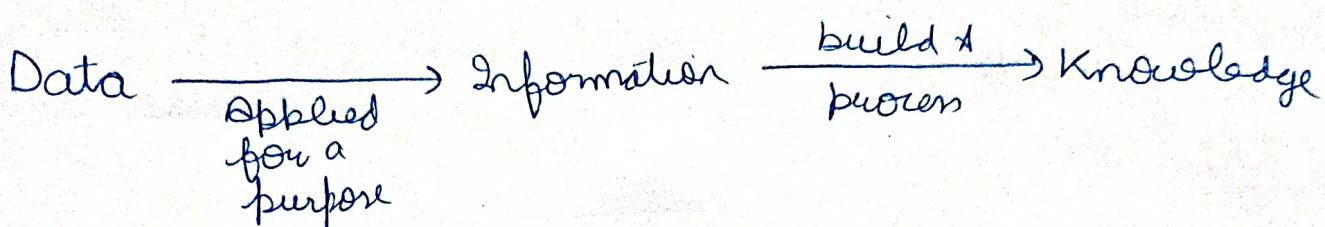


Machine learning

- Ability of algorithms to learn from data
- logic is not provided by programmer, but the machine understand it by itself
- We provide a ~~lot~~ input and output and train it with a large amount of inputs, then deploy it



input data → x-train
output data → y-train] Training data

input → x-test
output → y-test] testing data

Nominal Data

→ No ordering

→ Data having label but not order

e.g. Gender

→ Male

'm → Female

Ordinal Data

→ Data having natural order

e.g. How do you feel?

1 - Very happy

2 - happy

3 sad

4 Very Sad

Interval data
data having intervals blue each value

eg Temp?

30°
40°
60°

Binary data

Only 2 labels
0, 1

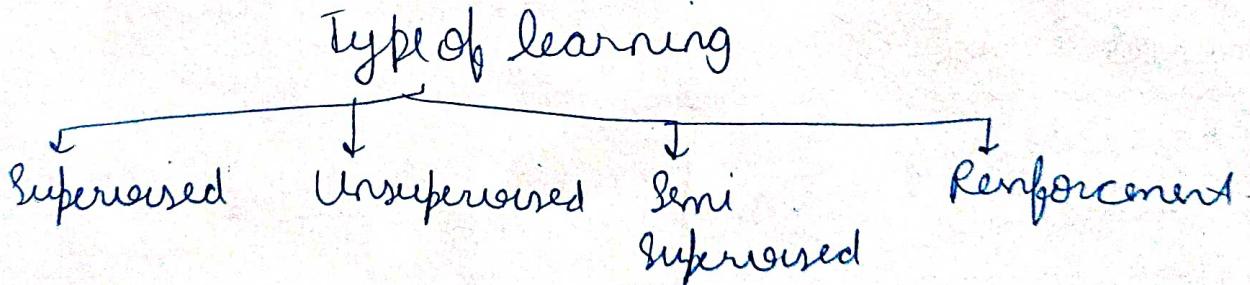
Yes/No

Time data

Date, year, month, hour

LEARNING for a MACHINE

→ A machine is said to be learning from past experience with respect to some class of task.



Supervised learning

When model is getting trained on a labelled dataset
labelled dataset is one having both input & output

Classification

↳ Discrete output

e.g. 0 or 1, Yes/No

Regression

→ Continuous output parameter

e.g. (0-5)

any value

0, 1, 0.2, 0.3, 1.3, 2.4

Unsupervised learning

→ No output / we don't give target to our model while training

Clustering

→ group the data sharing the similar properties

Association

→ try to look for best connection
blue parameters of large data

Semi-supervised learning

→ Some data has output labels, some have not
both labelled & unlabelled data

Reinforcement learning

In this model keeps on increasing its performance, using a reward feedback to learn the behaviour pattern

e.g. self driving car



* learn from mistakes

Classification

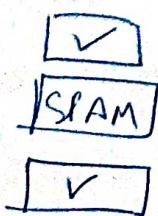
- Email spam detection (Yes/No)
- Image recognition (✓/✗)
- Stock price prediction

