

Supervised Machine Learning → ARJUN

Memo

Date:

ADHIKA RI

Machine learning.

How machine recognize?

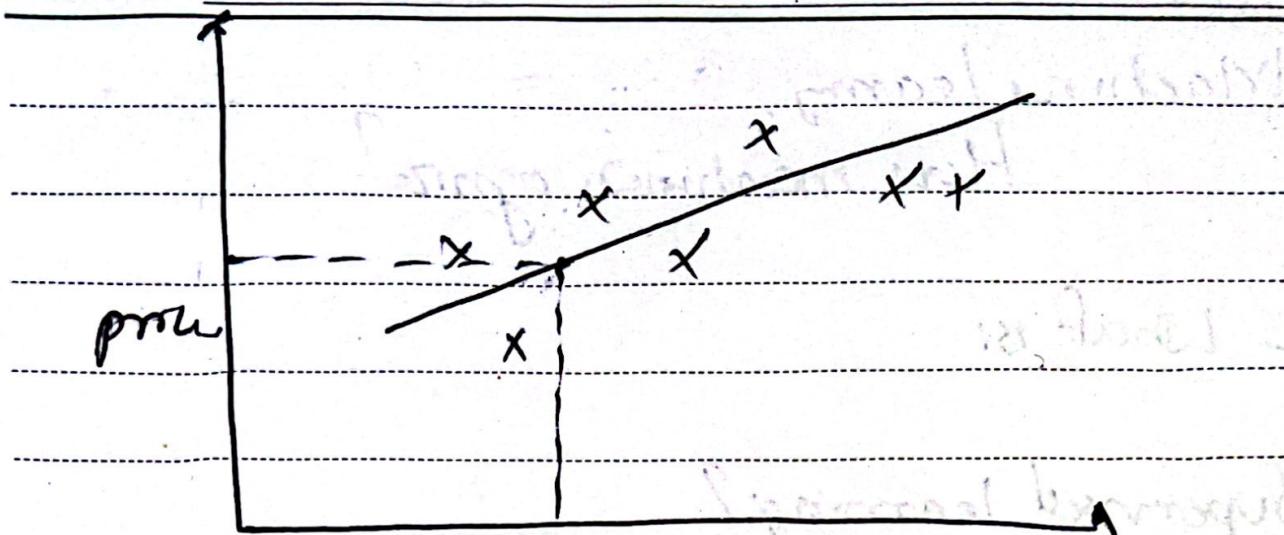
→ function.

What is:

Supervised learning?

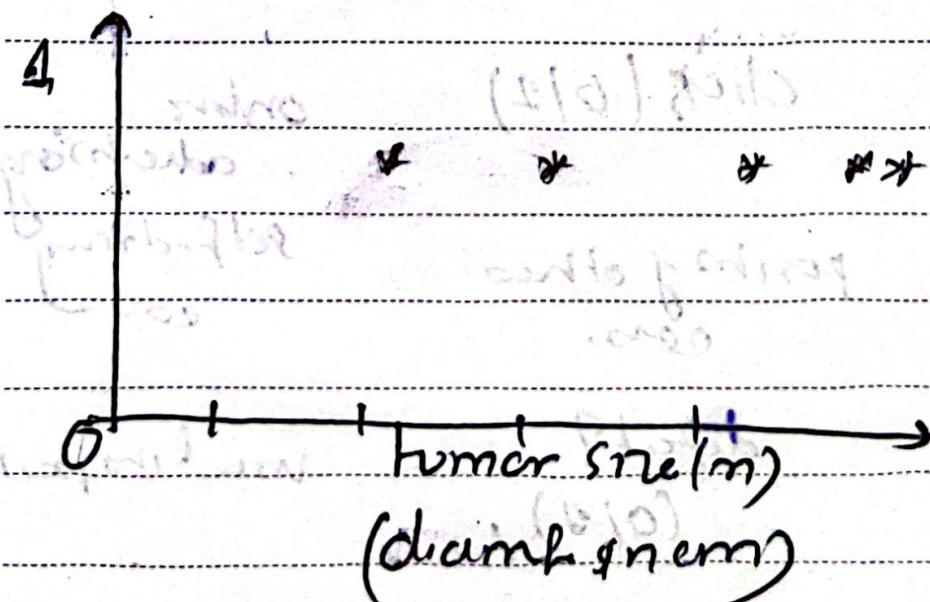
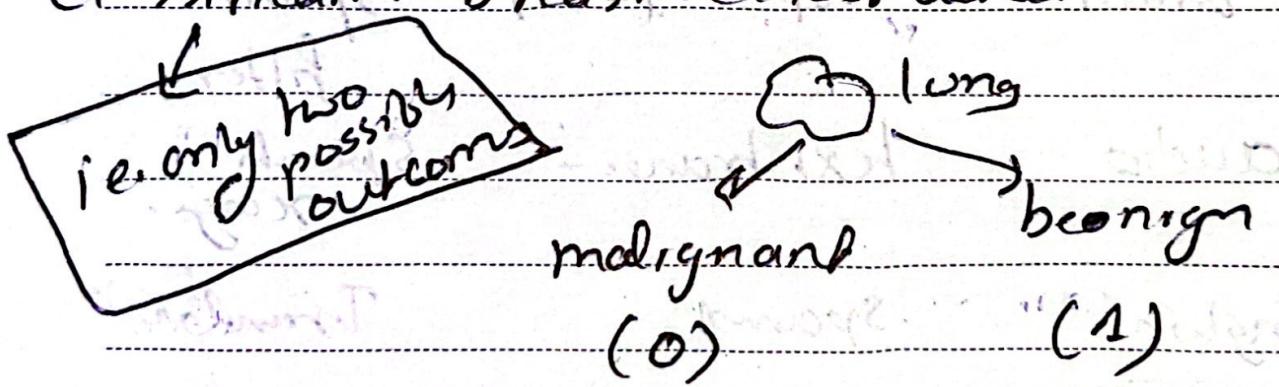
→ Algorithm to learn n to y .

Input (n)	Output (Y)	Applican
Email	Spam? (0/1)	Spam filters
audio	text transcr	Speech recog.
English	Spam	Translator
ad, user info	click (0/1)	online advertising
image, radar info	posting others cars.	self-driving car
image phon.	defect? (0/1)	visual inspection



Regression: House price prediction.

Classification: Breast Cancer detect?



Memo

Date:

1 0 0 0 0 x 0 x 0 x x x →

0 - benign

X Malign

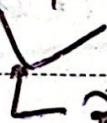
classifier:

predict word.

→ small number of possible outcomes

(0, 1, 2 - - only)

Predict categories / class



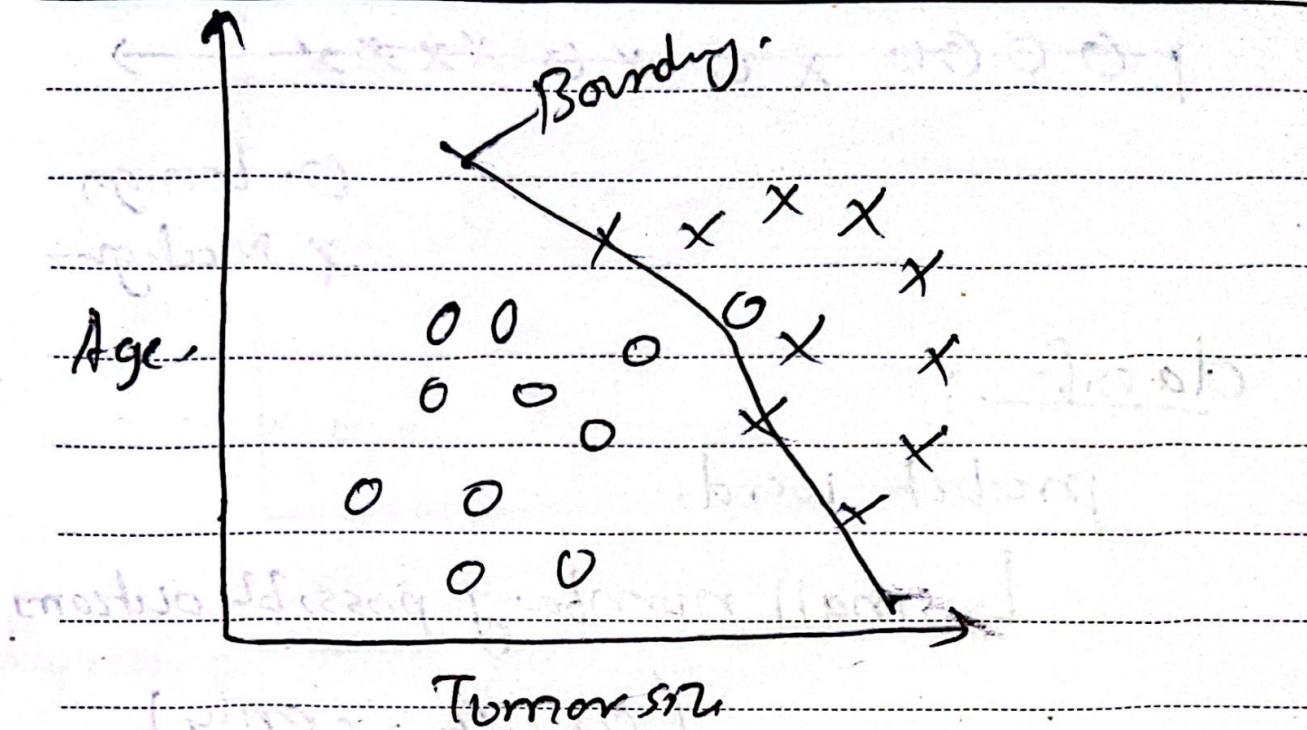
Def can be word than
number

Regression Cost, etc.

