_2] + b \\ \vdots \\ [x_{m1}, x_{m2}] \cdot [\omega_1, \omega_2] + b \end{bmatrix}$$

$$\begin{bmatrix} [x_{11}, x_{12}] \\ \vdots \\ [x_{m1}, x_{m2}] \end{bmatrix} \begin{bmatrix} [\omega_1] \\ \omega_2 \end{bmatrix} + b = \text{np.dot}(x, \text{weights}) + b$$

matrix multiplication

$$\text{np.dot}(A, B) = AB - \text{mat. multp}$$

$$\text{Model: } f_{\omega, b}(x^i) = \omega_1 x_{i1} + \omega_2 x_{i2} + b$$

* The cost function is given as:

$$J(\omega, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\omega, b}(x^i) - y^i)^2$$

$$x_{\text{train}} = \begin{bmatrix} 1 & x_1 & \dots & x_n \\ \vdots & \vdots & & \vdots \\ m & x_m & & \end{bmatrix}$$

Gradient descent algorithm:

$$\omega = \omega - \alpha \frac{1}{m} \sum_{i=1}^m (f(x^i) - y^i) x^i$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f(x^i) - y^i)$$

$$y_{\text{train}} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$y_{\text{pred}} = \begin{bmatrix} \omega^T x_1 + b \\ \omega^T x_2 + b \\ \vdots \\ \omega^T x_m + b \end{bmatrix}$$

$$\text{def } d\omega = \frac{1}{m} \sum_{i=1}^m (f(x^i) - y^i) x^i$$

dot product

$$d\omega = \begin{bmatrix} d\omega_1 \\ d\omega_2 \end{bmatrix}$$

repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

→ sum over all first components

2 features →

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}$$

...

}

$$\theta_0 = \omega_0$$

$$\theta_1 = \omega_1$$

$$\theta_2 = \omega_2$$

$y - y_{\text{pred}}$ =

$$\begin{bmatrix} \omega^T x_1 + b - y_1 \\ \omega^T x_2 + b - y_2 \\ \vdots \\ \omega^T x_m + b - y_m \end{bmatrix}$$

$$(x_{\text{train}})^T = \begin{bmatrix} x_1 \\ \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix} & \begin{bmatrix} x_2 \\ \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} \end{bmatrix} & \dots & \begin{bmatrix} x_m \\ \begin{bmatrix} x_{m1} \\ x_{m2} \end{bmatrix} \end{bmatrix} \end{bmatrix}$$

$$[\omega^T x_i + b - y_i] \quad \text{with} \quad [\omega^T x_1 \quad x_2]$$

$$(x_{\text{train}})^T \cdot (y_{\text{pred}}) = \begin{bmatrix} x_1 \\ \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix} & \begin{bmatrix} x_2 \\ \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} \end{bmatrix} & \dots & \begin{bmatrix} x_m \\ \begin{bmatrix} x_{m1} \\ x_{m2} \end{bmatrix} \end{bmatrix} \end{bmatrix} \cdot \begin{bmatrix} f(x_1) - y_1 \\ f(x_2) - y_2 \\ \vdots \\ f(x_m) - y_m \end{bmatrix} \quad \text{matrix multp}$$

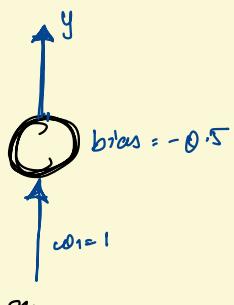
$$= \begin{bmatrix} x_{11} (f(x_1) - y_1) + \dots + x_{m1} (f(x_m) - y_m) \\ x_{12} (f(x_1) - y_1) + \dots + x_{m2} (f(x_m) - y_m) \end{bmatrix} = dW$$

$$\therefore d\omega = \left(\frac{1}{m} \right) n p \cdot \text{dot} ((x_{\text{train}})^T, (y_{\text{pred}}))$$

$$db = \left(\frac{1}{m} \right) n p \cdot \text{sum} (y_{\text{pred}})$$

Why do we need sigmoids?

↳ Thresholding logic used by perceptron is very harsh



like (Dislike a movie, bias = -0.5 \Rightarrow thresh: 0.5)

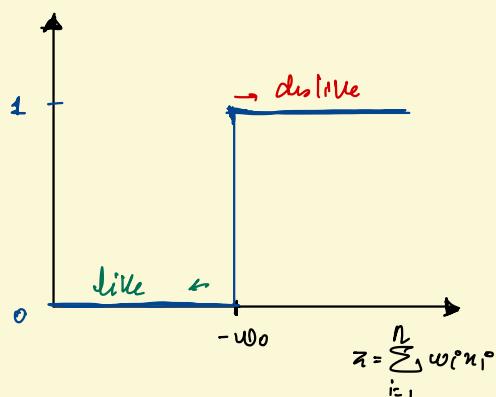
if critics = 0.51 \geq 0.5

if o/p = 0.49

like ✓
dislike ✗

decision is harsh.

Critics really



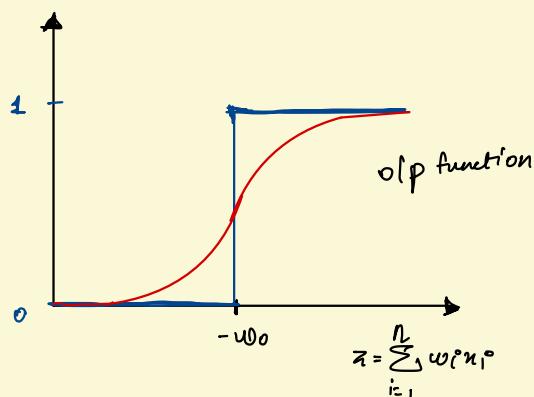
if $z \geq -w_0$ fire,
else don't fire.

This is harsh
Not continuous, not diff.

If else present

Sigmoid: is a family of functions

One such func is logistic function.



$$y = \frac{1}{1 + e^{-w^T x}}$$

No longer a sharp transition around the threshold
 $-w_0$
& value is b/w 0 & 1

Interpretation of sig. func. as a prob:

Smooth, cont & diff.

With 51% prob will like the movie

With 49% prob will like the movie