Regression

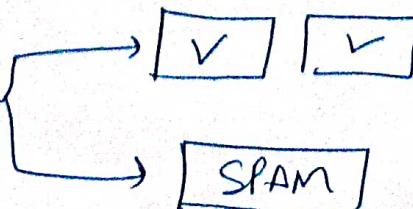
- House price prediction
- Health insurance assessment

e.g. email spam detection

Emails



} classifier



AI (Artificial Intelligence)

Intelligence

- ↳ ability of reasoning
- ↳ Reward based learning
- ↳ Generalized learning.

Subsets of A.I.

- Machine learning
 - Natural language processing
 - Computer vision (dealing with images/video)
 - Robotics (Hardware + Software)
- Deep learning
- Data science

* Neural network

- ↳ networks of neurons
- ↑
 single unit of an algorithm

Data scientist jobs

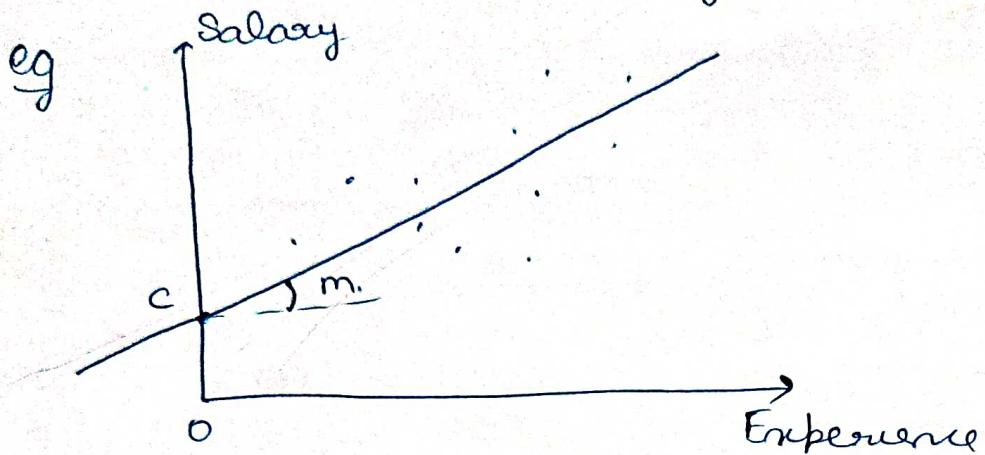
→ understanding/problem

→ (0)

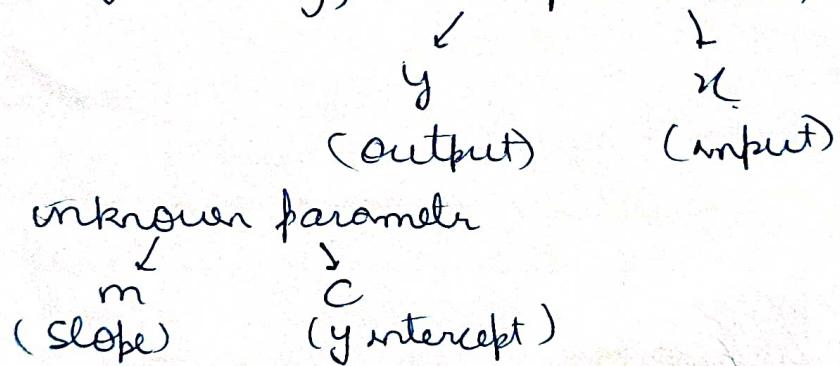
Ist ML ALGORITHM

* LINEAR REGRESSION

→ Hypothesis, $y = mx + c$ straight line



At time of training, known parameters



Target → To have least fit line, will predict output with least error

↳ Steps

- ① Forward propagation (pick random value of m, c)
- ② Calculate error (or cost function)

$$\boxed{\text{Cost fn} = \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{2n}}$$