Regression

whether Cat or dog?
pictures

Malignant or benign?



Tu

Supervised learning

learn from being given right answer

Regression

predict a number

Inherently many

possible outcomes.

classification

predict categories

small number

Memo

Date:

Unsupervised learning

→ clustering

→ Dimensionality reduction
→ Anomaly detection.

→ Data only comes with inputs x_i , but not output labels y_i .

* Algorithm has to find structure in data.

clustering

Groups similar data

points together

Anomaly detection

find unusual data points.

Dimensionality reduction

Compresses data using fewer numbers.

Linear regression Model

predict numbers.

Training data
Train and Model.

Terminology:

① → Training set

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

② learning algorithm,

↓ produce

f

→ take new input n

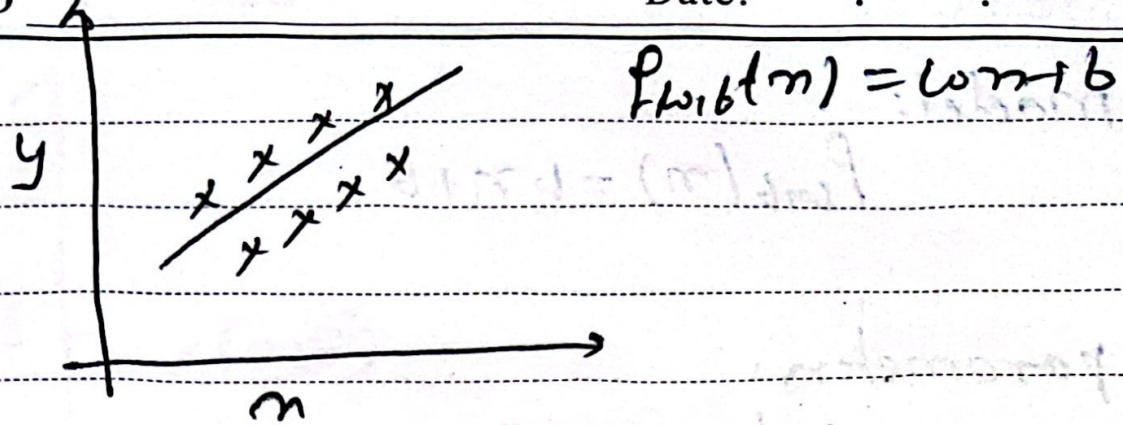
and output (prediction)

\hat{y}

Estimated prediction for model

~~f_{linear}~~

$$f_{w,b}(n) = w_n + b$$



linear regress with one variable.

Univariate Linear regression

Cost function:

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

m : number of training examples

Squared error cost function

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

For multiple,

$$J(\vec{w}, b)$$

Memo

Date:

model:

$$f_{w,b}(n) = w_n + b$$

parameters:

$$w, b.$$

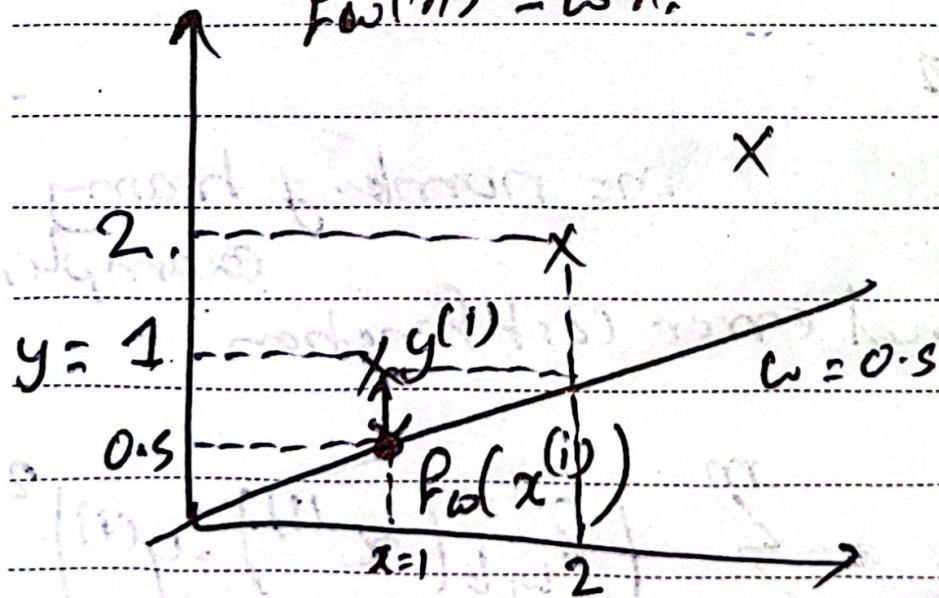
cost function:

Model pred - Actual value

goal.

(Minimize $J(\theta)$)

