

# Machine Learning Algorithms

- Supervised ("has 'right answer'")

→ Tools are important  
→ Learning how to apply is even more important.

- Unsupervised

(Given Output Labels)

✗

Supervised

Majorly

Classification

Regression

Predicting a number from infinitely many possible outputs.

Predict finite set of possible output - categories.

✗

Unsupervised

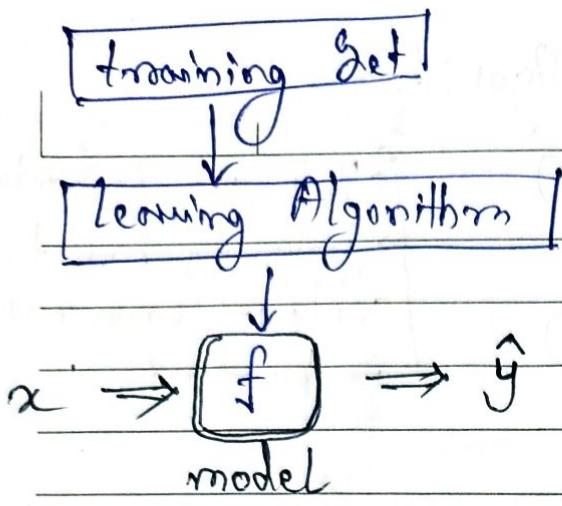
Clustering

Ex: Clustering customers into homogeneous market segments.

❖ Terminologies: (Regression)

i) Training Set: Data used to train the model.

ii)  $\{x^{(i)}, y^{(i)}\}$  single training example  
 $x^{(i)}$  features "input" variables ( $x$ )  
 $y^{(i)}$  target / "output" variable ( $y$ )  
 $m$  No. of training examples  
 $i$  index



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

or,  $f(x) = wX + b$

$w, b$ : Parameters

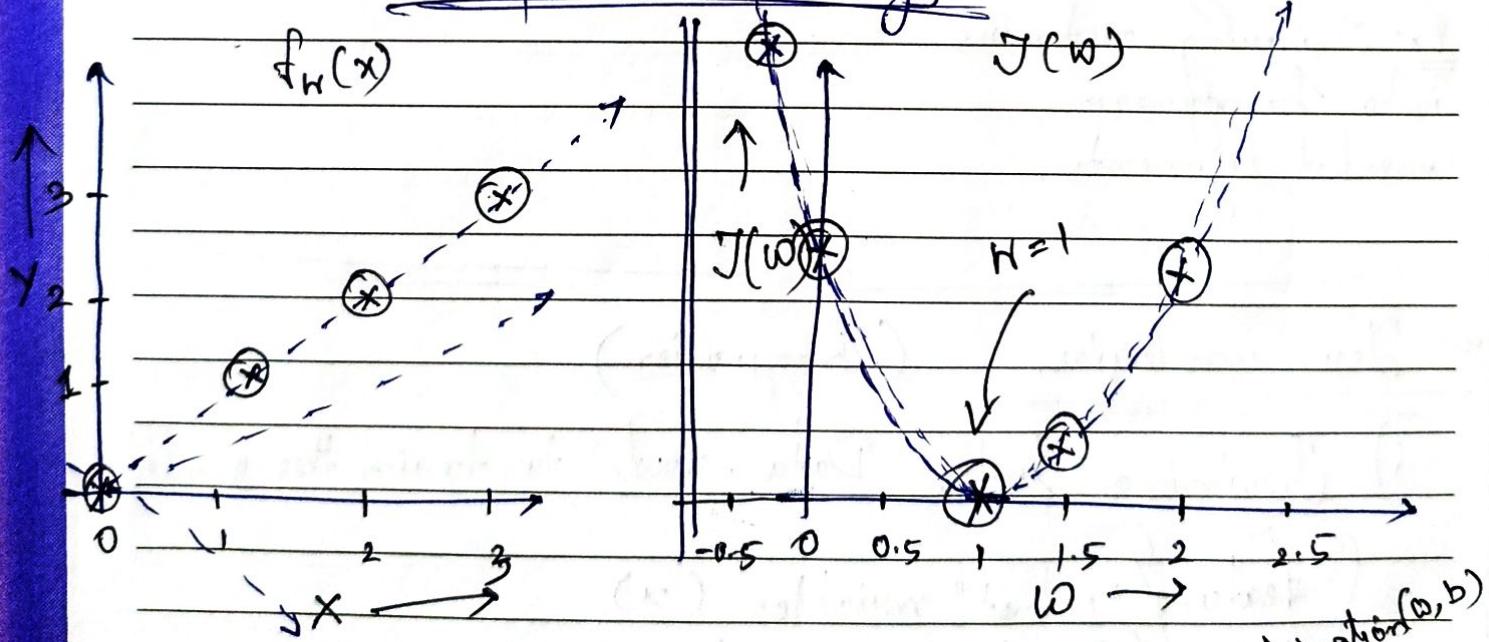
Cost function:

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

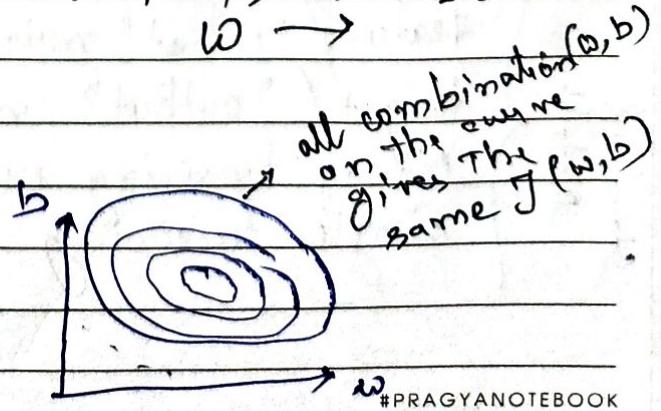
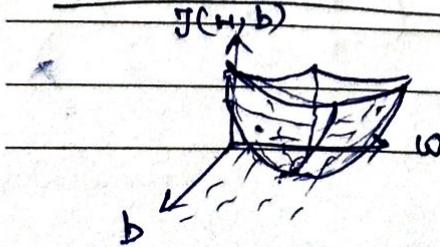
Hence,  $m$  is the No. of training examples.

goal: minimize  $J(w, b)$

Graphical Analogy (say  $b = 0$ )



(\*) Contour Plots



~~repeat until convergence,~~

## Gradient Descent

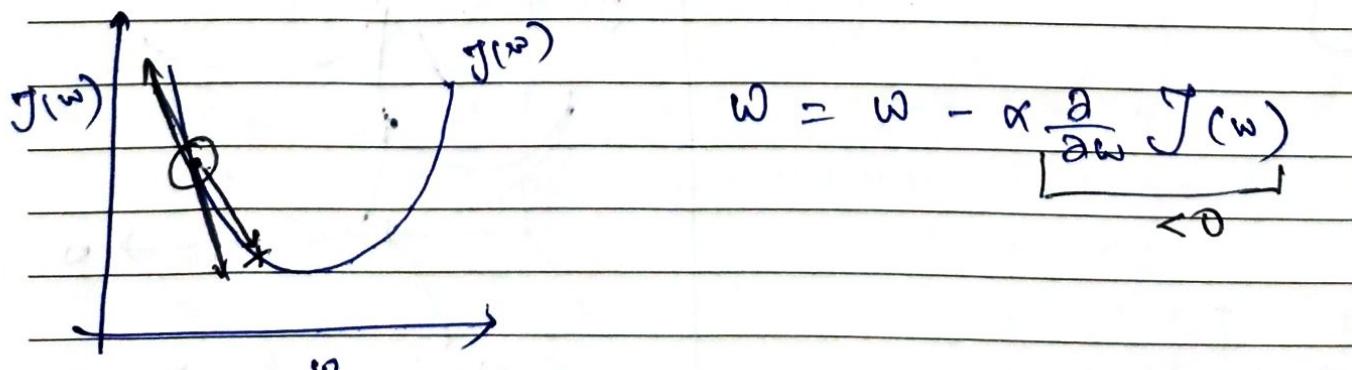
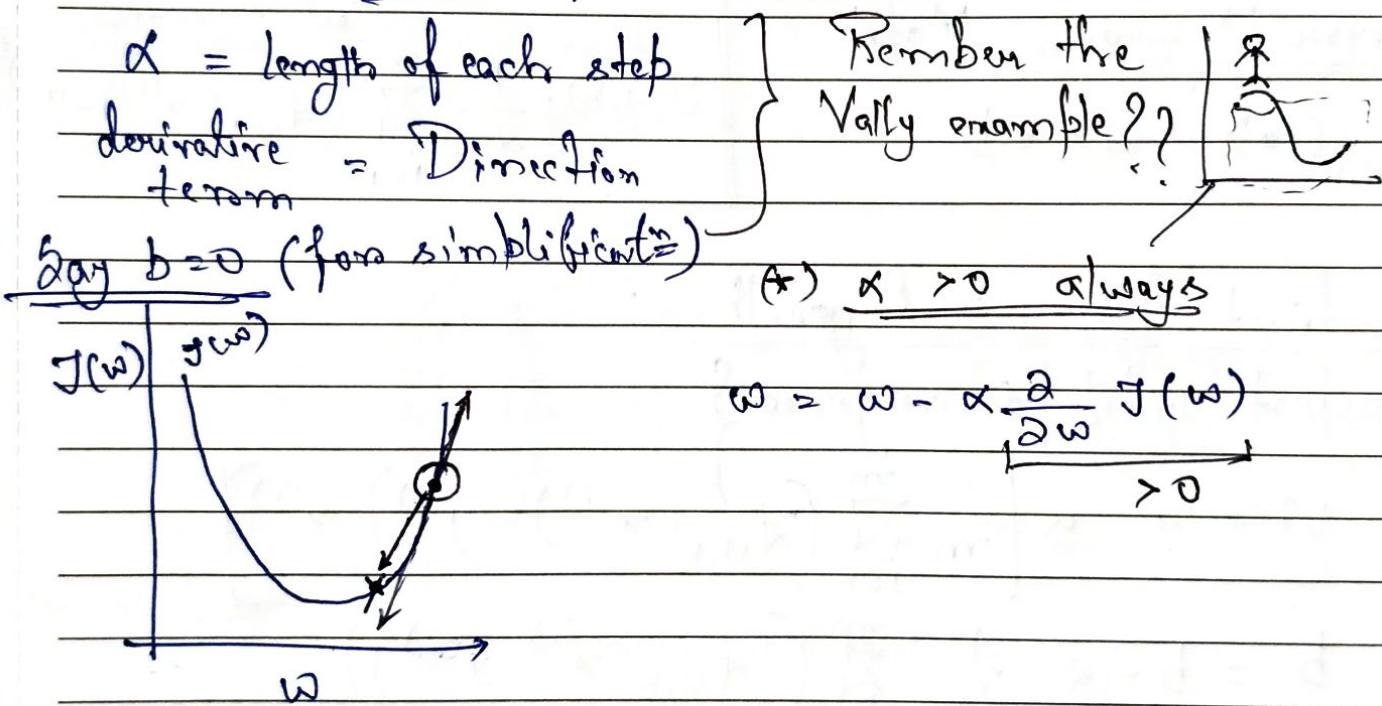
function used to minimize any function