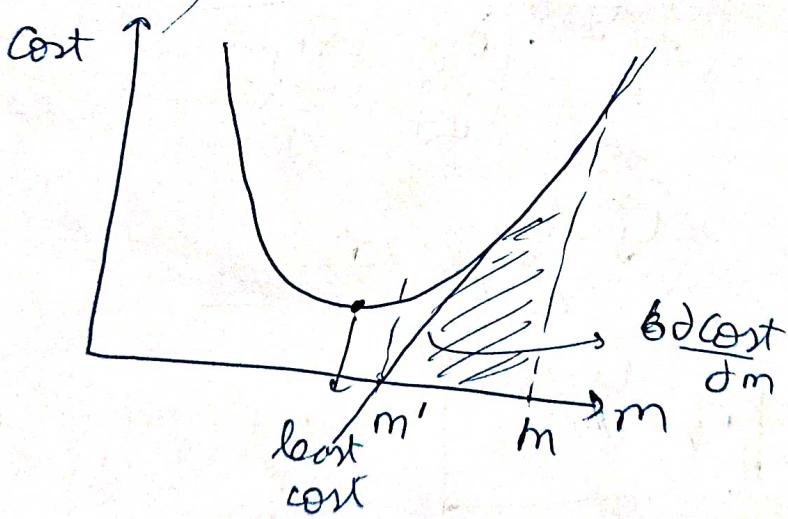
\hat{y}_i → predicted output
 y_i → actual output
 n → no. of instances

3- Gradient Descent algorithm

↳ used to get optimal value of m and c
so that we can get least error or cost function

$$m' = m - \alpha \frac{\partial \text{Cost}}{\partial m} \rightarrow \text{To get minimum cost value}$$
$$c' = c - \alpha \frac{\partial \text{Cost}}{\partial c} \rightarrow \text{W.r.t. } m \text{ and } c$$

(Learning rate) α is used to make gradient fast or slow accordingly.



Forward propagation

↓
Cost func'

↓
Gradient descent

Backward propagation

Maths behind linear regression

$$m' = m - \alpha \frac{\partial \text{Cost}}{\partial m}$$

$$\Rightarrow \frac{\partial \text{Cost}}{\partial m} = \frac{\partial \text{Cost}}{\partial f} \cdot \frac{\partial f}{\partial m} \quad \text{①} \quad [f: \hat{y} = mx + c]$$

$$\Rightarrow \frac{\partial f}{\partial b} = \frac{\partial \text{Cost}}{\partial f} \Rightarrow \frac{\partial \frac{(\hat{y} - y)^2}{2n}}{\partial \hat{y}}. \quad (\text{ coz } \hat{y} = y)$$

$$\Rightarrow \frac{\partial}{\partial n} (\hat{y} - y) \cdot \frac{\partial (\hat{y} - y)}{\partial \hat{y}}$$

$$\Rightarrow \left(\frac{\hat{y} - y}{n} \right) \cdot \left(\frac{\partial \hat{y}}{\partial \hat{y}} - \frac{\partial y}{\partial \hat{y}} \right)$$

$$\Rightarrow \left(\frac{\hat{y} - y}{n} \right) (1 - 0)$$

$$\boxed{\frac{\partial f}{\partial b} : \frac{\hat{y} - y}{n}} = \frac{\partial \text{Cost}}{\partial f} - \checkmark$$

$$\text{Now } \frac{\partial f}{\partial m} = \frac{\partial (mx + c)}{\partial m}$$

$$\Rightarrow \frac{\partial mx}{\partial m} + \frac{\partial c}{\partial m} \Rightarrow x + 0 \Rightarrow x$$

$$\frac{\partial f}{\partial c} \Rightarrow \frac{\partial mx + c}{\partial c} \Rightarrow 1$$

$$\delta m = \frac{\partial \text{Cost}}{\partial m} = \frac{\partial \text{Cost}}{\partial f} \cdot \frac{\partial f}{\partial m}$$