$$f_w(n) = w_n.$$

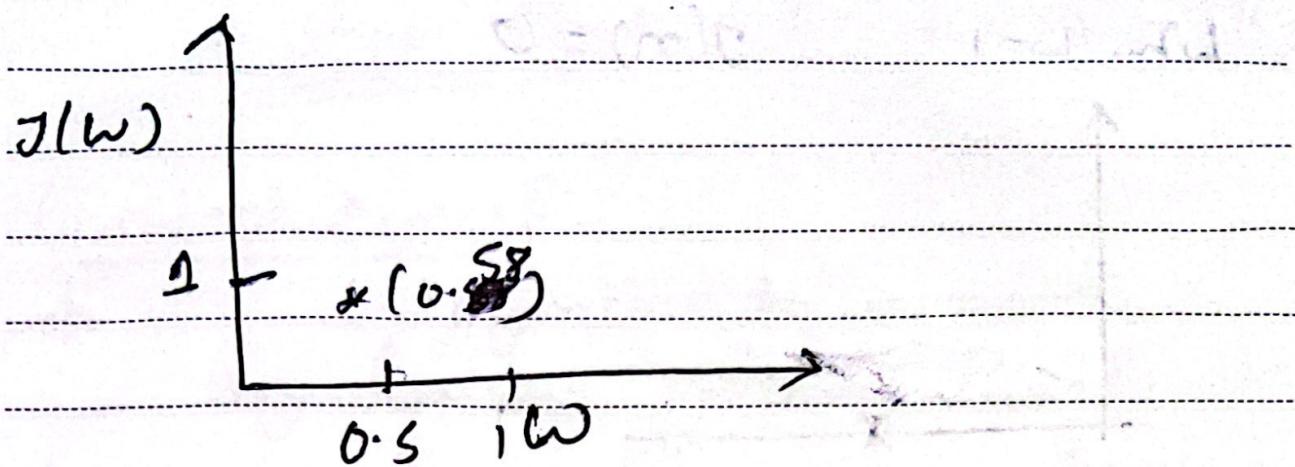


$$J(0.5) = \frac{1}{m} [(0.5 - 1)^2 + (1 - 2)^2]$$

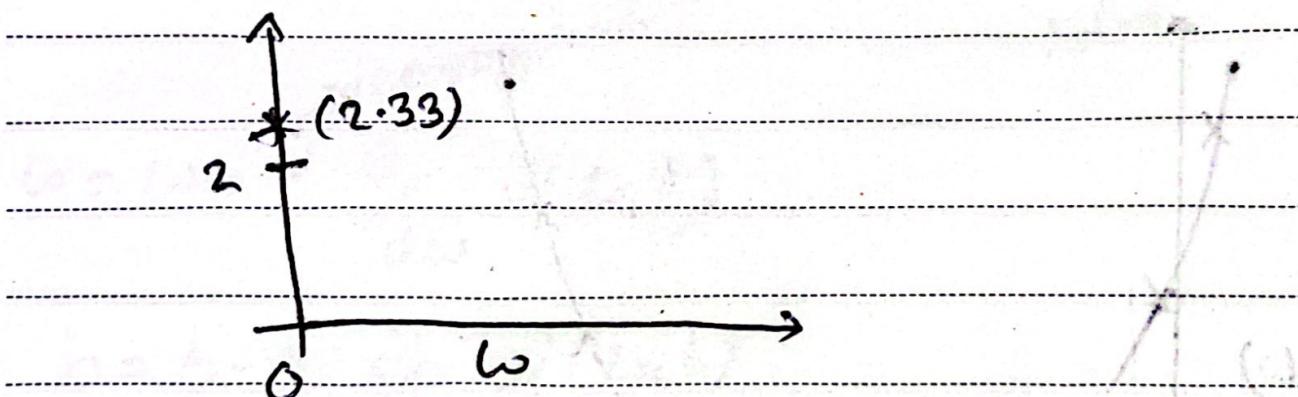
$$\approx 0.58$$

Memo

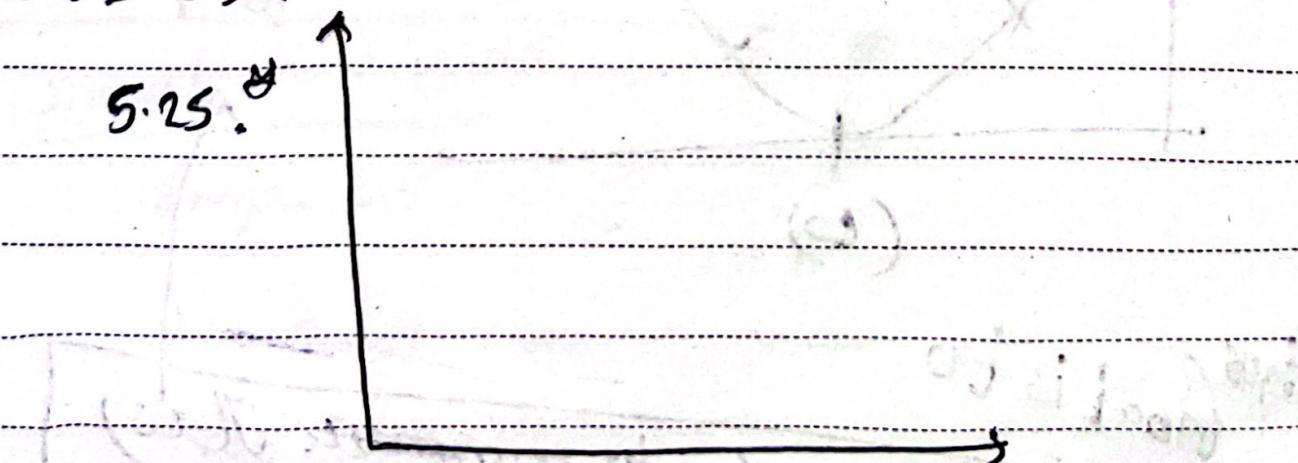
Date: 27.7.2022



When $\omega = 0$



When $\omega = -0.5\pi$

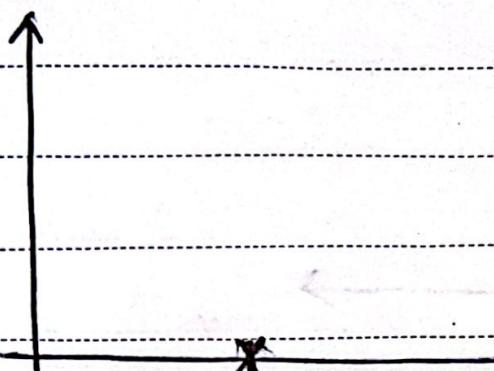


Memo

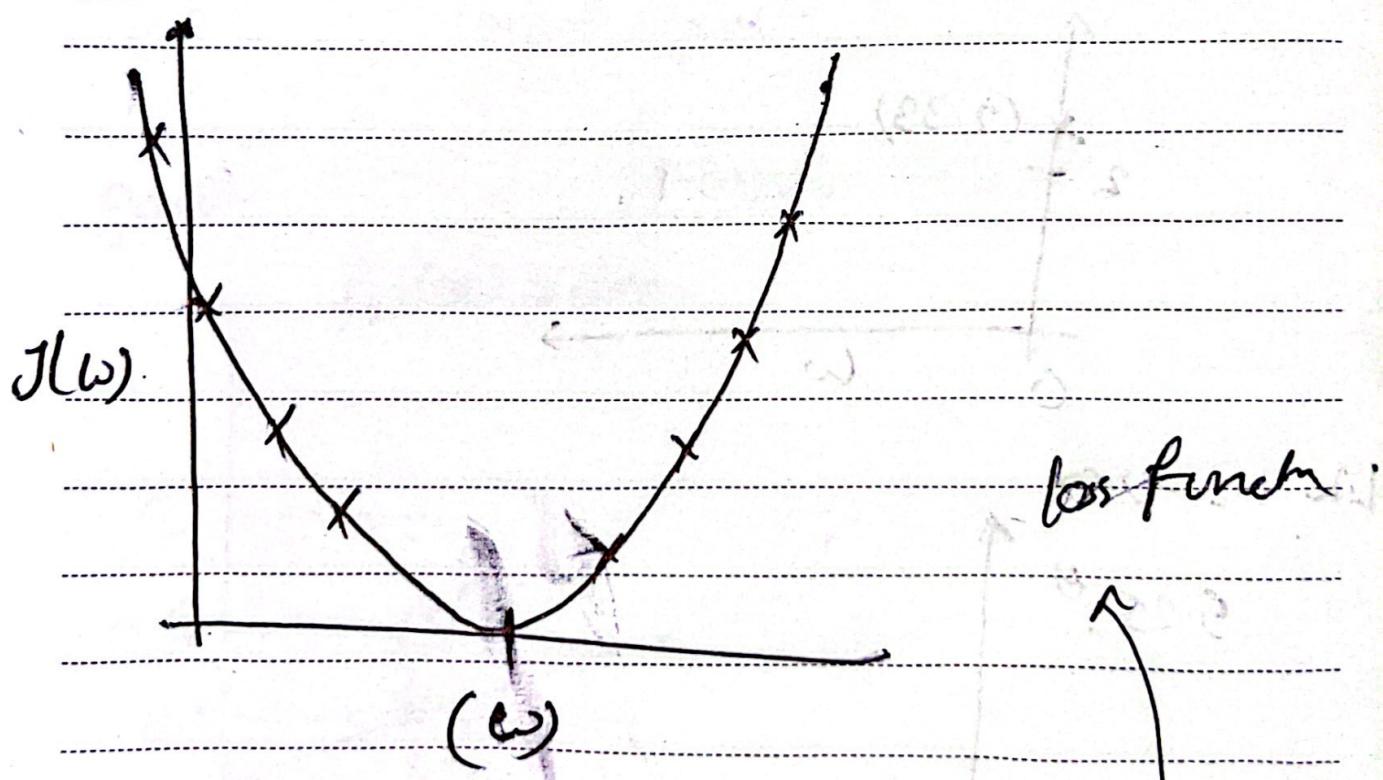
Date:

When $\omega=1$

$$J(\omega) = 0$$



Try



This goal is to

Choose ω to minimize $J(\omega)$

Memo

Date:

Shape of bowl (For two parameters)

Gradient Descent

* \rightarrow Algorithm to minimize
any function including cost functn.

~~w~~ \downarrow One dimension $\frac{\partial J}{\partial w}$

learning rate. Idea,

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

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

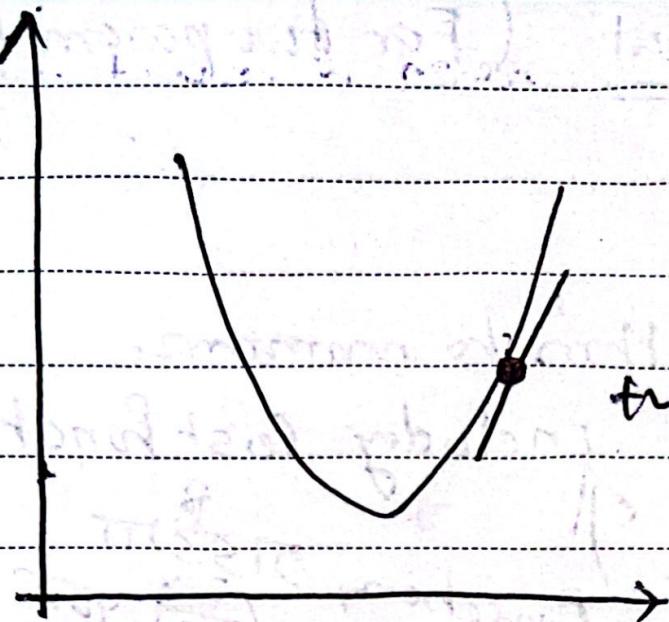
Simultaneous update

$$\text{tmp_} w =$$

$$\text{temp_} b$$

$$w = \text{temp_} w$$

$$b = \text{temp_} b$$



$$\Rightarrow \omega = \omega - d \cdot i(\theta)$$

sine function

= decrease



$$\omega = \omega + d$$

sine +ve slope

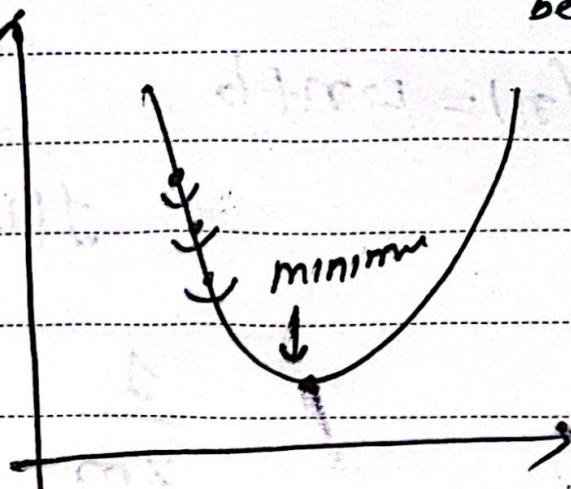
$$= \omega + d$$

= increase

Memo

Date:

If α is too small, Gradient descent may
be slow.