Algorithm: (Simultaneous Update)

$$\left. \begin{array}{l} \text{tmp\_w} = w - \alpha \frac{\partial J(w, b)}{\partial w} \\ \text{tmp\_b} = b - \alpha \frac{\partial J(w, b)}{\partial b} \\ w = \text{tmp\_w} \\ b = \text{tmp\_b} \end{array} \right\}$$

$\alpha = \text{learning rate} (\geq 0)$

## Intuition of Gradient Descent



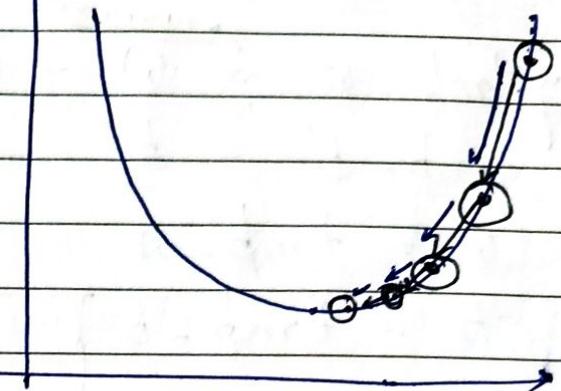
Once reached a local minima  $J(w) = 0$  so,  $w$  will not change.

# Choosing efficient learning rate ( $\alpha$ )

## Near a local minimum:

- Derivative becomes smaller.
- Update step becomes smaller.

$\therefore$  Can reach minimum without decreasing learning rate.



## Linear Regression Model

$$f_{w,b}(x) = w_1 x + b$$

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

## Cost Function

## Gradient Descent Algorithm:

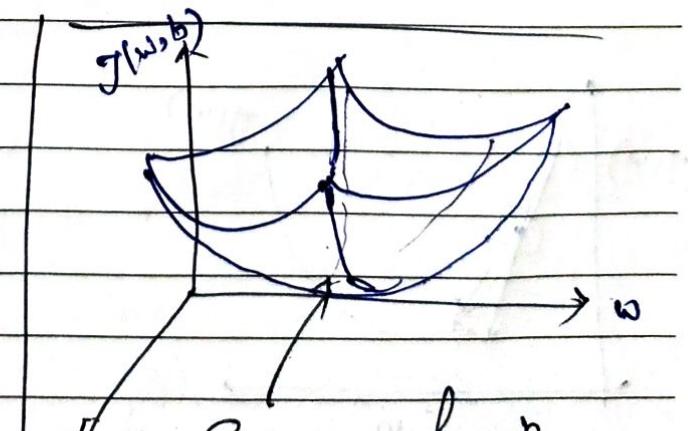
repeat until convergence }

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

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

{}

$$g(w, b)$$



Convex fun  
like this will always converge to a global minimum

## Multiple Features:

$x_j \rightarrow j^{\text{th}}$  feature

$\vec{x}^{(i)} \rightarrow$  features of  $i^{\text{th}}$  training example

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

## New Model Notation:

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

Dot Product

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

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

## Vectorization:

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

shorter code

faster computation  
due to parallel processing

Beyond the  
Scene:

## Gradient Descent

$$\vec{w} = (w_1, w_2, \dots, w_b)$$

$$\vec{d} = (d_1, d_2, \dots, d_b)$$

for  $j$  in range(16):

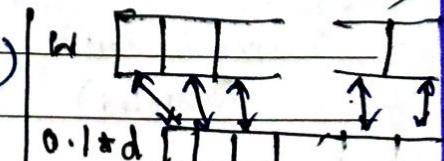
$$w[j] = w[j] - (0.1 * d[j])$$

$$w = np.array([0.5, 1.3, \dots, 2.1])$$

$$d = np.array([0.3, 0.2, \dots, 0.1])$$

$$w = w - (0.1 * d)$$

improves  
16 steps each  
at a time



$\Rightarrow$  parallel multiplication  
+ addition in a single step.

# Gradient Descent

One Feature

repeat {

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

$$\leftarrow \frac{\partial}{\partial w} J(w, b)$$

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

repeat {

$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

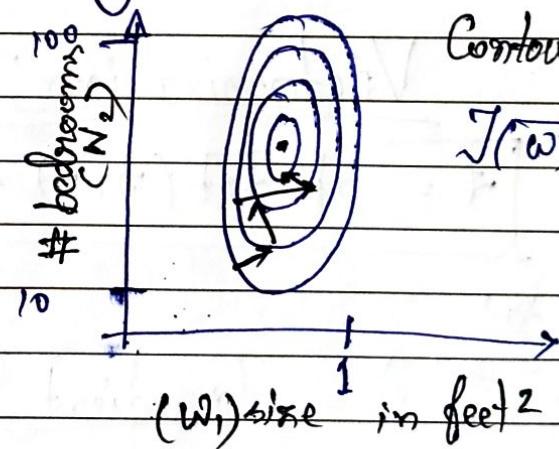
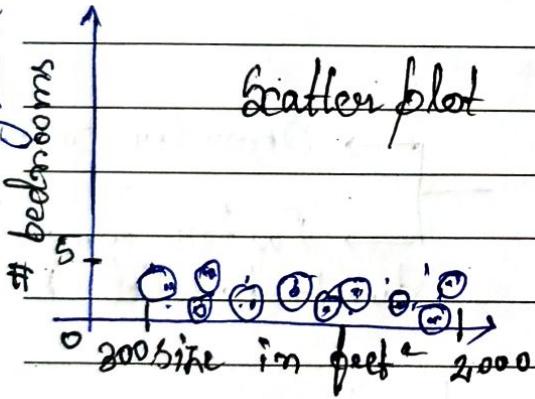
$$\leftarrow \frac{\partial}{\partial w_1} J(\vec{w}, b)$$

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

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

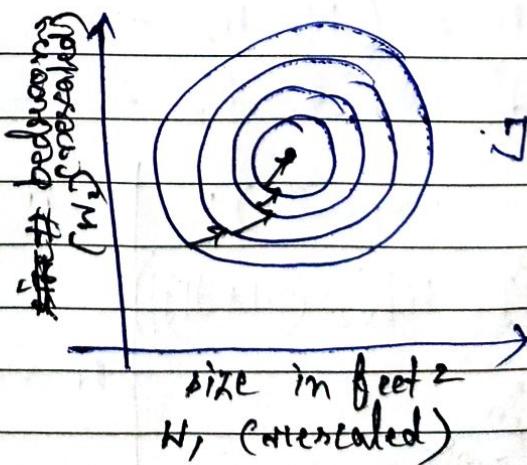
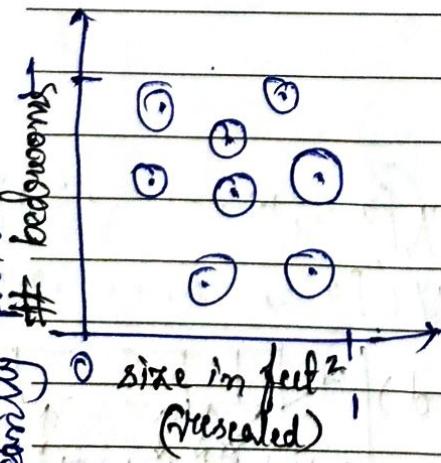
shows down gradient descent  
different ranges of values

Feature scaling



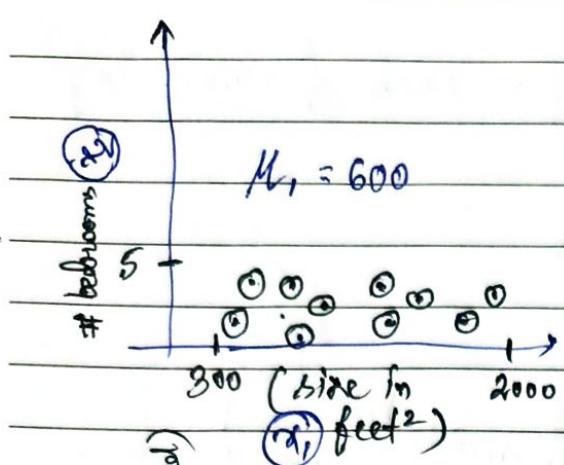
$$J(\vec{w}, b)$$

standardizing features  
significantly reduces  
gradient descent



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

## Mean Normalization



$$300 \leq x_1 \leq 2000$$

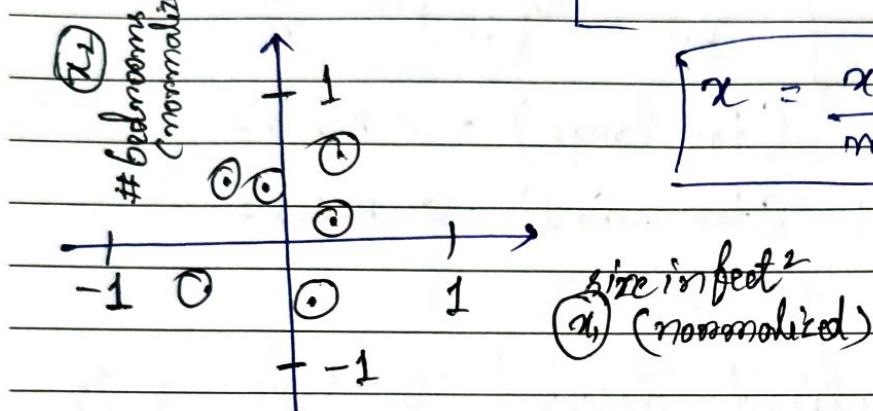
$$x_1 = \frac{x_1 - \mu_1}{2000 - 300}$$