Use equation

$$\left. \begin{aligned} \delta m &= \left(\frac{\hat{y} - y}{n} \right) x \\ \delta c &= \left(\frac{\hat{y} - y}{n} \right) \cdot 1 \end{aligned} \right\}$$

Gradient descent

$$\left. \begin{aligned} \delta f &= f(x) - y \\ \delta m &= \delta f * x \\ \delta c &= \delta f \end{aligned} \right\}$$

$$m' = m - \alpha \delta m$$

$$c' = c - \alpha \delta c$$

$$\left. \begin{aligned} m' &= m - \alpha \left(\frac{y' - y}{n} \right) x \\ c' &= c - \alpha \left(\frac{y' - y}{n} \right) \end{aligned} \right\}$$

Updating parameters

Multiple Linear Regression

→ Same thing as linear but with multiple inputs

$$\rightarrow y = \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3 + b$$

b - bias / constant value

$$\partial m = \frac{\partial \text{Cost}}{\partial b} \times \frac{\partial b}{\partial m} \Rightarrow \frac{\hat{y} - y}{n} \times x$$

$$\partial \omega_1 = \frac{\hat{y} - y}{n} (x_1)$$

$$\partial \omega_2 = \frac{\hat{y} - y}{n} (x_2)$$

$$\vdots \quad \vdots \quad \vdots$$

* Linear Regression Assumptions

→ linear relation b/w input & output
 $y = mx + c$

.5 → Homoscedasticity

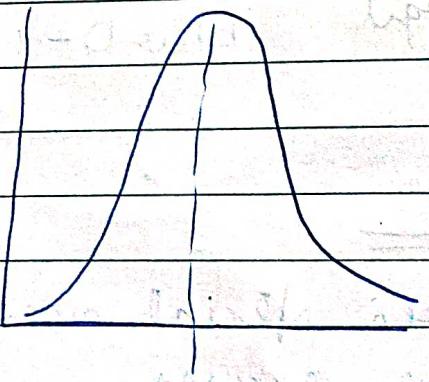
As, prediction = $\hat{y} = mx + c$ ①

residuals, error = $(\hat{y} - y)$ ②

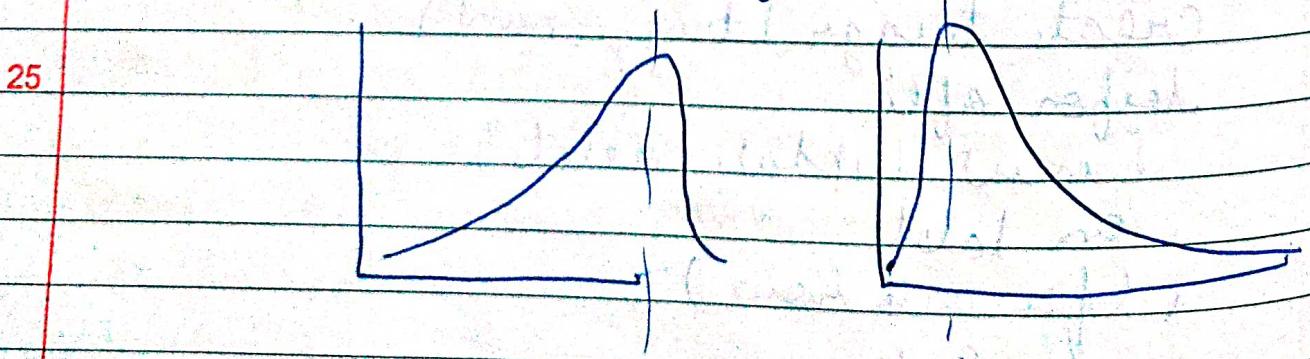
10 so graph b/w ①, ② will be random
no relation b/w ①, ②

→ Normality in residuals

15 should have normal (gaussian) distribution (symmetric)



skewed distribution (high error)



* OLS method

→ Ordinary least square ←

5 R-squared value

$$R^2 = 1 - \frac{RSS}{TSS}$$

10 $R^2 \rightarrow$ coeff. of determination
RSS → sum of squares of residuals
TSS → Total sum of square

15

20

25

30

Polynomial linear Regression

→ If data is not linear

5

eg →

10



As in linear.