If α is large fail to converge.

At flat minimum,

$$\alpha = 0$$

$$w = w_0$$

Training Linear Regression.

Gradient Descent for Linear
Regression

Linear regression Model.

$$f_{\omega, b}(x) = \omega_0 + \omega_1 x$$

Cost function

$$J(\omega, b) =$$

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

Gradient descent algorithm

repeat until convergence {

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

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

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

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

"Batch" gradient descent.

↳ Each step of gradient descent
use all training examples.

Weeks - 2

Previously, ~~for~~ single variable

But now,

Multiple variables

Multiple variables

Price / m²
per m²

$$f_{w,b}(n) = w_1 n_1 + w_2 n_2 + \dots + w_n n_n + b$$

↑ ↑
size rooms

$$\vec{w} = [w_1, w_2, w_3, \dots, w_n] \Rightarrow$$

vector vector

vector parameter

$$\vec{x} = [x_1, x_2, x_3, \dots, x_n]^T \text{ model}$$

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

This is multiple linear regression

Vectorization Part-1.

$$\vec{w} = [w_1, w_2, w_3] \quad n=3$$

b is number

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

Numpy

$$w = np.array([1.0, 2.5, -3.3])$$

$$b = 4$$

$$x = np.array([10, 20, 30])$$

Vectorize,

$$\# f = np.dot(w, x) + b$$

Unvectorized

↳ For j in range(0, n)

$$f = f + w[j] * x[j]$$

Memo

Date: . . .

Vector

notations \times (bold)

`np.array([1, 2, 3])`

`a = np.arange(10)`

`c = a[2:7:2]`

`print("a[2:7:2] = ", c)`

start

stop

↓

2, 3, 4, 5, 6

end

`a[3:]` = 3, 4, 5, -

`a[:3]` = 0, 1, 2

`a[:]` = all.

`np.sum(a)` => sum all.

Memo

Date: . . .

Fn wr.

dot ():

return

a = np.array ([])

b = np.array ([])

c = np.dot (a, b).

$$a \cdot b = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = a_0 b_0 + a_1 b_1 + \dots$$

np.zeros(4) \Rightarrow 0, 0, 0, 0,

np.zeros(4,) \Rightarrow ''

a.shape \Rightarrow (4,)

np.zeros((1, 5)) \Rightarrow 0, 0, 0, 0, 0

a.shape \Rightarrow (1, 5)

But,

np.array ([[5], [4], [3]])

↑ a = [[5]

[4]

3 2 1

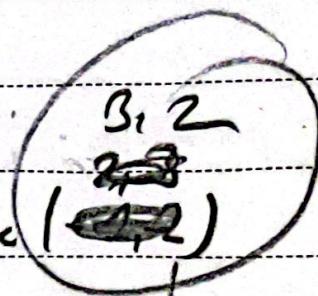
Memo

Date:

Creating matrix

reshape

$a = np.arange(6).reshape$



$a.reshape(3, 2)$

$a = [[0, 1]]$

$[2, 3]$

$[4, 5]]$

$np.reshape(a, [2, 3])$

as
 $np.arange(6).reshape(-2, 2)$

P(-1) says; given column and
no. of array,
calculate yourself.

Memo

Date:

Slicing

$a = np.arange(20).reshape(-1, 10)$

$a[0:2:2] \Rightarrow [2, 3, 4, 5, 6]$

Array are separated by commas

array = np.array([[[1, 2, 3], [4, 5, 6],
[7, 8, 9]]])

array.shape $\Rightarrow (3, 3)$

~~1D shape~~

a.shape \rightarrow no. of elements for 1D array

\rightarrow row / column for 2D array

Gradient descent for Multiple Regress.-

Memo

Date:

For n, features

Cost Function: $J(\vec{w}, b) \rightarrow [w_1, w_2, w_3]$

~~GD~~ Gradient descent

repeat {

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

$$w_n = b_n - \cdot \cdot \cdot$$

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

~~bb~~

Memo

Date:

$$X_Vec = X_train[0, :]$$

0 row at a time

X-vec value: [- - -],

~~Grade~~
 w, b

\downarrow
new value of f'

Feature and parameter value

#

Memo

Date:

Feature scaling is a method to scale numeric features in the same scale or range (like -1 to 2, 0 to 1)

Also called
Data normalization.

Apply on independent variable

" " train data and transform
train and test.

Why?

If height 140 cm and 8.2 feet

→ ML algorithm understand only
numbers but not feature unit

So according to ML, 140 > 8.2

#

Memo

Date:

$$\# \text{ price} = w_1 x_1 + w_2 x_2 + b$$

↓ ↓

x_1 bedrooms.

x_1 : size (feet²) x_2 : bedrooms.

range: 300 - 2,000 range 0.5.

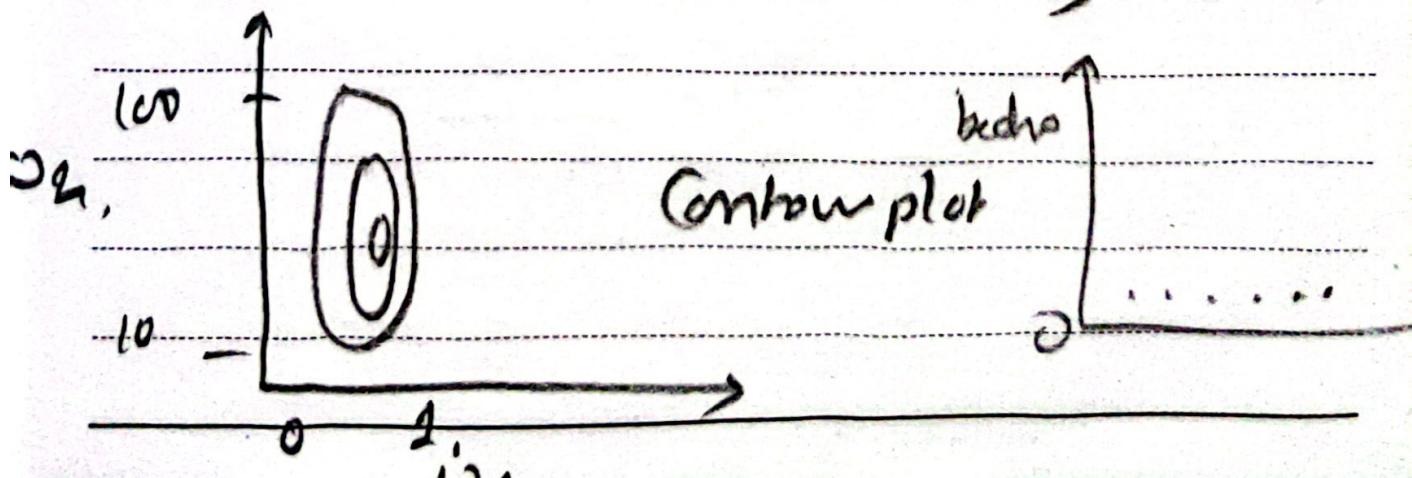
(large)

Actual: $x_1 = 2000$, $x_2 = 5$, price = \$500K.

Size of parameter w_1, w_2 ?

If $w_1 = 0.1$ (small), $w_2 = 50$ (large, $b = 50$ const.)