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

$$0 \leq x_2 \leq 5$$

$$x_2 = \frac{x_2 - \mu_2}{5 - 0}$$

$$-0.46 \leq x_2 \leq 0.59$$



$$x = \frac{x - \text{mean}}{\text{max} - \text{min}}$$

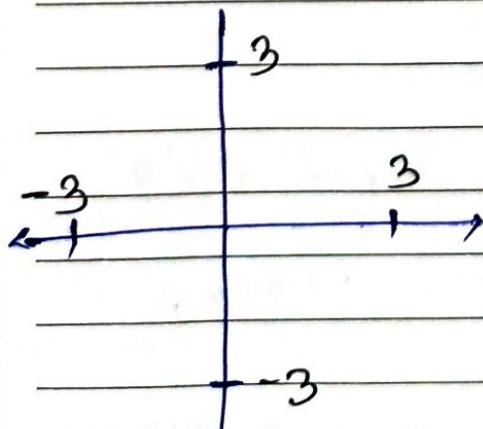
## Z-score Normalization

$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$$-0.67 \leq \frac{x_1 - \mu_1}{\sigma_1} \leq 3.1$$

$$-1.6 \leq \frac{x_2 - \mu_2}{\sigma_2} \leq 1.9$$

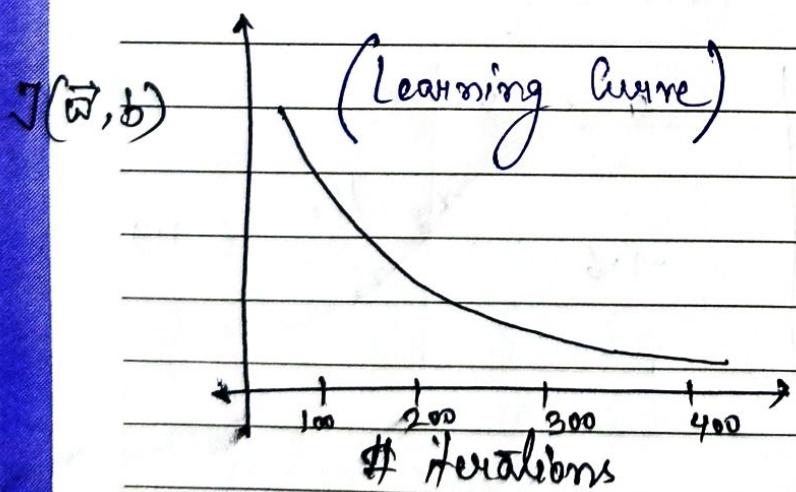


## Where to scale & where not needed

Aim for about acceptable ranges	$-1 \leq x_j \leq 1$ for each feature $x_j$
	$-3 \leq x_j \leq 3$
	$-0.3 \leq x_j \leq 0.3$
	✓
$0 \leq x_1 \leq 3$	$\left\{ \begin{array}{l} \text{not need} \\ \text{but still} \end{array} \right.$ <span style="border: 1px solid black; border-radius: 50%; padding: 2px;">not train in rescaling</span>
$-2 \leq x_2 \leq 0.5$	
$-100 \leq x_3 \leq 100$ (too large)	$\rightarrow$ Rescale
$-0.001 \leq x_4 \leq 0.001$ (too small)	$\rightarrow$ Rescale

Is your gradient descent converging??

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



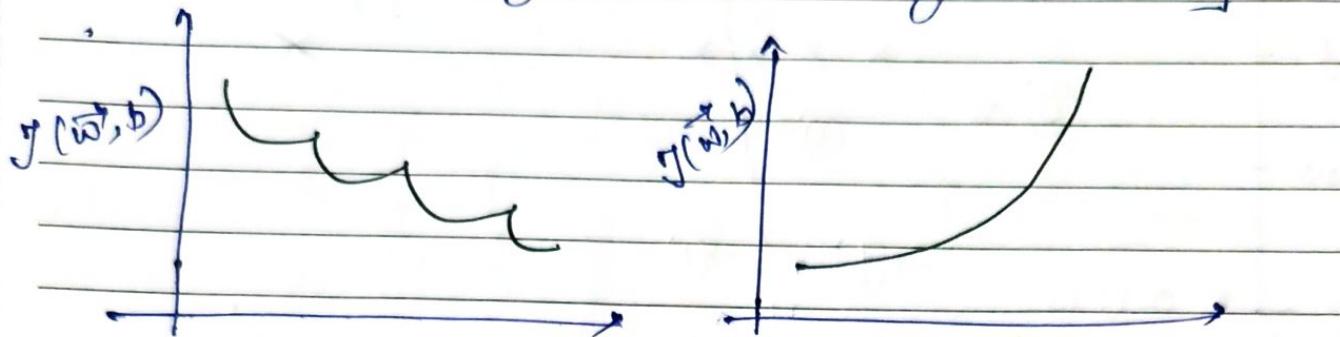
- \* If  $J(\vec{w}, b)$  ever increases
  - $\Rightarrow$  i)  $\alpha$  (learning rate) is wrongly chosen.
  - ii) Bug in the code.

After: Automatic convergence test, def,  $\epsilon = \alpha'$   
 if  $J(\vec{w}, b)$  decreases by  $\leq \epsilon$  in one iteration,  
 declare convergence.

## Choosing learning rate:

Problematic learning rate:

- [ i)  $\alpha$  is too large ]
- [ ii) Bag in the code ]



~~(\*) Choose smaller  $\alpha$  and that may take more time to converge but has to decrease Cost function ( $J(w, b)$ ) to decrease in each iteration.~~

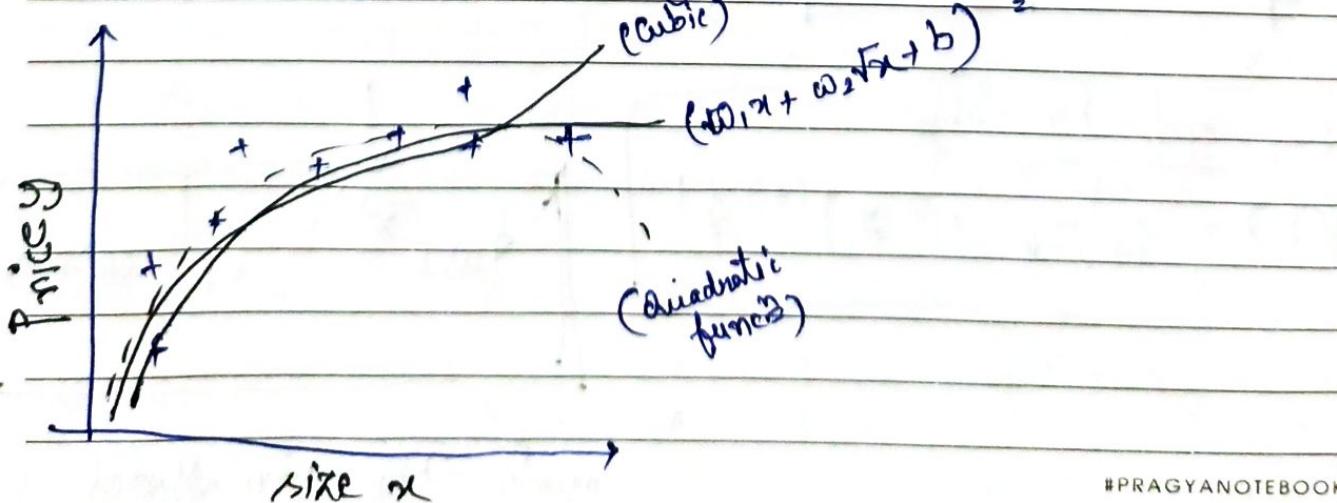
## Feature Engineering (making custom features)

$$f_{\vec{w}, b}(\vec{x}) = \underbrace{w_1 x_1}_{\text{length}} + \underbrace{w_2 x_2}_{\text{width}} + b$$

$$x_3 = x_1 x_2 \quad (\text{new feature})$$

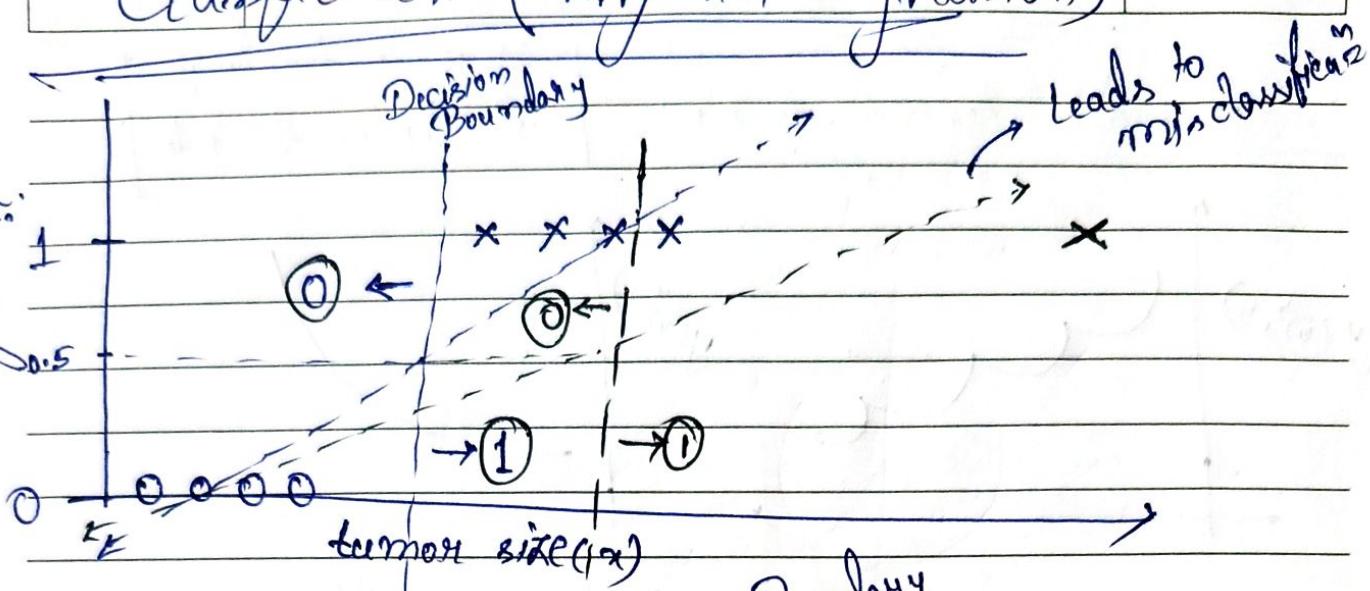
$\hookrightarrow$  (area)

