

Supervised learning

$x \rightarrow y$
 Input output

* Type

1) Regression — predict number

2) classification — predict categories

* Structure

training set



learning algorithm

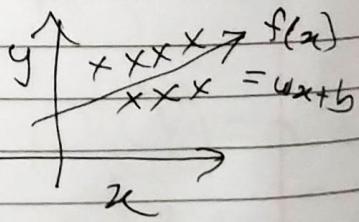
$x \rightarrow [f] \rightarrow \hat{y}$
 feature model prediction
 (estimated y)

for regression with straight line and one variable
 we can define with

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

or $f(x)$

univariate linear regression

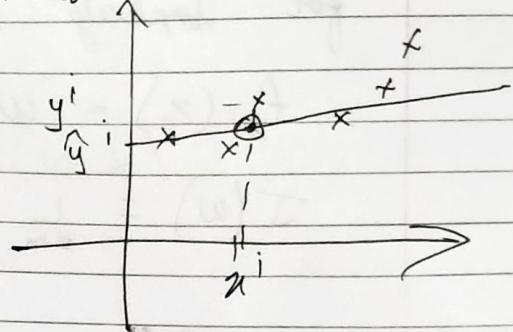


$y^{(i)}$ \rightarrow actual value

$\hat{y}^{(i)}$ \rightarrow Prediction our model made

$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

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



* cost function

$$\left(\frac{\hat{y}^{(i)} - y^{(i)}}{\text{error}} \right) \left[\frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 \right] = J(w, b)$$

Mean Squared Error cost function

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

cost function is To check how well
predictor line
the cost function fit the dat

so our goal is to minimize cost func.

i.e minimize $J(w, b)$

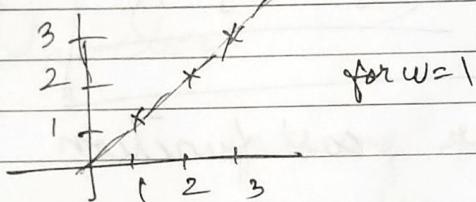
for simplicity lets take $b = 0$

$$\therefore f_w(x) = w x \quad \cancel{f_w(x)}$$

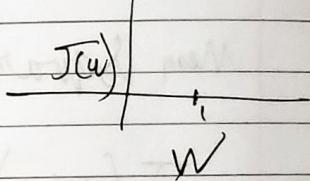
$$\therefore J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

$f_w(x)$ if w is fixed, function of x
is depended of w

$J(w)$ depends on w ie parallel

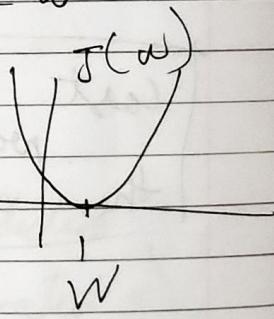


$$J(w=1) = 0$$



for all w

so as we can see
our goal is to reduce
 w as much as possible to
get the best fitted line



for both w and b you will get a
3d plot like a bowl those are very
wavy

Gradient descent

It helps us find the local minima
easily instead of doing all things
manually

GD algorithm -

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

α = learning rate

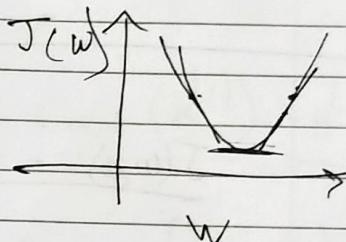
$$b_f = b_i - \alpha \frac{\partial}{\partial b} J(w, b)$$

again consider $b=0$

$$J(w) = w^2$$

Then for G.D

$$\begin{aligned} w_f &= w_i - \alpha \frac{\partial}{\partial w} J(w) \\ &= w_i - \alpha w \end{aligned}$$



$$f_{w,b} = w_x + b$$

$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$ — prediction of our model
 $y^{(i)}$ — actual value

$$\hat{y}^{(i)} - y^{(i)} \rightarrow \text{error}$$

$$\text{cost} = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

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

Our job is to find (w, b) so cost = 0

so to minimize cost
we use gradient descent in which we
update w, b simultaneously

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

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

$\alpha \Rightarrow$ learning rate

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

def compute_cost(x, y, w, b):

~~m = n~~ shape [0]

cost = 0

for i in range(m):

$$f_{wb} = w \times x[i] + b$$

$$y_i = y[i]$$

$$\text{cost} += (f_{wb} - y_i)^2$$

return cost / 2m

compute_cost gives us $J(w, b)$

In gradient descent we do

repeat until convergence : {

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

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

}

$$\text{So } \frac{\partial J(w, b)}{\partial w} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$\frac{\partial J(w, b)}{\partial b} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{w,b}(x^{(i)}) - y^{(i)})$$