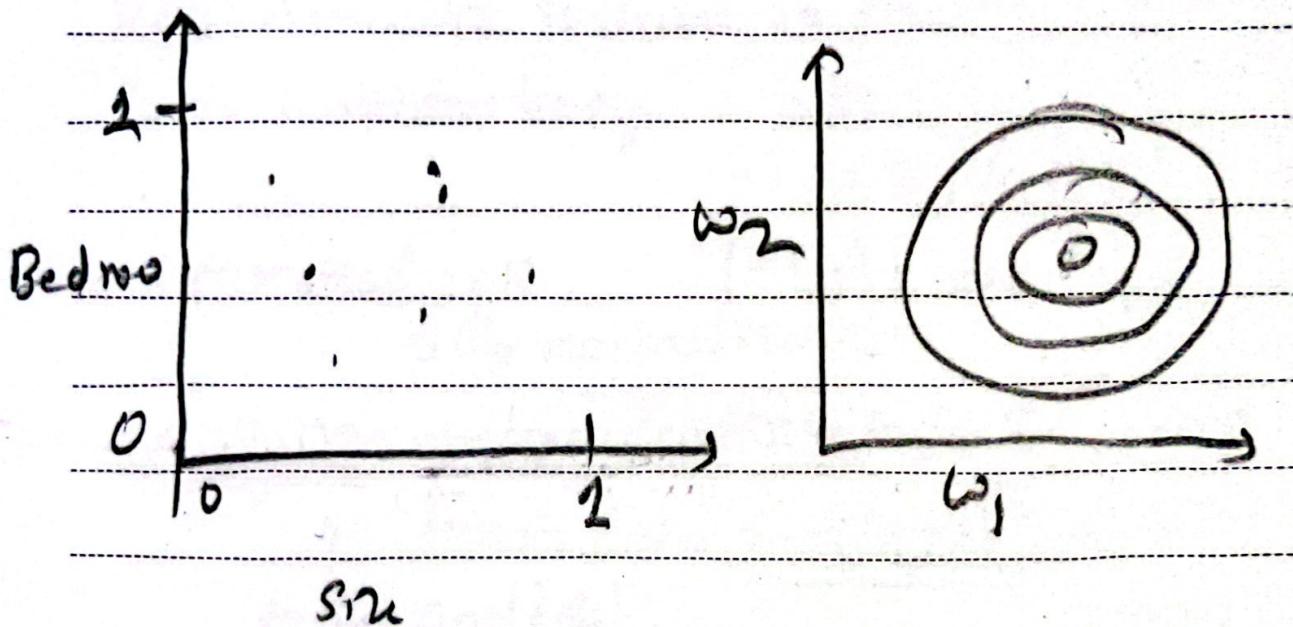
$$\begin{aligned} \text{price} &= 0.1 \times 2000 + 50 \times 5 + 50 \\ &= 500K \quad (\text{reasonable}) \end{aligned}$$



Memo

Date:

But rescaling.

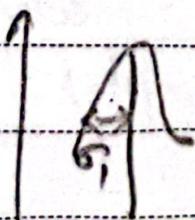


This rescaling helps in
decreasing gradient descent.
run fastly.

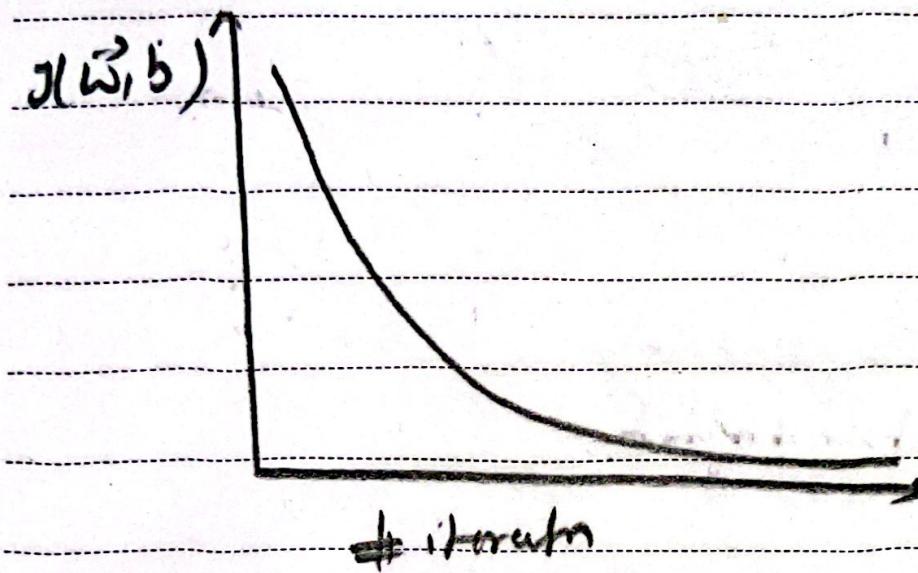
Rescaling.

- ① Dividing by max
- ② Mean normalization.
- ③ Z-score normalization.

$$x_i = \frac{x_i - \bar{x}_i}{\sigma_i}$$



Learning curve



Adamah converges if

$$\text{if } \varepsilon = 0.001$$

and

$J(\vec{w}, b)$ decrease by $\leq \varepsilon$

Then convergence or look into graph.

Memo

Value of α to try

0.001 0.003 0.01

x_3

x

Feature Engineering:
 $V = 1 \times b$

Polynomial regression.

Choice of features:

$f(\vec{x}, b | m) = \boxed{\quad}$

Linear Regression using

Scikit learn.

from sklearn.linear_model import

SGDRegressor

from sklearn.preprocessing import

StandardScaler

from joblib import

load_hour_data

Memo

Date:

Classification

→ logistic Regression.

↳ 0 and 1

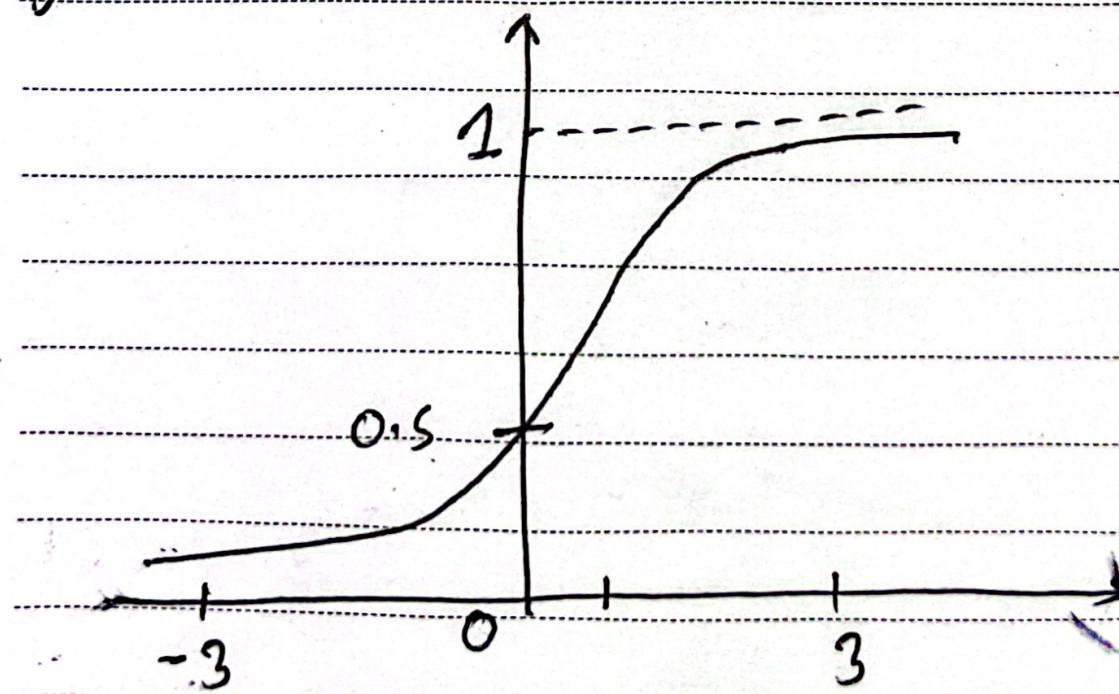
binary classification

Class = category.

& Decision boundary.

though
regression
is used
for classification

Sigmoid function.



Sigmoid function
outputs b/w 0 and 1

Memo

Date:

$$g(z) = \frac{1}{1+e^{-z}} \quad 0 < g(z) < 1$$

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

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



z

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

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

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

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

"

Logistic Regression."