Polynomial regression



# Classification (Why not Regression)

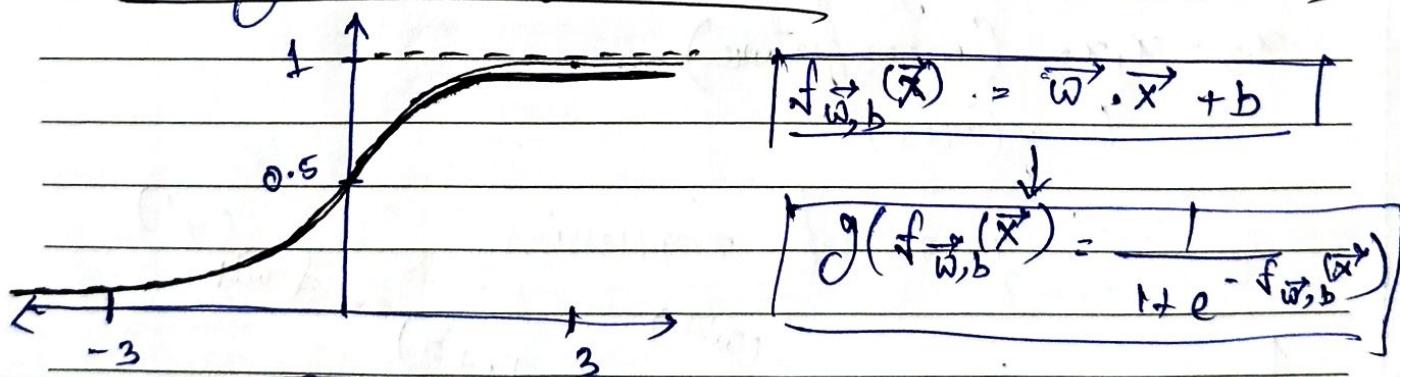
malignant??



$$\begin{aligned} \text{if } f_{\vec{w}, b}(\vec{x}) & \leq 0.5, \hat{y} = 0 \\ \text{if } f_{\vec{w}, b}(\vec{x}) & \geq 0.5, \hat{y} = 1 \end{aligned}$$

Decision Boundary  
(shifted)

Logistic Regression (~~not used for classification~~)



Sigmoid Function

$$g(x) = \frac{1}{1 + e^{-x}}, 0 \leq g(x) \leq 1$$

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

" Logistic Regression Model "

# Interpretation of Logit Output

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

"Probability" that class is 1

alternatively,  $f_{\vec{w}, b}(\vec{x}) = P(y=1 | \vec{x}; \vec{w}, b)$

Probability that  $y$  is 1, given  
input  $\vec{x}$ ; parameters  $\vec{w}, b$

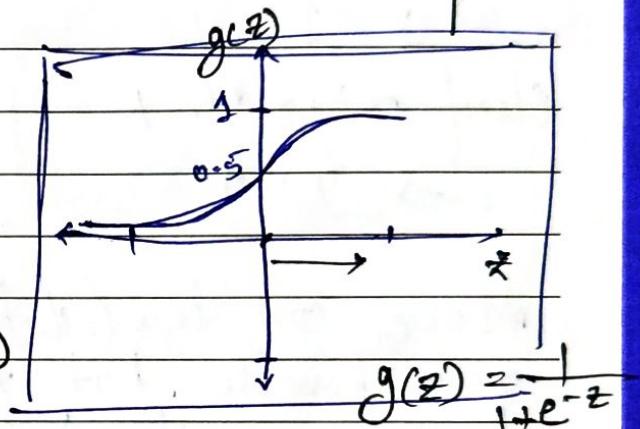
Q:  $f_{\vec{w}, b}(\vec{x}) \geq 0.5$  ?

Yes:  $\hat{y} = 1$

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

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

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



Hence, algorithm is,

if  $\vec{w} \cdot \vec{x} + b \geq 0.5$ :  $\hat{y} = 1$

elseif  $\vec{w} \cdot \vec{x} + b < 0.5$ :  $\hat{y} = 0$

~~elseif~~

Decision Boundary

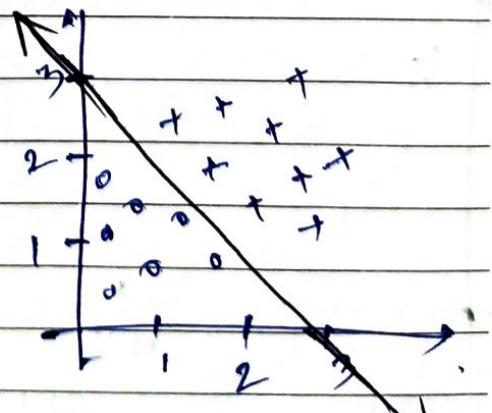
Consider,  $w_0 = 1; w_1 = 1, w_2 = 3$

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_0 x_0 + w_1 x_1 + w_2 x_2 + b)$$

Decision Boundary:  $\vec{z} = \vec{w} \cdot \vec{x} + b = 0$

$$\vec{z} = x_1 + x_2 + -3 = 0$$

$$\Rightarrow x_1 + x_2 = 3$$



ad,  $g(\vec{z}) = 0$

## Non-linear Decision Boundaries:

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

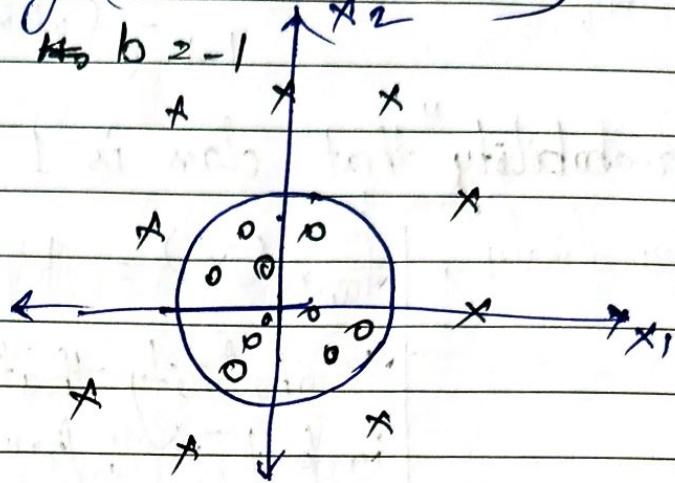
Consider,  $w_1 = 1, w_2 = 1, b = -1$

Decision Boundary,

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

$$\begin{aligned} \text{if } x_1^2 + x_2^2 \geq 1 \\ \Rightarrow y = 1 \end{aligned}$$

$$\begin{aligned} \text{elseif } x_1^2 + x_2^2 < 1 \\ \Rightarrow y = 0 \end{aligned}$$



More complex (ellipses, etc.) decision boundaries are possible, but for linear model its decision boundary will always be linear.

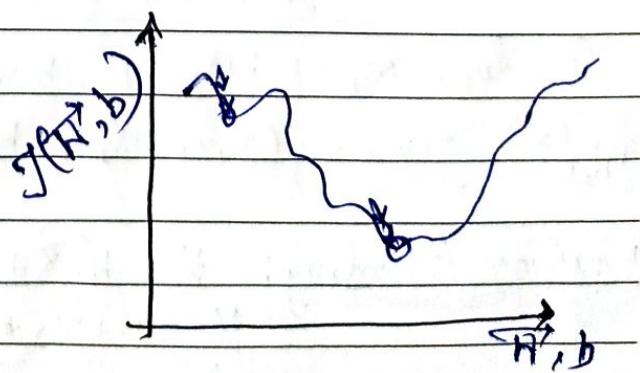
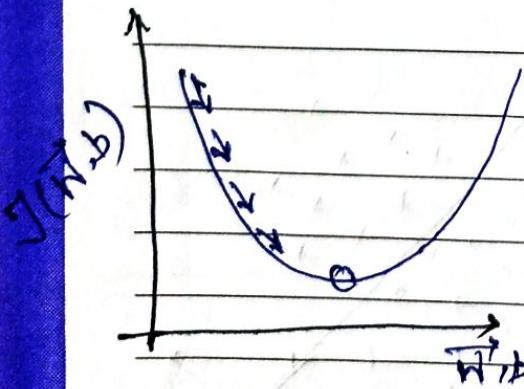
## Cost Function

(Why Squared Error cost function is not the best for Logit Model)

Linear Regression:

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

$$\begin{cases} \text{Logistic Regression} \\ f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}} \end{cases}$$



Convex

→ Global  
 $\min \Sigma = \text{local} \min \Sigma$

Non-Convex

ANOTEBOK

Cost function for Linear Regression:

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

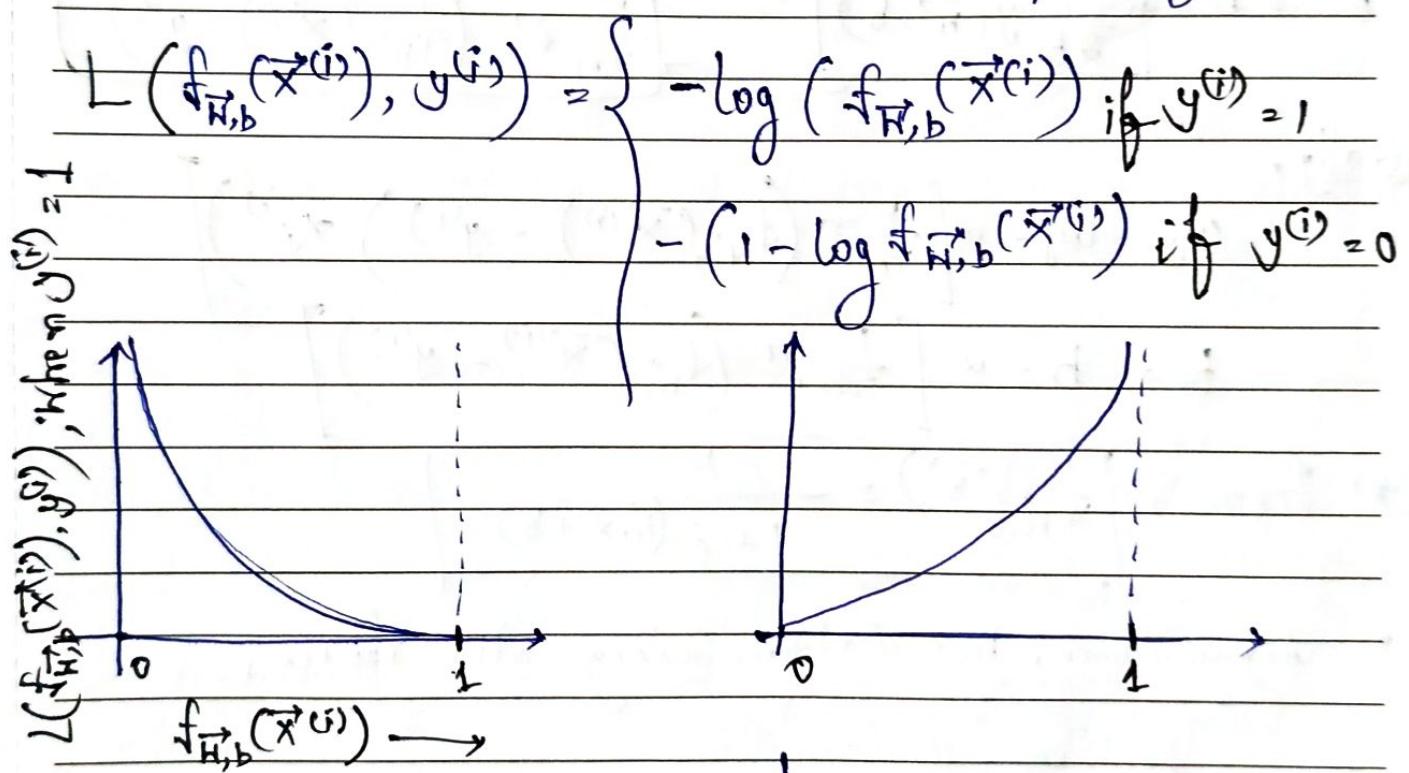
frame,  
 $f_{\vec{w}, b}(\vec{x}^{(i)})$   
 $= \vec{w} \cdot \vec{x} + b$

$$\rightarrow \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f(\vec{x}^{(i)}) - y^{(i)})^2$$

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

→ Called loss funct<sup>n</sup>

Now let's construct loss-function for logit-mode:



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

Simplified Cost function:

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

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

But why particularly this function??

$\Rightarrow$  Maximum Likelihood Estimation.

## Gradient Descent

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

Repeat {

$$w_j = w_j - \alpha \left[ \frac{\partial}{\partial w_j} J(\vec{w}, b) \right]$$

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

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

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

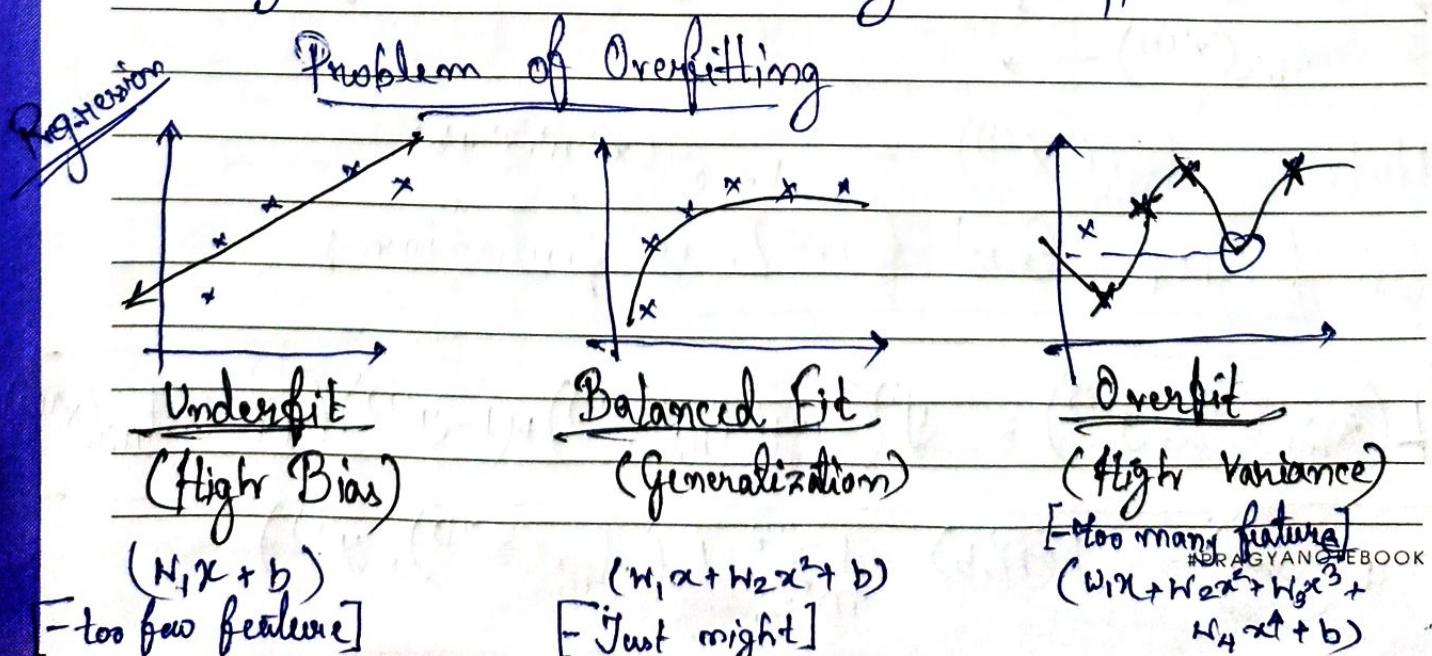
finally,

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

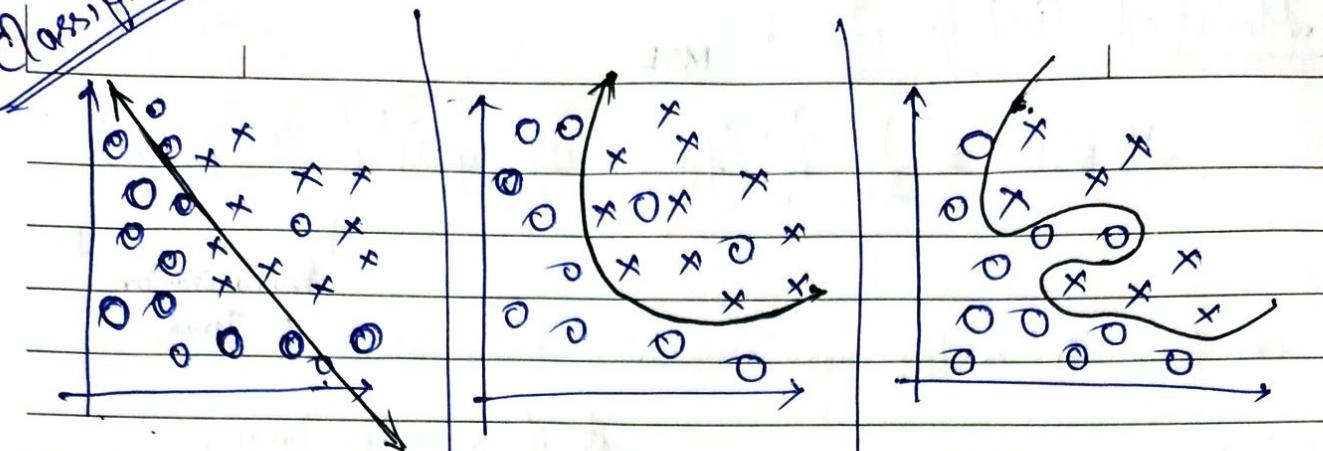
$$b = b - \alpha \left[ \frac{1}{m} \sum (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \right]$$

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

N.B. Convergence, Vectorization, Scaling still applies.



# Classification



$$\Rightarrow \hat{y} = w_0 x_0 + w_1 x_1 + b$$

$$\Rightarrow f_{w,b}(x) = g(\hat{y})$$

where,

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

(High Bias)

$$\hat{y} = w_0 x_0 + w_1 x_1 + w_2 x_2$$

$$+ w_3 x_1^2 + w_4 x_2^2$$

$$+ w_5 x_1 x_2 + b$$

(Balance fit,  
generalized)

$$\hat{y} = w_0 x_0 + w_1 x_1 +$$

$$w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2$$

$$+ w_5 x_1^2 x_2^3 + w_6 x_1^3 x_2$$

$$+ \dots + b$$

C

~~Overfitting~~

~~Regularization~~

~~Solutions~~

~~Overfitting~~

i) Data Augmentation:

→ Use more data to fit the model.

ii) More features + Insufficient data = Overfitting

→ Use fewer features (Feature Selection)

iii) Regularization (regulating the impact of features)

N.B.: Simpler models usually don't overfit.

b) Penalizing all the features usually results in a more effective model less prone to overfitting.

c) Penalizing specific features using its cost function can also be done!