Move onto ~~multivariable~~, multiple linear regress when we have :

$$f_{w,b}(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$$

$$f_{w,b}(x) = w_1 x_1 + w_2 x_2 \dots w_n x_n + b$$

take \vec{w} as $[w_1, w_2, w_3, \dots, w_n]$
and \vec{x} as $[x_1, x_2, x_3, \dots, x_n]$

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

functions in vector

* Vectorization

$$\text{Let } \vec{w} = [w_1 \ w_2 \ w_3]$$

$$\vec{x} = [x_1 \ x_2 \ x_3]$$

$$f = \text{np. dot}(w, x) + b$$

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

$$J_{\vec{w}, b} = \frac{1}{2m} \sum_{j=0}^{m-1}$$

way faster than writing manually or
using a for loop

* Feature scaling

* Mean Normalization (for a dummy dataset)

if average of $x_1 \rightarrow 600 = \bar{x}_1$

$$\text{we do } x_1' = \frac{x_1 - \bar{x}_1}{2000 - 300}, \quad x_2' = \frac{x_2 - \bar{x}_2}{5 - 0}$$

$$-0.18 \leq x_1' \leq 0.82$$

$$-0.46 \leq x_2' \leq 0.52$$

* Z-score normalization
find standard deviation σ

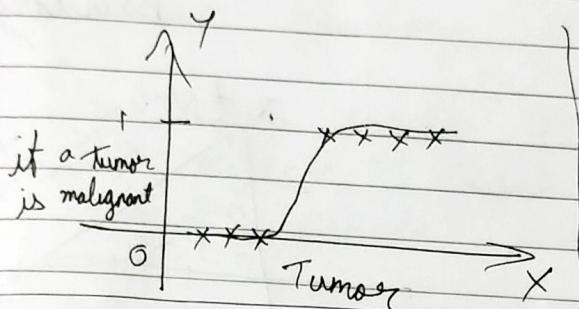
$$\text{so } x_1' = \frac{x_1 - \bar{x}_1}{\sigma_1} \quad \text{and } x_2' = \frac{x_2 - \bar{x}_2}{\sigma_2}$$

do feature scaling if the any
feature is too large / small compared
to other feature

* Feature engineering

Combining features to make new feature

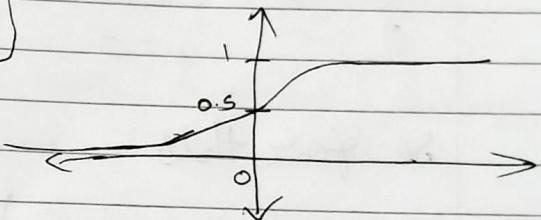
Logistic regression → classifying algorithm



to map this lines we use the sigmoid function

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

$$0 < g(z) < 1$$



we take
we know $w \cdot x + b = z$

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

we take a threshold

$$\text{if } f_{\vec{w}, b}(\vec{x}) \geq 0.5 \text{ then } \hat{y} = 1$$

yes $\rightarrow \hat{y} = 1$

else

no $\rightarrow \hat{y} = 0$

when is $f_{\vec{w}, b}(\vec{x}) \geq 0.5$

$$g(z) \geq 0.5$$

$$z \geq 0 \Rightarrow \vec{w} \cdot \vec{x} + b \geq 0$$

so when
 $\vec{w} \cdot \vec{x} + b \geq 0$
 $y = 1$

else
 $\vec{w} \cdot \vec{x} + b < 0$
 $y = 0$

When $z = \vec{w} \cdot \vec{x} + b = 0 \Rightarrow$ decision boundary line
 cl

example \rightarrow

$$\begin{aligned} x &\rightarrow y = 1 \\ 0 &\rightarrow y = 0 \end{aligned}$$

so for this

$$f_{w,b}(\vec{x}) = g(z) = g(w_1x_1 + w_2x_2 + b)$$

as we have two features

$$\text{lets take } w_1 = 1, w_2 = 1, b = -3$$

$$\text{so we get } x_1 + x_2 - 3 = z$$

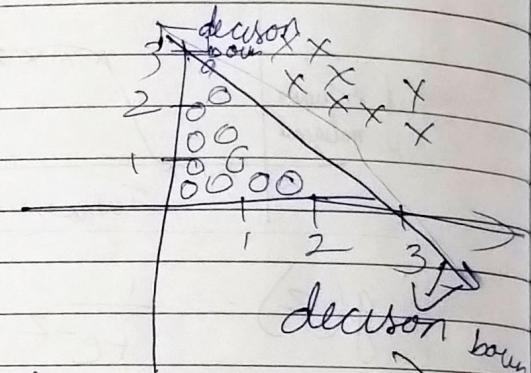
$$\text{decision boundary } \rightarrow x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3 \quad \text{positive}$$

anything above the line will be $1 \rightarrow \text{true}$
 and below will be $0 \rightarrow \text{negative}$