Memo

Date:

$$P(\vec{w}_1, b)(\vec{x}) = 0.7$$

i.e. Model is 70% chancey one

Numpy,

`exp()`

`q = np.array([1, 2, 3])`

`np.exp(a)`

$g = \frac{1}{1 + np.exp(-z)}$



$\frac{1}{1 + e^{-z}}$

`np.set_printoptions(precision=3)`

`np.c_ = [np.array([1, 2, 3]), np.array([4, 5, 6])]`

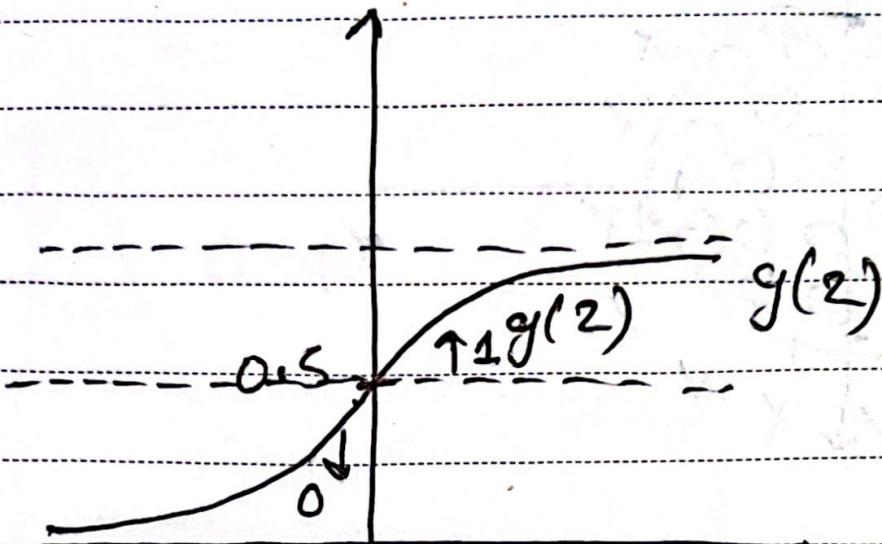
\Rightarrow `array([1, 2], [4, 5], [3, 6])`

Memo

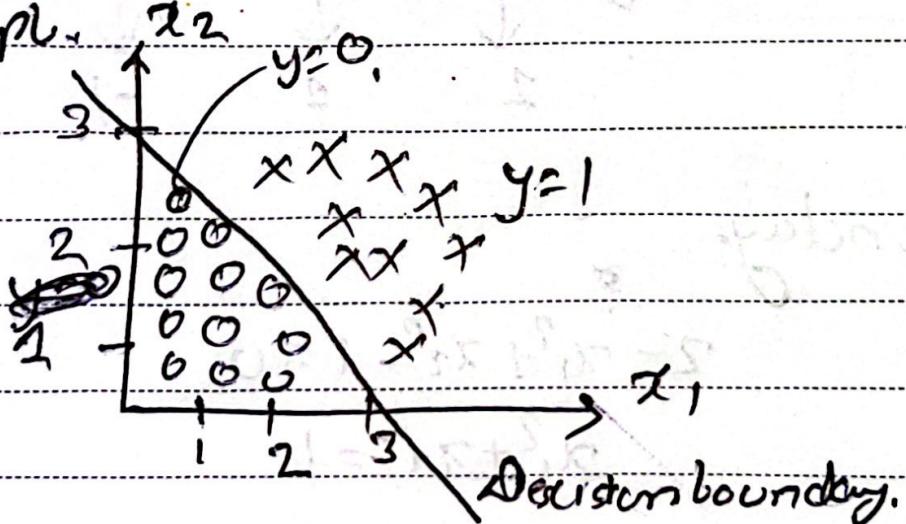
Date:

draw_vthresh(ax, o)

#



Exmpl.



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

and if $w_1 = 1, w_2 = 1, -3$

$$z = \vec{w} \cdot \vec{x} + b = 0 \quad \left. \right\} \text{decision boundary.}$$

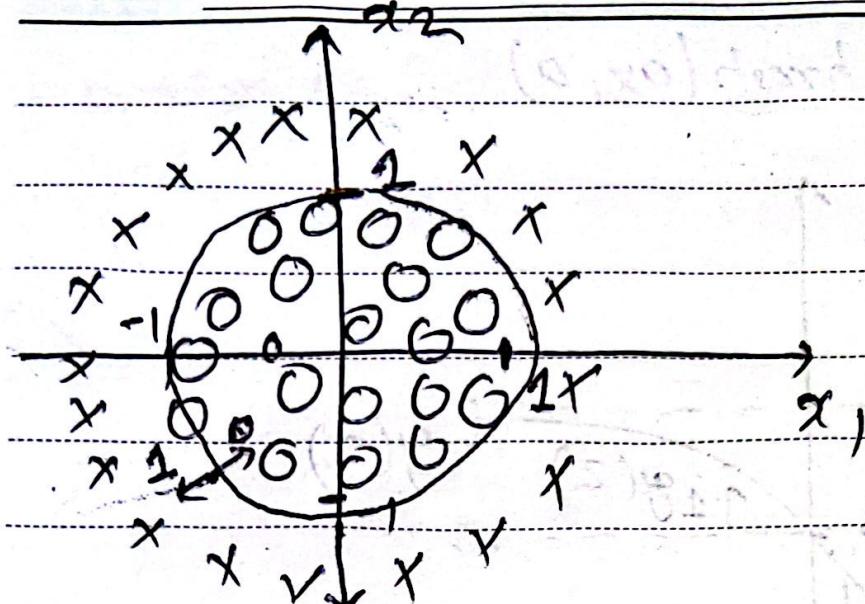
$$\bullet z = x_1 + x_2 - 3 = 0 \quad (\text{put } w_1 = 1, w_2 = 1, b = -3)$$

$$\bullet x_1 + x_2 = 3$$

i.e.

Memo

Date:



$$P_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1^2 + w_2 x_2^2 + b)$$

↓ ↓ ↓
 1 2 -1

decision boundary:

$$z = x_1^2 + x_2^2 - 1 = 0$$

$$x_1^2 + x_2^2 = 1$$

Squared Error Cost function

All not be in logistic function.

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

Memo

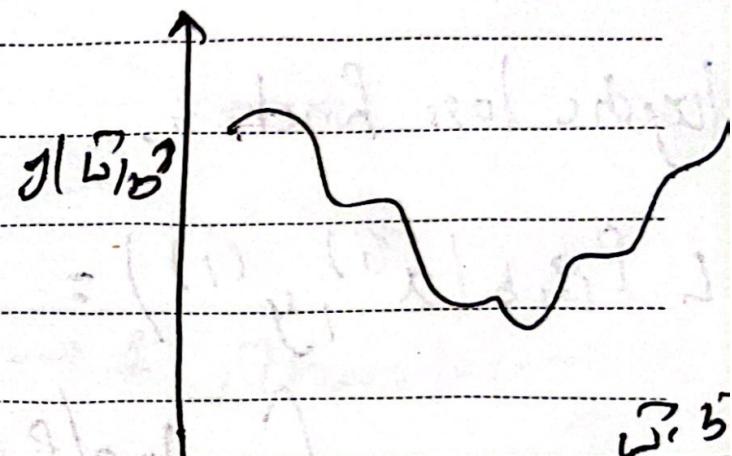
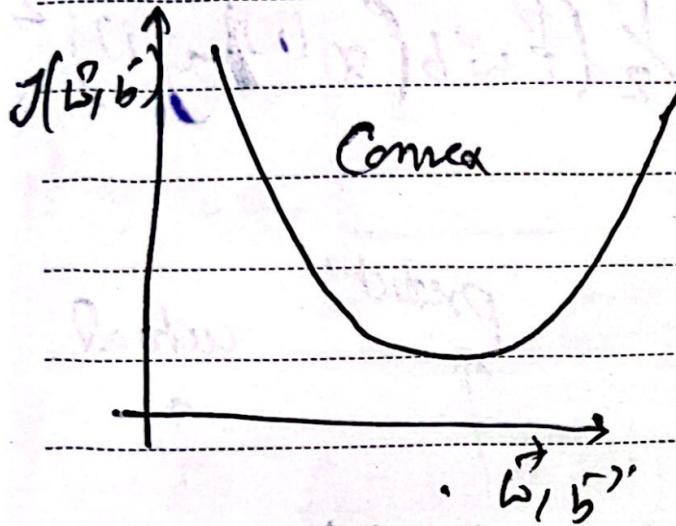
Date: . . .

linear regres.

logistic regres

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

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



This is not easy
curve so,

we redefine cost function as.

#

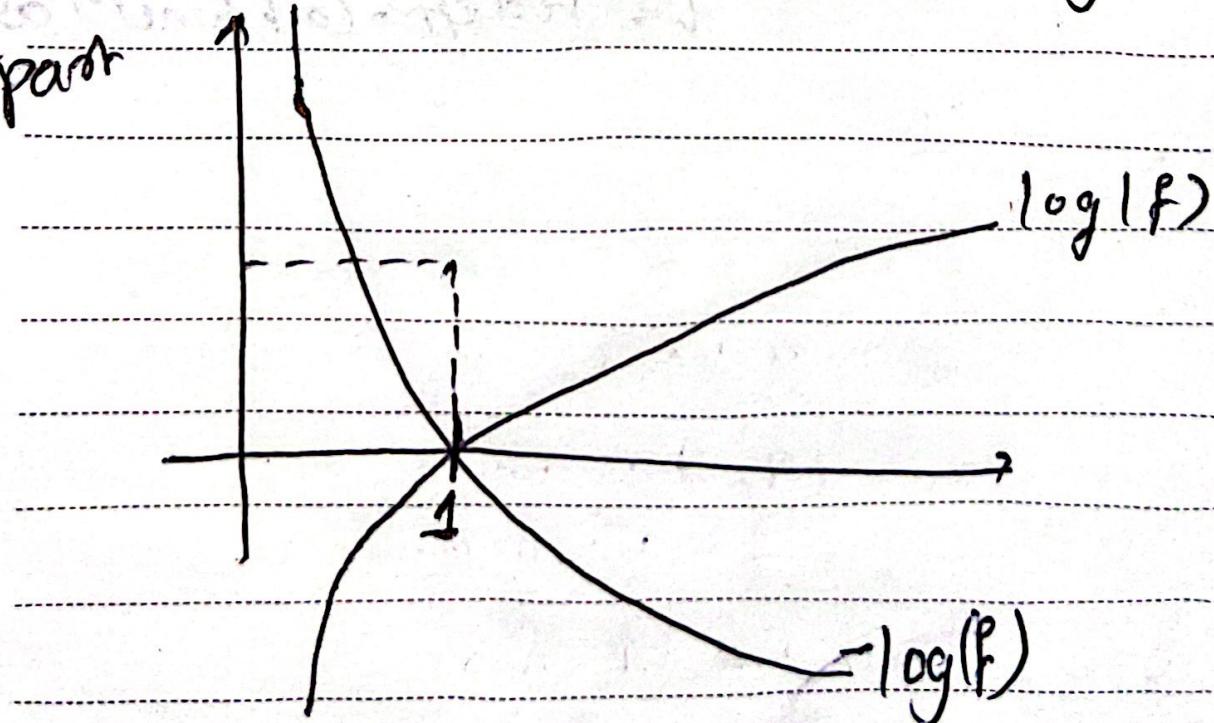
$$\text{loss} \Rightarrow L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) =$$