15

$$y = mx + c$$

So eqn can be

20

$$\hat{y} = m_0 x^0 + m_1 x_1 + m_2 x_2 + m_3 x_3$$

We can get curves using degrees, x_1^2, x_1^3, \dots

25

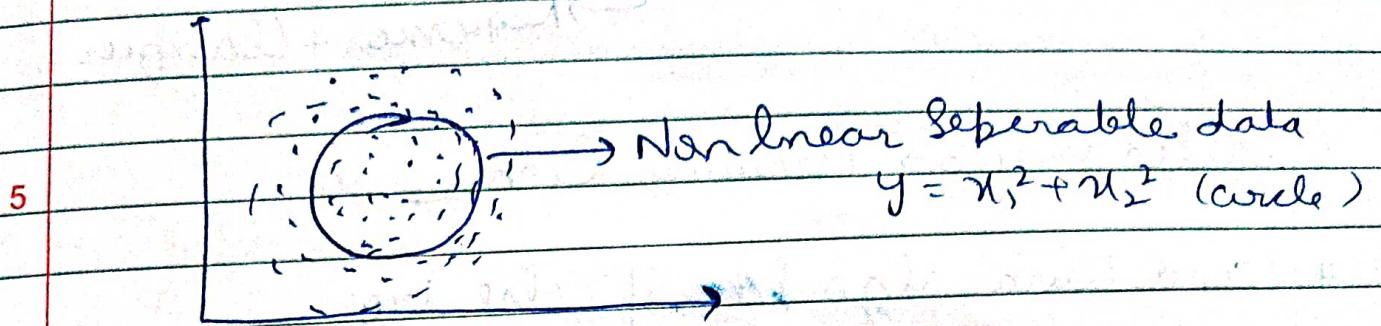
* LR $\rightarrow \hat{y} = mx + c$

* MLR $\rightarrow \hat{y} = m_1 x_1 + m_2 x_2 + m_3 x_3 + c$

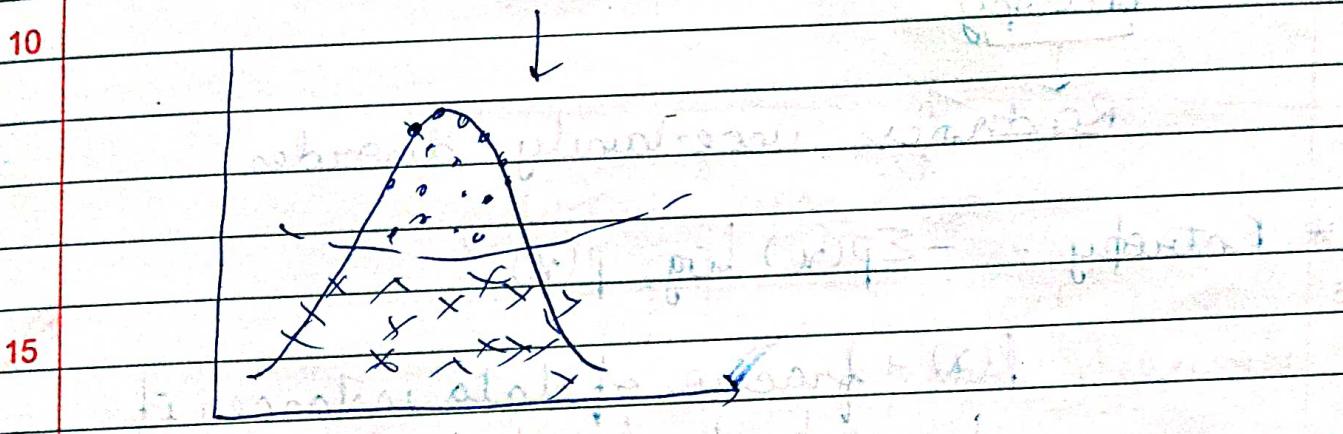
* PLR $\rightarrow \hat{y} = m_1 x_1 + m_2 x_1^2 + m_3 x_1^3 + c$