As with linear regression, we can do logistic regression with non linear functions to get non linear decision boundaries

Example

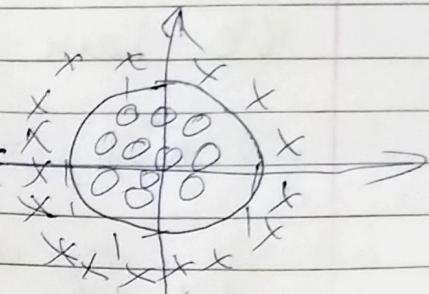


here

$$f_{\vec{w}, b}(\vec{x}) = g(z)$$

$$= g(w_1 x_1^2 + w_2 x_2^2 + b)$$

$$\text{Let } w_1 = 1, w_2 = 1, b = -1$$



: decision boundary $\Rightarrow x_1^2 + x_2^2 - 1 = 0$
 $x_1^2 + x_2^2 = 1$

anything inside the close boundary will be ~~false~~ 0 \rightarrow negative and any point outside circle will be 1 \rightarrow positive

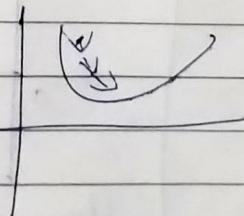
We can have even more complex decision boundaries by having more complex Z

Mean Squared cost function is not a good choice for regression

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

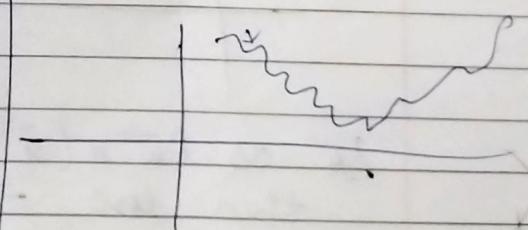
for linear regression

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



for logistic reg.

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



too many local minima

We have

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

Let's take $\frac{1}{2} f(\vec{w}, b)(\vec{x}^{(i)}) - y^{(i)} \Rightarrow$

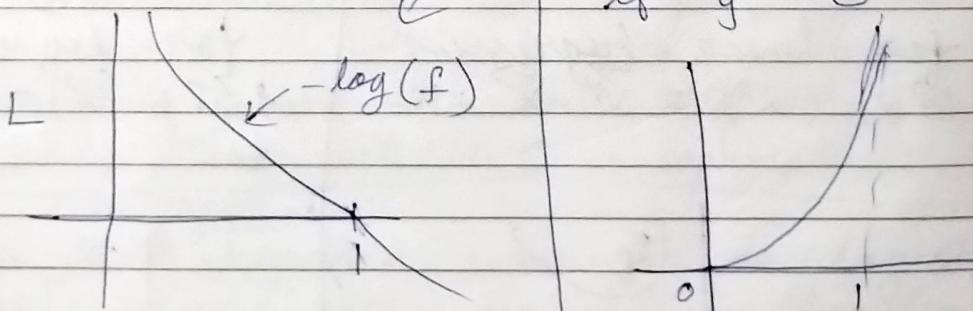
$$L(f(\vec{w}, b)(\vec{x}^{(i)}), y^{(i)})$$

So for logistic regression the cost function we will use is

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

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

If $y^{(i)} = 1$, then graph looks like



So as $f(\vec{w}, b)(\vec{x}^{(i)}) \rightarrow 1$

then loss $\rightarrow 0$

as it's the correct answer

else loss will be very
more

Here as $f(\vec{w}, b)(\vec{x}^{(i)}) \rightarrow 0$

loss approaches ∞

as it's the correct
label

Or we can simply write loss as

$$L(\vec{f}_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(\vec{f}_{\vec{w}, b}(\vec{x}^{(i)})) - (1-y^{(i)}) \log(1-\vec{f}_{\vec{w}, b}(\vec{x}^{(i)}))$$

if when $y^{(i)} = 0$ this becomes zero, else this

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

gradient descent

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

We do the same for as linear regression
for gradient descent

repeat {

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

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

} simultaneously updates

different from linear regression that

In L.I.R. $\vec{f}_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

In L.O.R. $\vec{f}_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$