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

logistic loss function,

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

First part



Memo

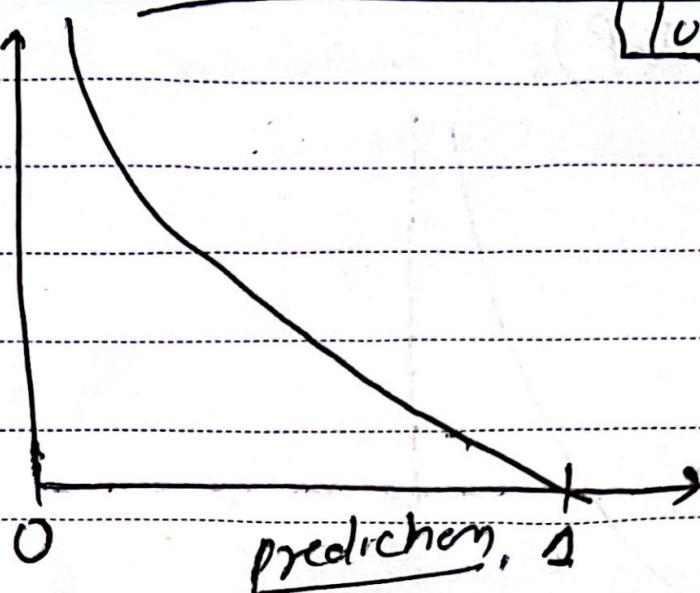
First part

Date:

Loss function

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

$$\text{If } y^{(i)} = 1,$$



As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 1$, Then loss $\rightarrow 0$.

algorithm predicts 1 loss = 0

As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 0$, loss $\rightarrow 1$

prediction
Then.

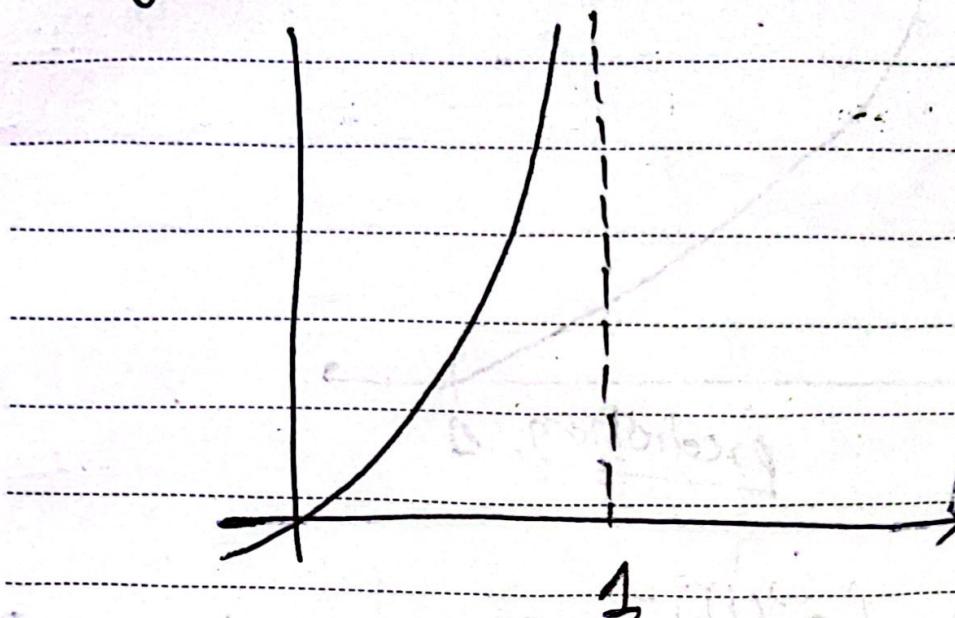
If $y = 1$, and model predicts 0,
then, loss function becomes infinite
i.e. cost fn is infinite.

If $y = 1$ and model predicts 0.8,
loss fn becomes 0.0968
i.e. near to 0.

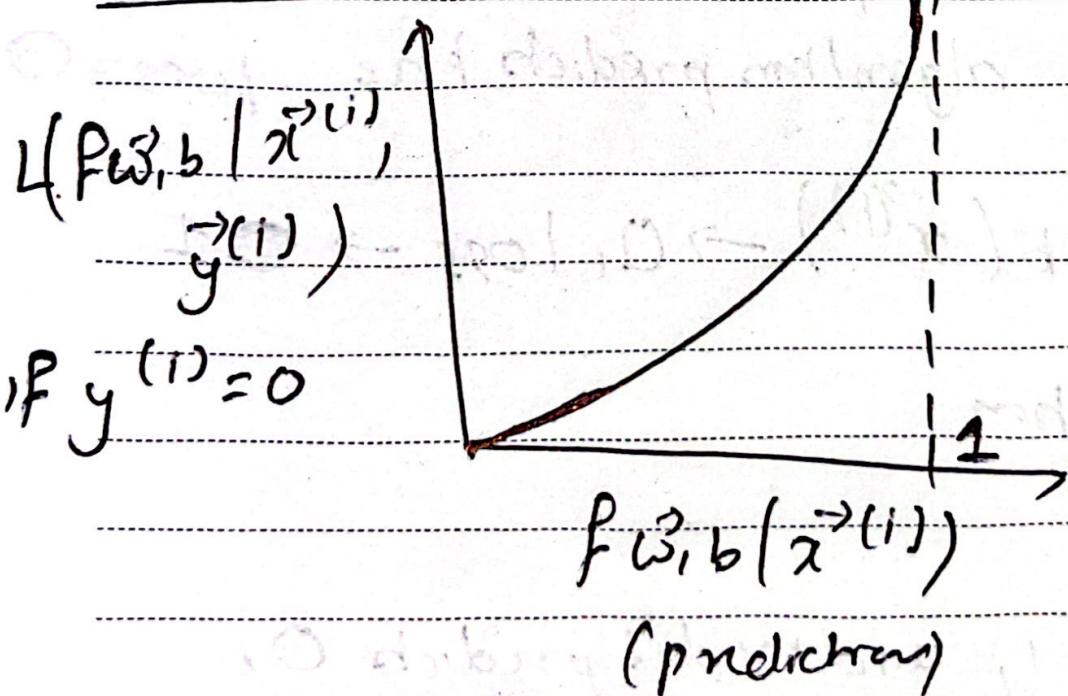
Memo

Date:

$$-\log(\textcircled{2}) \cdot 1 - f$$



zoo mm



i.e. For actual value = 0,
if prediction = 1,

loss $P_1 = \infty$
and cost $P_0 = 0$.

Memo

Date: . . .

But if prediction = 2, for actual val = 1,
loss $P_1 = 0$, sum cost $P_1 = 0$

Thus,

Cost Fn:

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

making
convex,

$$= \begin{cases} -\log(f_{\vec{\omega}, b}(\vec{x}^{(i)})) \\ \text{if } y^{(i)} = 1 \end{cases}$$

$$-\log(1 - f_{\vec{\omega}, b}(\vec{x}^{(i)})) \\ \text{if } y^{(i)} = 0.$$

loss is measure of the difference a single example to its target value while,
cost is measure of losses over training set.

