

# Machine Learning

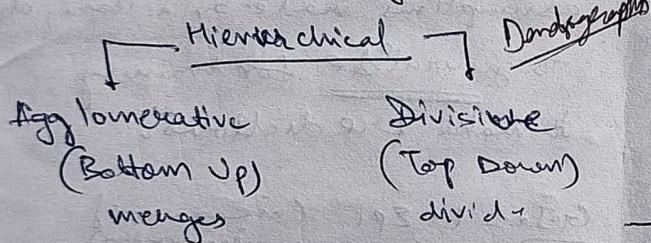
Learning → Change in behaviour

USML → coins sorting

training model w/o any labels, identify patterns on its own.

Machine learning is a subset of AI enabling computers to learn & make decisions based on data, w/o being explicitly programming.

Clustering → learning through observation based on common features.



## Types of Clustering

DBSCAN → cluster points that are closely packed together (point density)

Grid Based: grid structure env. performs clustering on grid. Large spatial data, weather forecasting

Spectral: (multidimensional) reduce dimensions, social network analysis, eigenvalues of similarity matrix

$$P(C_i | X) = \frac{P(X | C_i) \cdot P(C_i)}{P(X)}$$

$$P(X | C_i) = \prod_{j=1}^n P(x_j | C_i)$$

## SMC

input / output, training, error

Training model with <sup>Labels</sup> data

for every input output is provided

→ Image Recognition → Medical Data  
Speech Recognition, Diagnosis Process  
→ Automation complex data.

Reinforcement Learning → agent learns from environment, makes decisions by performing actions and receiving feedback in form of rewards / penalties

## Environment, Action, Reward

### Centroid Based Clustering

inter cluster similarity

monopole, accuracy, data req.

Regression: statistical technique to relate dependent var. ( $y$ ) to one or more independent vars. ( $x$ )

$$y = \bar{w}\bar{x} + b \quad \text{Linear}$$

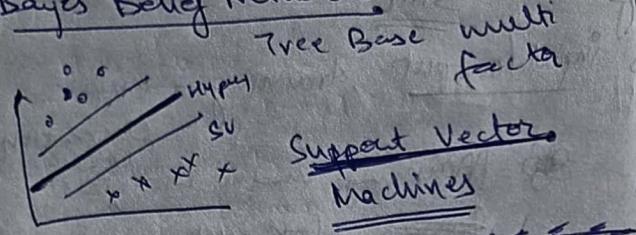
Logistic Regression → probability input flows to certain category

Concept Learning → inference from label training data, patterns

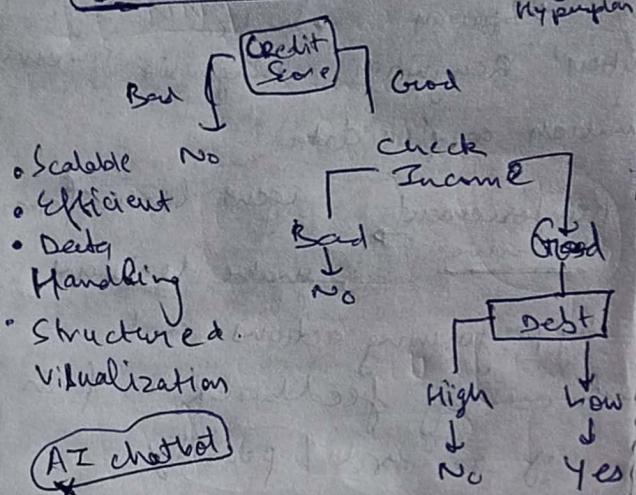
Bayes Classifier → highest posterior probability  
→ spam / 'offer' / 'promotion'. theoretical problem

$$P(Y | X_1, X_2, \dots, X_n) = \frac{P(Y) \cdot P(X_1 | Y) \cdot P(X_2 | Y) \cdots P(X_n | Y)}{P(X_1) \cdot P(X_2) \cdots P(X_n)}$$

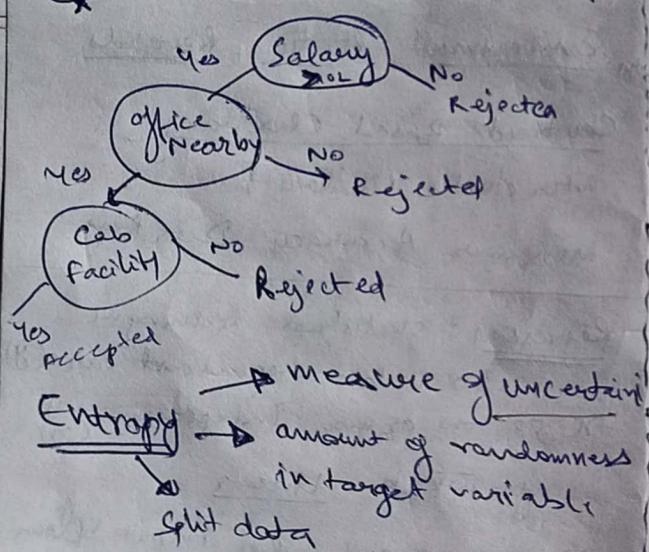
## Bayes Belief Networks



## Decision Trees



## AI chatbot



$$\begin{array}{lll} E=0 & E=1 & 0 < E < 1 \\ \text{min} & \text{max} & \text{mix} \\ \text{uncerta} & \text{certain} & \end{array}$$