30

Support vector machine algo



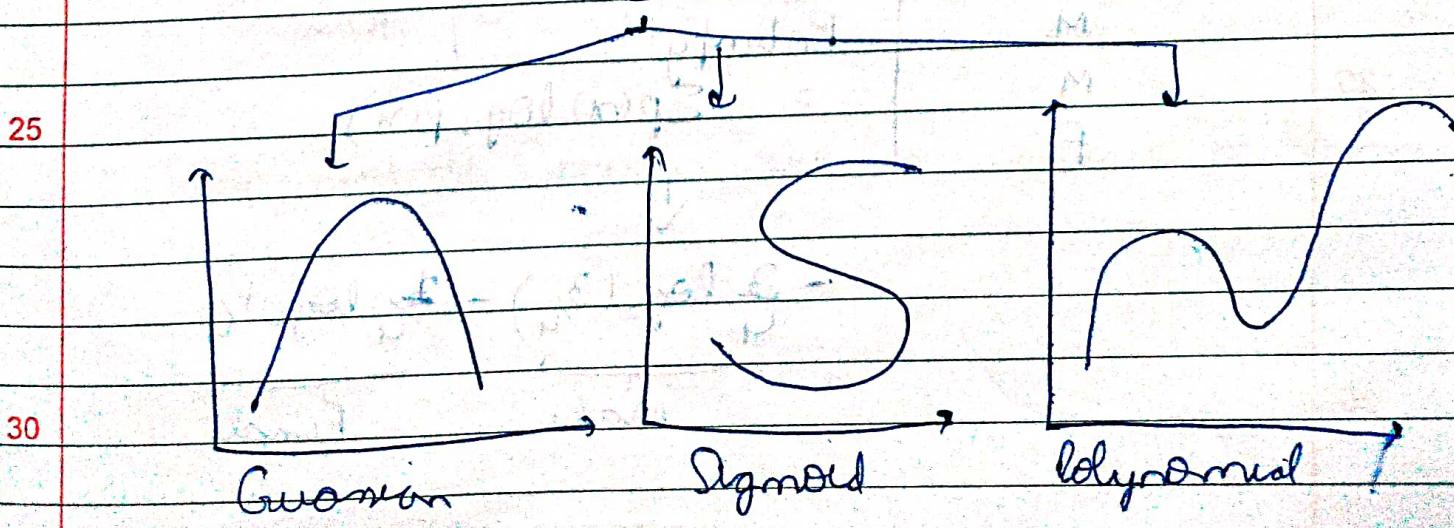
→ We map our data to kernel



* SVM takes data and add dimensionality to make data linearly separable using kernels

20 5 * Then you use hyperplane

Kernels



Decision Tree Algo

↳ Regression + Classification

5 called as ID3 → Iterative Dichotomiser

* conditional algo like if-else tree

↳ Info gain

↳ Entropy

↳ Randomness, uncertainty, disorder

* Entropy :- $-\sum p(x) \log_2 p(x)$

$p(x) \rightarrow$ fraction of data instances of
↓ particular type

Probability of a particular data x

eg Gender]
M
M
M
F

- Entropy
 $= -\sum p(x) \log_2 p(x)$

↓
- $\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \frac{1}{4}$
Male Female

* Information Gain

→ how much info a feature / input parameter gives us about output / target parameter

5

$$\text{Info Gain}_{(P_i)} = \text{Entropy}_{(\text{Target})} - \text{Entropy}_{(P_i, \text{Target})}$$

P_i - Parameter

10

Random forest Algo

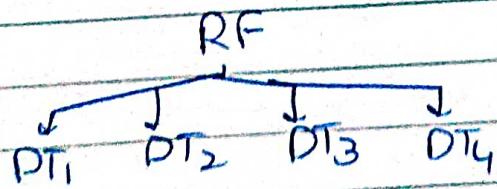
15

works as both
regressor & classifier.