Simplified loss function:

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

Cost

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

Gradient Descent

repeat {

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

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

linear reg: $\vec{f}_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

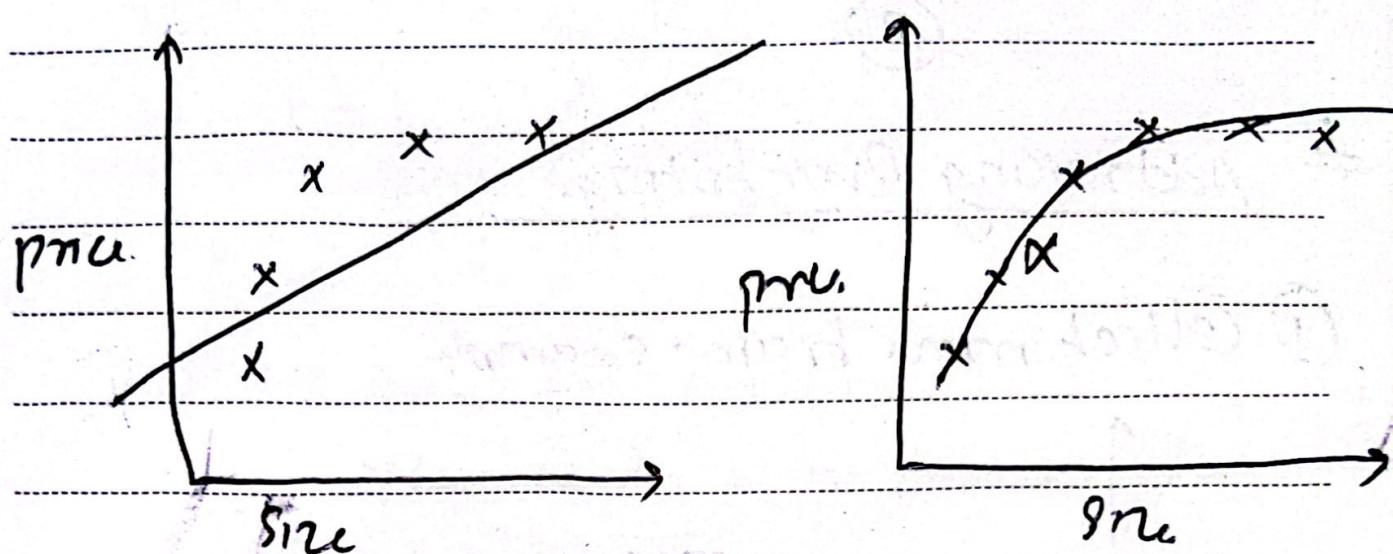
logistic

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

Memo

Date:

Just right



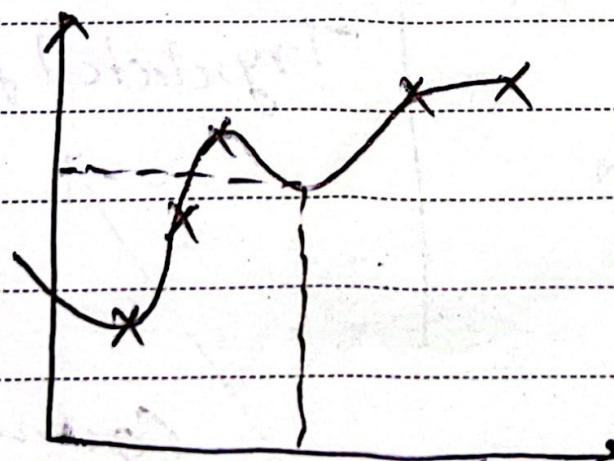
$$w_1 n + b$$

$$w_1 z + w_2 z^2 + b$$

→ Under Fit

High bias. Set well

generalization,
pretty well



$$w_1 n + w_2 n^2 + w_3 n^3 + w_4 n^4 + b$$

Over fit

High variance

fit pretty extremely well.

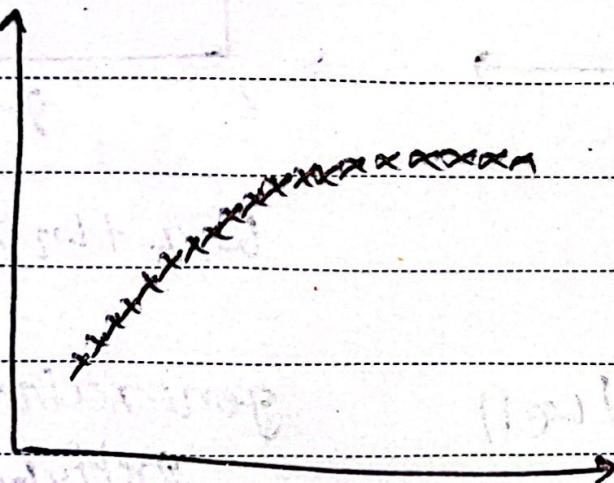
Memo

Date:

→ neither high bias nor high variance
②.

Addressing Overfitting.

① Collect more train example



②. all feature +
insufficient data

↓
overfit.

Try selected features.

Regularization,

Column 2

③. Reduce size of parameters

Regularization to reduce Overfitting

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

If we have 100 features

Then we have 100 parameters.

$$w_1 x_1 + w_2 x_2 + \dots$$

Now Cost fn is, mean sq. Error

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

$$\frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \quad \begin{cases} \text{regularization term.} \\ \text{regularization term.} \end{cases}$$

λ = regularization parameter.

$$\lambda > 0.$$

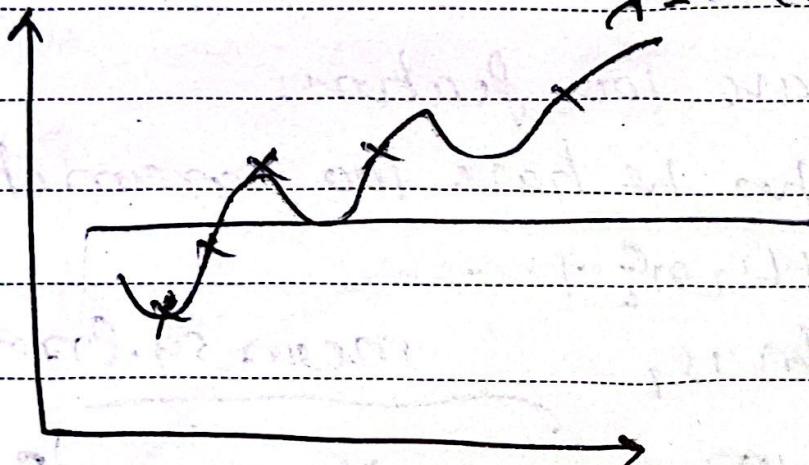
Lesser λ value param.,
to now will be fn.

⑥. is not regularized but some do

If λ is larger, it will be under fit

If λ is small, it will be over fit

$$\lambda = 0 \text{ (OverFit)}$$



{ Increasing λ , reduces overfitting by reducing size of parameters }

Implementing Gradient descent

repeat {

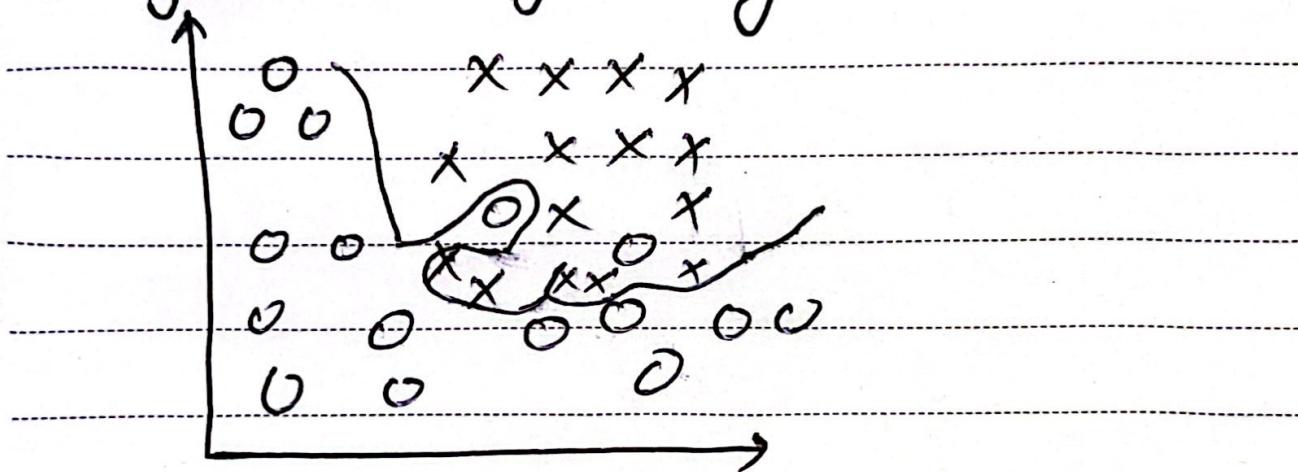
$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m [f_{\beta, b}(x^{(i)}) - y^{(i)}] x_j^{(i)} + \frac{\lambda}{m} w_j \right]$$

Memo

Date: . . .

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

Regularized logistic Regression.



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 +$$

$$w_4 x_1^2 x_2^2 + \dots + b$$

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

Cost J :

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1-y^{(i)}) \log(1-f_{\vec{w}, b}(\vec{x}^{(i)}))] + \frac{\lambda}{2} \sum_j w_j^2$$

Memo

Date:

Graham descent.

Same as logistic.