## Modified Cost Function:

MSE

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

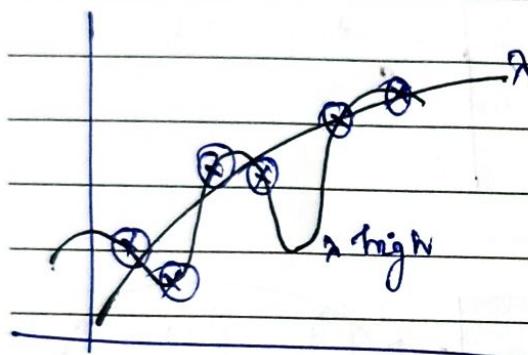
Regularization Term

(\*) here ' $\lambda$ ' is the regularization parameter.

Note:

$$\lambda \uparrow \Rightarrow \downarrow w_i's$$

if  $\lambda$  is too large,  
model shrink to a linear  
one, if  $\lambda = 0$ , will  
result in a wiggly one  
(having no regularization effect)



## Modified Gradient Descent to

Derivative term is

$$\frac{\partial}{\partial w_j} (J(\vec{w}, b)) = \frac{1}{m} \sum_{i=1}^m (\hat{f}_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$$

Simultaneous Update:

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

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

Note, here

$$\hat{f}_{\vec{w}, b}(\vec{x}^{(i)}) = \vec{w} \cdot \vec{x}^{(i)} + b$$

similarly,

Regularized Logistic Regression follows --

with distinction in the functional defn.

# Stat - ML Playlist (codebasics)

## → Univariate Linear Regression:

(from sklearn import linear\_model)

library

» lin-reg = ~~linear\_model.LinearRegression()~~

» lin-reg.fit([[1, 1], ])

» lin-reg.predict([[1, 1]])

» lin-reg.coef\_ # fitted parameter value

» lin-reg.intercept\_ # ,

Scatter plot

import matplotlib.pyplot as plt

» plt.scatter()

## → Multiple Linear Regression

\* » df[...] # Creates pandas.series object.

\* » df[[...]] # creates pandas.datframe object.

» df.fillna(val) # fill up NA values with val given.

» W2n.word\_to\_num() # converts words to number

» maf('fn', na\_action = "ignore") # ignores na values.

### (\*) Instantiating LinearRegression object

» [lin-reg = linear\_model.LinearRegression()]

$$\# \hat{y}^{(i)} = w^{(i)}x + b$$

Algorithmic Parameters

→ Iteration

→ Learning rate

→ Gradient Descent :

$w = p \geq 1$ ,  $m = \text{len}(x)$ , cost

def grades(x, y, iter=100, l=0.001) :

for i in range(iter) :

$$y_p = (w \cdot x) + b$$

$$W = (\frac{1}{m}) * \text{np.sum}((y_p - y) * x)$$

$$B = (\frac{1}{m}) * \text{np.sum}(y_p - y)$$

$$w = w - (l * W)$$

$$b = b - (l * B)$$

Numpy

» np.array()

» np.sum()

» Broadcasting

ndarray \* scalar  
= ndarray

→ Saving and Loading a model: (pickle/joblib/cloudpickle)

import pickle

with open("file-name", "wb") as f:

pickle.dump(model-name, f)

with open("file-name", "rb") as f:

pick m; = pickle.load(f)

import joblib

joblib.dump(model, "file")

joblib.load("file")

m2 = joblib.load("file")

~~manually~~ → Dummy variable:

» pd.get\_dummies(df["col-name"],  
dtype='int')

» model.score(X, y) → ind. training set

→ dependent labels

» from pandas import

scatter\_matrix,

scatter\_matrix(df[[ ]], figsize=c))

Colors [Red Blue Green]

Red 1 0 0

Red 1 0 0

Blue 0 1 0

Green 0 0 1

→ dummy variables  
To avoid multi collinearity we need to drop one of the dummy variable columns

fitted model

~~label encoder  
OneHotEncoder~~

```
>> from sklearn.preprocessing import LabelEncoder  
>> leben = LabelEncoder() # Instantiation  
>> leben.fit_transform(dt[:, 1])
```

# Column Transformer

```
from sklearn.compose  
import ColumnTransformer
```

```
ct = ColumnTransformer([("name", OneHotEncoder(), [col_num])],  
                      remainder = 'passthrough')
```

```
x = ct.fit_transform(x1)
```

df.column  $\Rightarrow$  Dataframe

df.column.values  
 $\Rightarrow$  np.array

$\rightarrow$  train-test-split:

```
from sklearn.model_selection import train_test_split
```

```
X-train, X-test, y-train, y-test
```

```
= train_test_split(X, y, test_size = 0.2,
```

```
random_state = value)
```

$\rightarrow$  Logistic Regression [Prediction of Categories]

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

```
>> from sklearn.linear_models import LogisticRegression
```

```
>> model = LogisticRegression()
```

```
>> model.score(X-test, y-test)
```

Confusion matrix:

```
>> from sklearn.metrics import confusion_matrix
```

```
>> cm = confusion_matrix(y-test, model.predict(X-test))
```

## B. Visualization:

\* import matplotlib.pyplot as plt  
 import seaborn as sns

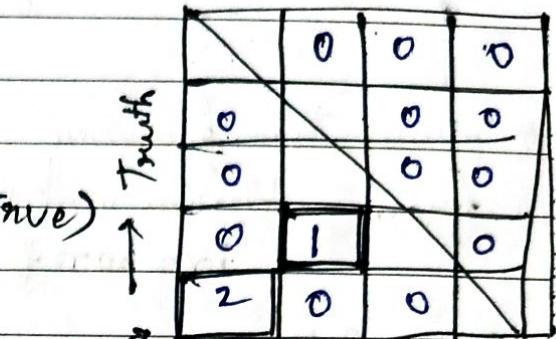
%matplotlib inline

plt.figure(figsize=(7,5))

sns.heatmap(cm, annot=True)

plt.xlabel("Predicted")

plt.ylabel("Truth")



Truth was something and predicted something else.

## Decision Tree:

» from sklearn.tree import DecisionTreeClassifier

(\*) Only label encoding suffices, one Dummy variable not needed.

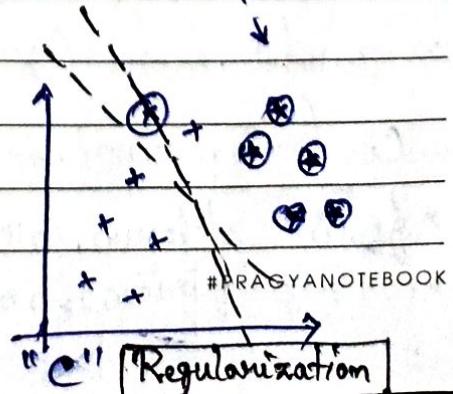
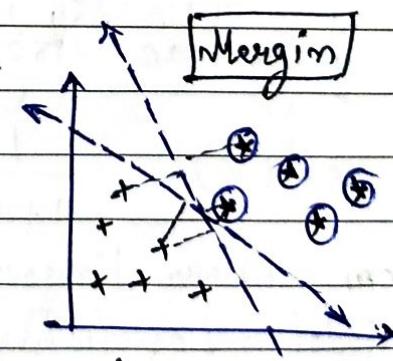
» from sklearn.tree import DecisionTreeClassifier

## Support Vector Machine

» from sklearn.svm import SVC

# Instantiation through Constructors

sv = SVC(c= , kernel='')



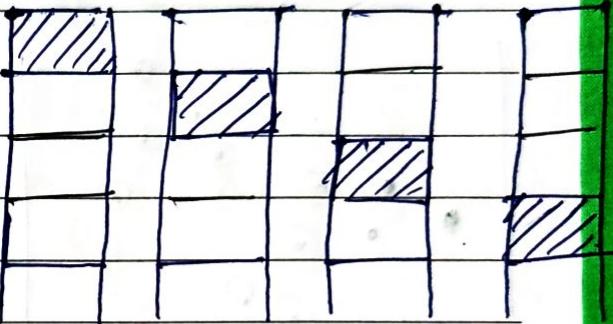
→ Random Forest (i) Generating random decision trees  
ii) Take majority vote

↳ from sklearn.ensemble import RandomForestClassifier

(\*) Based on Bagging + feature sampling

What ML model  
to use?

→ K-fold Cross Validation:



Parameter Vs hyperparameter

i) Weights in neural network	ii) Learning rate, no. of iterations in neural network.
iii) learned from data	iv) To be set heuristically by the practitioner.
v) intercept, coefficients of Linear Regression model.	vi) fit_intercept : bool = True, is a hyperparameter in LinearRegression Model.