→ Ensemble model → Bagging model

20

↓
we can have n no. of decision trees



25

→ Enhancement of Decision Tree to increase performance

30

CLASSIFICATION

ALGORITHMS

DATE : / /
PAGE :

→ logistic Regression algo

↳ uses a sigmoid function

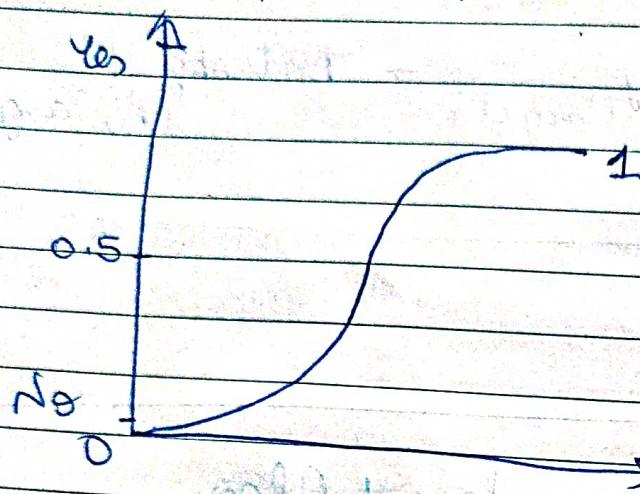
5

10

15

20

25



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

$$z = \hat{y} = mx + c \quad (\text{y predicted})$$

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

* Case-1

\hat{y} is very large

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

$$\hat{y} \uparrow \rightarrow -\hat{y} \downarrow \rightarrow e^{-\hat{y}} \downarrow$$

$$\underline{g(z) \approx 1}$$

Case ②

$$\text{if } y \downarrow \rightarrow -\hat{y}^T \rightarrow e^{-\hat{y}^T}$$

5

$$\textcircled{a} \quad g(2) \approx 0$$

Cost function

10

$$J = \text{Cost} = \frac{1}{n} \sum_{i=1}^n [L(y_i - \hat{y}_i)]$$

$\rightarrow L = \text{loss}$

15

$$\text{loss} = [y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

* If $y = 1$

$$L = -y \log \hat{y} \rightarrow -\log \hat{y}$$

20

If we want loss to minimum,

$$\log \hat{y} \rightarrow$$

$$\hat{y}^T \rightarrow 1 \text{ (max prob)}$$

25

* If $y = 0$

$$L = -(1-y) \log (1-\hat{y})$$

$$= -\log (1-\hat{y})$$

30

loss min $\rightarrow \log (1-\hat{y})^T$

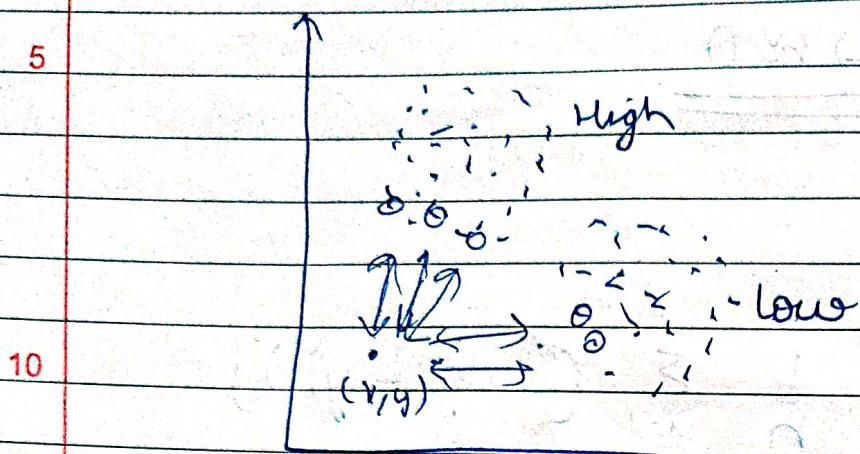
$$(1-\hat{y})^T$$

$$\hat{y} = 0$$

(max prob)

K-NN algorithm

K-nearest neighbours



→ uses euclidean distance

$$\sqrt{(x_{\text{test}} - x_{\text{train}})^2 + (y_{\text{test}} - y_{\text{train}})^2}$$

It will find k nearest neighbour and determine whether pt. lies where in class

~~K-Nearest Neighbour~~

Naïve Bayes Classifier

2 toss
 Coin → (Event)
 $T \quad H \rightarrow \text{Sample space}$
 $\mapsto \{HH, HT, TH, TT\} \quad \{T, H\}$

$$P(\text{at least one head}) = 3/4$$

5
 $P(2 \text{ head}) = \frac{1}{4}$

$P(1 \text{ head}, 1 \text{ tail}) = \frac{1}{4}$

* Bayes Theorem

10 $P(A/B) = \frac{P(B/A) * P(A)}{P(B)}$

15 Event
 ↓
 $E_1 \quad E_2$

$P(\text{both } E_1, E_2 \text{ happening}) = P(E_1) * P(E_2)$

20 $P(\text{either } E_1 \text{ or } E_2 \text{ happened}) = P(E_1) + P(E_2)$

25

30

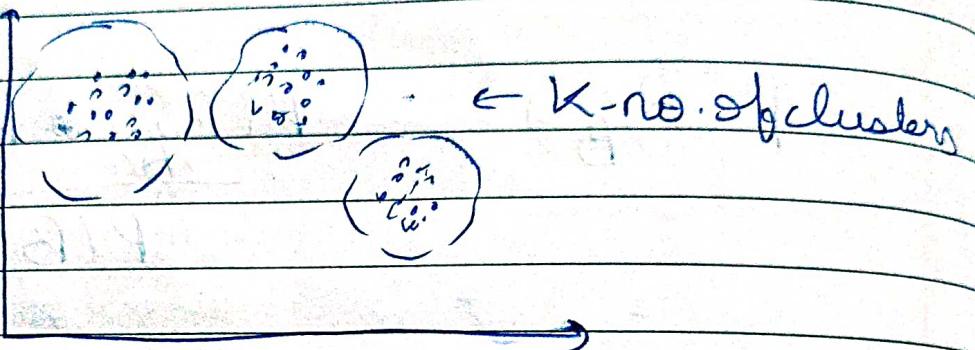
* Unsupervised learning

→ Clustering Algorithms

5 ① K-means Clustering

↳ also called hard clustering

10



15

Steps

→ Random Initialization

→ Assign pt. to nearest neighbour

→ Updating mean coordinate

* 20

→ Repeating process until cluster is found.

* → This algo. will categorize the items into K groups or clusters of similarity
To calculate similarity, use Euclidean distance.

25

30

* Feature engineering

→ data cleaning

→ Scaling of data

5

↳ from sklearn.preprocessing

import MinMaxScaler

↳ $x = pd.DataFrame(MinMaxScaler().fit_transform(x))$

→ Correlation

↳ heatmap or

10

↳ Extra Tree classifier

(from sklearn.ensemble import ExtraTreeClass)

(model = ETC())

(model.fit(x, y))

(features_inh = pd.Series(model.feature_importances_))

(feature_inh.nlargest(1))

15

↳ Select K best method

20 * PCA (Principal component Analysis)

in

→ To reduce data parameters to efficient amount

* suppose we have a data of 100 dimensions,

25 we can not plot it as it is, we use
PCA

30