Extraction of rows using index:

df.iloc [index]

Manual cross-K-fold:

from sklearn.model\_selection import StratifiedKFold // KFold  
sk = StratifiedKFold (n\_splits = )

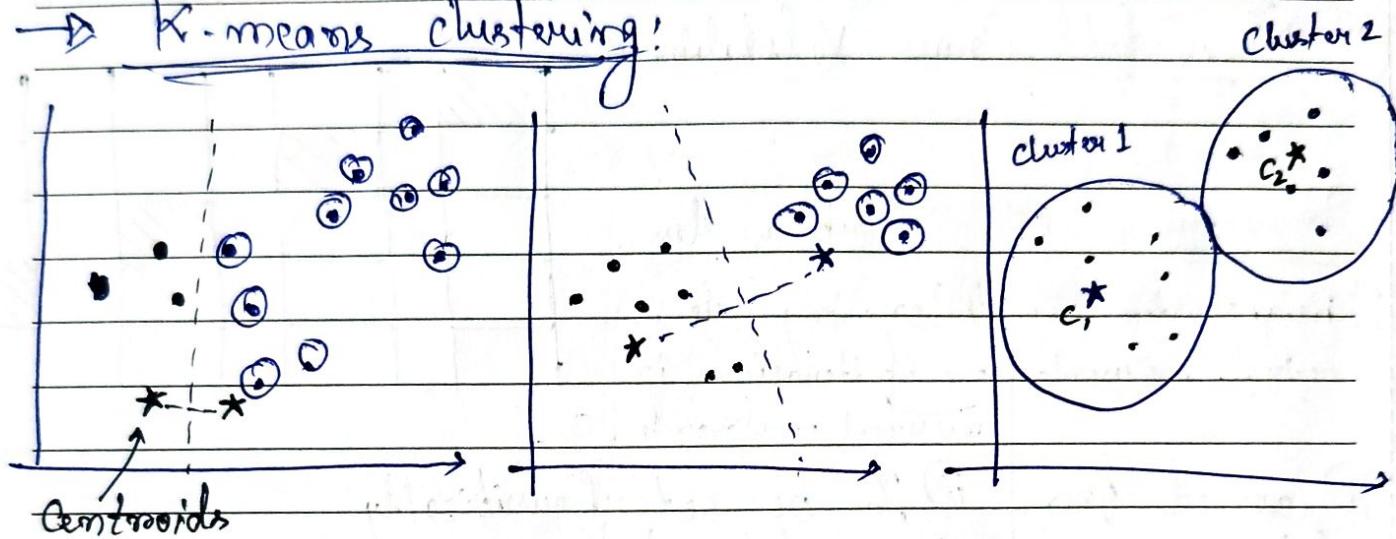
for i, (train\_index, test\_index) in

enumerate (sk (X, y)) :

Cross-validation direct implementation

from from sklearn.model\_selection import  
cross-validation-score

→ # of clusters  
→ K-means clustering!



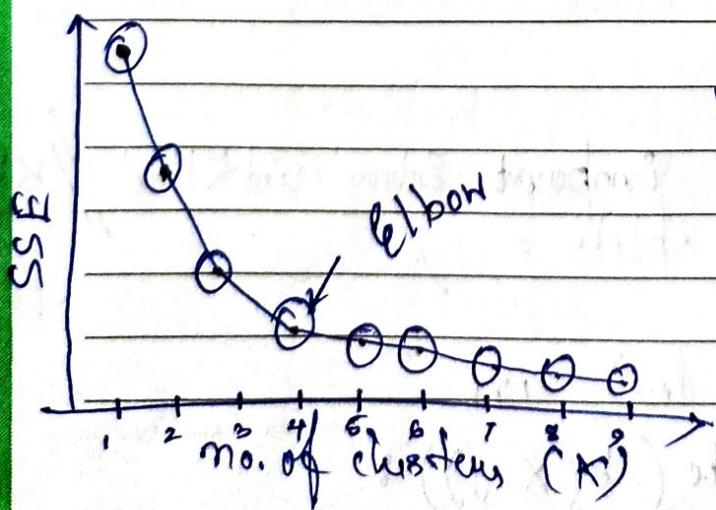
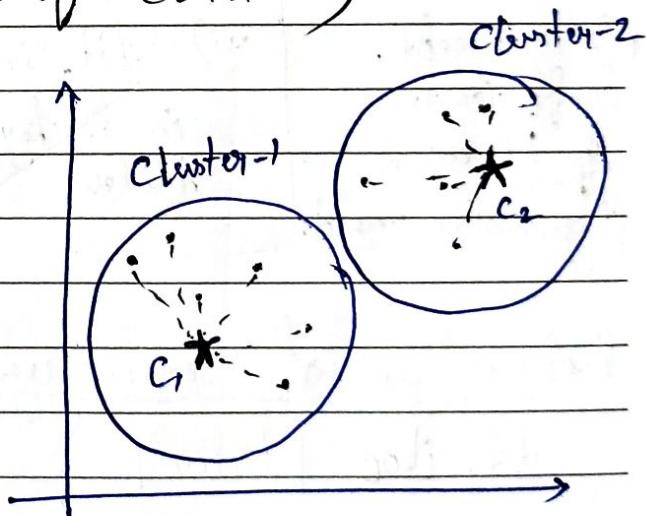
Determination of K (no. of clusters)

In Elbow Method:

Sum of squared errors:

$$SSE_1 = \sum_{i=1}^n \text{dist}(x_i - c_1)^2$$

$$SSE = \sum_{i=1}^k SSE_i$$



\* Feature scaling matters a lot  
from sklearn.preprocessing import MinMaxScaler

sc = MinMaxScaler()

sc.fit(df[[" "]])

sc.transform([ " " ]) = sc.fit\_transform([ " " ])

## KMeans Implementation:

```
from sklearn.cluster import KMeans
```

```
km = KMeans(n_clusters = )
```

```
df["cluster"] = km.fit_predict(df)
```

## below Methods:

k-range = range(1, 10)

SSB = [ ]

for i in k-range:

```
km = KMeans(n_clusters = i)
```

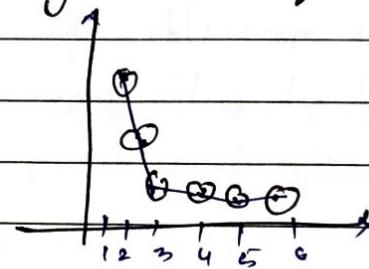
```
km.fit(df)
```

```
sse.append(km.inertia_) # gives sum of squared errors
```

import matplotlib.pyplot as plt

```
plt.scatter(k-range, sse)
```

```
plt.plot(k-range, sse)
```



## Näive Bayes:

(continuous classes)      (discrete classes)

```
from sklearn.naive_bayes import GaussianNB/MultinomialNB/
```

Pipeline → transform all columns

BernoulliNB → (Binary class)

```
from sklearn.pipeline import Pipeline
```

Column - Transformer

```
steps = [ ("est", Estimator()),
```

```
    ("clf", ColumnTransformer(attribute (column names)),
```

("tgt", Predictor())]

], remainder = "ignore"/  
(\*) transforms only specified columns. "booleans"

↳ should have the  
transform() method.

Estimator  
(fit() method)

Transformer  
(Estimator + transform()  
method)

Predictor  
(transformer +  
predict())

## Grid SearchCV & Random SearchCV

```
from sklearn.model_selection import  
GridSearchCV
```

```
from sklearn.model_selection import  
RandomSearchCV
```

```
# suppress warnings  
import warnings  
warnings.filterwarnings("ignore")
```

Consume more time

```
gcv = GridSearchCV( estimator = SVC() ,  
param_grid = { 'kernel': ('linear', 'rbf') ,  
'C': [1, 10] } , cv = 7 )
```

```
gcv.fit(X,y)
```

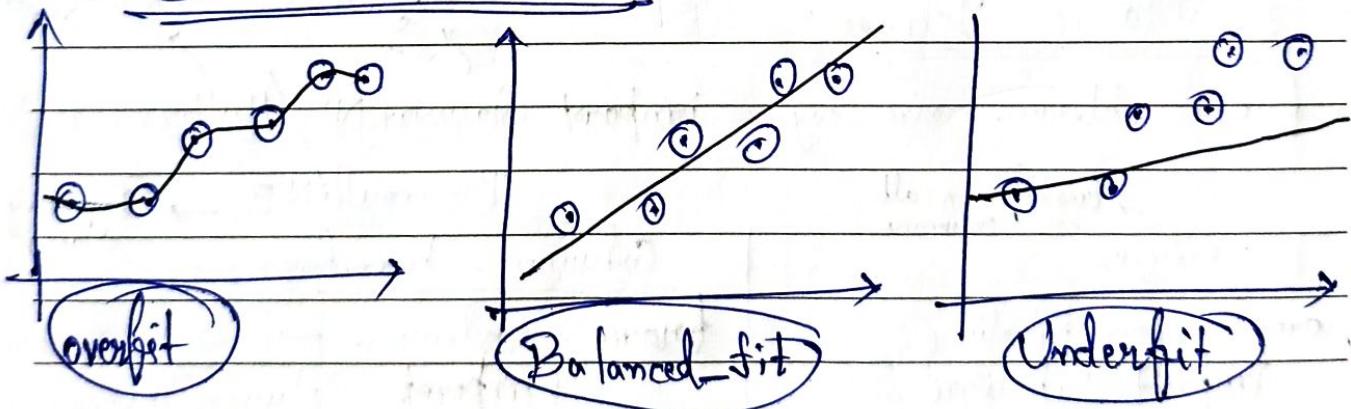
```
gcv.cv_results
```

```
gcv.best_estimator
```

```
frcv = RandomSearchCV  
( - . (+) , n_iter = 25 )
```

less time consumed

## Bias & variance:

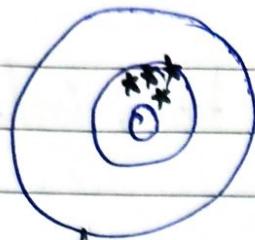


► Bias & Train error.

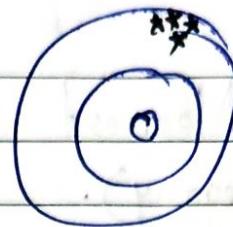
Variance: Test error, due to difference in selection of "test-set".

► Bath's eye diagram:

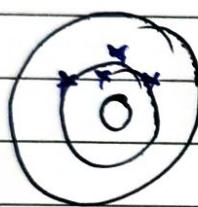
## → Bull's eye diagram:



low bias and  
low variance



high bias and  
low variance



low bias and  
high variance



high bias and  
high variance

★ Near to the inner-circle (low bias)

★ Cluster more compact (less variance)

How to balanced fit model:

i) K-fold Cross Validation.

ii) Regularization.

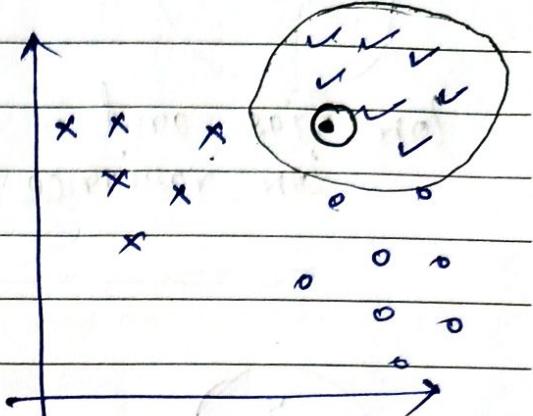
iii) Dimensionality Reduction.

iv) Bagging and boosting.

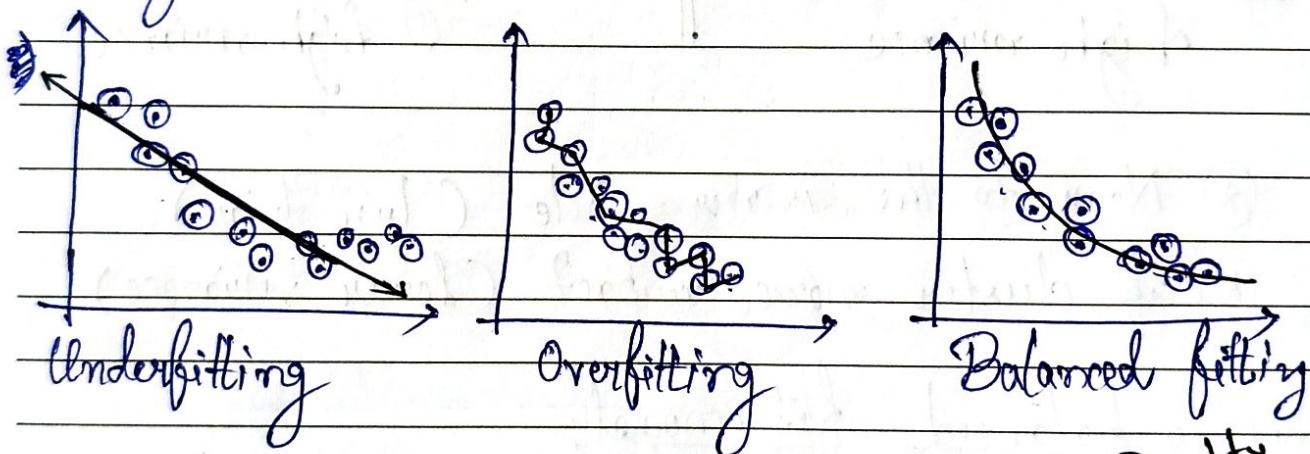
## K-Nearst-Neighbour (KNN)

supervised classification technique

Philosophy: Nearest K-vectors of a vector, indicates its class (category).  
from stlearn.neighbors import  
KNeighborsClassifier



## Regularization



### i) Ridge Regression ( $\ell_2$ )

→ Cost function:  $J(\theta) = \text{MSE}(\theta) + \alpha \sum_{i=1}^n \theta_i^2$   
↳ need to be minimized  $\Rightarrow$  smaller values of  $\theta$ ; i.e. are allowed.

» from stlearn.linear\_model import Ridge

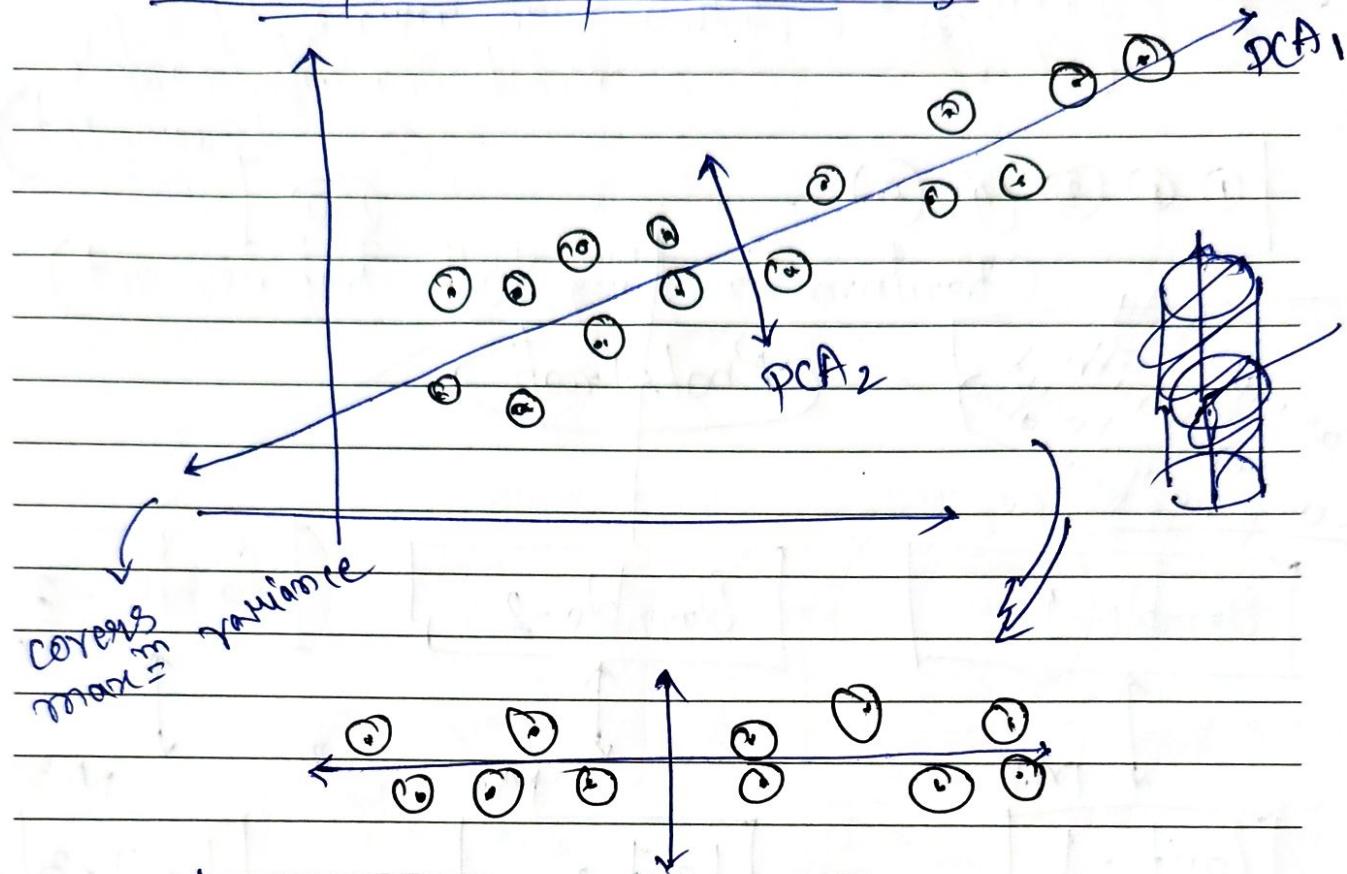
### ii) Lasso Regression ( $\ell_1$ )

Cheat Absolute Shrinkage and Selection

→ Cost function:  $J(\theta) = \text{MSE}(\theta) + \alpha |\theta| = \alpha \sum_{i=1}^n |\theta_i|$   
» from stlearn.linear\_model import Lasso

(all about variation max<sup>n</sup>)

## → Principal Component Analysis:



[\*) Scaling is very important

>>> from sklearn.decomposition import PCA

>> pca = PCA (n\_components = 4/0.95) # Transformer Class

>> X\_final = pca.fit\_transform(df)

>> pca = pca.fit(df[[1]])

>> pca.components\_

>> pca.explained\_variance\_ratio\_

~~Majority Vote~~

~~→ Bagging : (Jacking high variance by considering diff. sample train-set)~~

[① ② ③ ④ ⑤]

[100]

(Random Sampling With Replacement)

Datasets Only  
In contrast with  
Random Forest  
(feature also)

Bootstrap

Sample - 1

Sample - 2

Sample - 3

Logit - 1

Logit 2

Logit - 3

↓ Prediction

↓ Prediction

↓ Prediction

1

0

1

Majority vote  
(Aggregation)

1

∴ Bagging [Bootstrap - Aggregation]

Difference  
Bagging → Underlying model can be anything (SVM, KNN, etc.)

Bagged Tree → each model is a tree.

From sklearn.ensemble import BaggingClassifier  
/ Regressor

# Gradient Boosting

(Sequential Stagewise Addition)

$$(*) \quad [M_1] + [M_2] + \dots + [M_n] = \sum H_B$$

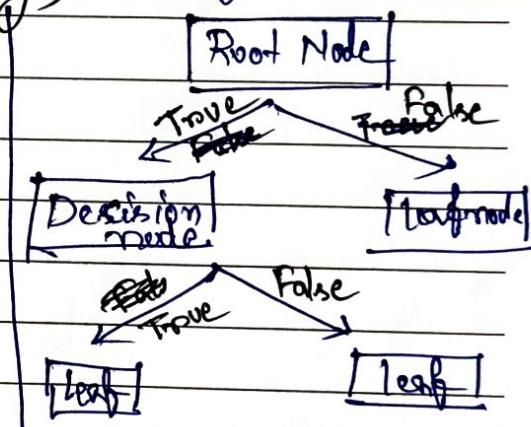
Weak-learner (High Bias)

↓ to get

Generalized Model  
(Balanced fit)

→ Basically training models (typically Decision Trees) on data training (residuals as target)  
(till last stage)

→ Prediction: learning rate  
 $M_1 + \alpha(M_2)$   
 $+ \alpha(M_3) + \dots$



→ Pseudo Residual:

$$\text{Actual} - \text{Predicted}$$

~~Actual~~      ~~Predicted~~

(\*) Base Model

a) Regression → Simple mean

b) Classification →  $\log(\text{Odds}) \rightarrow \text{Probability}$

(\*) Ada Boost

■ Decision stump  
(2 leaves)

■ "Weighted" Addition  
of models

Gradient Boost

■ 8-32 Leaf-nodes

■ Same learning-rate  
throughout

Sigmoid(log(Odds))

(\*) This leads to the extreme application XGBoost.