I<sub>G</sub> → measure change in entropy due to factors.

$$\begin{aligned} \text{Subset } H(S) &= P(A)H(A) + P(B)H(B) \\ A \quad | \quad & \quad B \quad | \quad C \\ H(A) & \quad H(B) \quad H(C) \end{aligned}$$

$$(I_G) = H - H(S)$$

Highest  $I_G$  is first Selected.

## Bayes

$$P(y=c/x_i) \propto \prod_{i=1}^n P(x_i/y=c) \cdot P(y=c)$$

ID3 → top down, greedy, High IG first, more overfitting

C4.5 → extension of ID3

## Gain Ratio

$$(1 - \sum P_i^2)$$

Prunes the tree, handles missing values & removes overfitting, categorical & continuous.

CART → Gini Impurity, Binary Splits, pruning Random Forest, Robust

## Inductive Bias

assumptions made by a learning alg. to generalize for training to make predictions.

$$G_I = 1 - \sum P_i \quad (\text{measure of impurity})$$

$$G_I = 0, \text{ pure}$$

## KNN

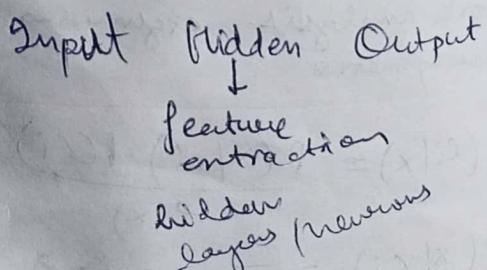
read self non parametric SMLE Euclidean Distance.

Lazy learning.

$$K \quad \left| \begin{array}{l} \text{imbalanced data} \end{array} \right.$$

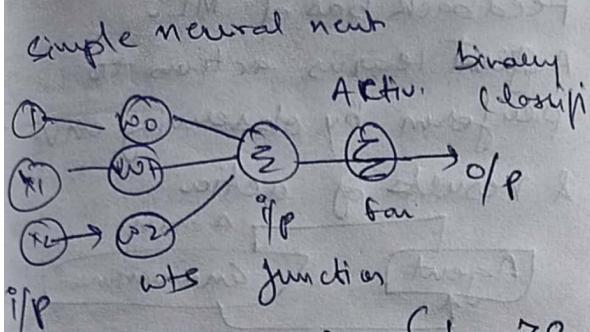
ANN → computing system to simulate human brain

data processing system, neuron, self learning capability



## Recurrent Network / Perceptron

Artificial neuron.



$$\sum(w_i x_i + b) = \begin{cases} 1, > 0 \\ 0, \leq 0 \end{cases}$$

And gate

$$x_1 + x_2 - 1 = y$$

$$\Delta w_i = \alpha(t - o) \times u_i$$

$$\Delta b_i = \alpha(t - o) \quad b = b + \Delta b$$

Linearly Separable Data.  
perfectly separated by ~~straight~~ line  
adjust wt. & bias to minimize error

$$\text{New value} = [x - f'(x)] x$$

$$E = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$

$$\Delta w = -\alpha \nabla E(w)$$

$$\nabla E(w) = \left( \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right)$$

$\frac{\partial E}{\partial w_i}$   $\alpha$ -slow imp

$\alpha$  + (slow)

$\alpha$  + (very fast jump)

Stochastic (Random) | Mini-Batch  
Smaller Batches  
or

Hierarchical clustering is connectivity based clustering model that groups data points together that are close to each other based on measure of similarity (distance)

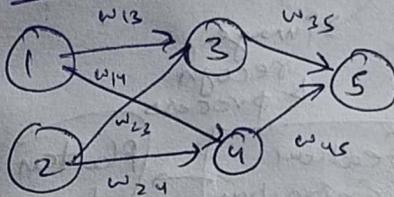
## Backpropagation

Repeatedly adjust weights of connections to minimize difference b/w actual op & desired op

Question

[1] Fwd Pass

B/W Pass



$$\Delta w_{35} = n S_5 H_3$$

$$S_5 = (y - o_5) o_5 \cdot (1 - o_5)$$

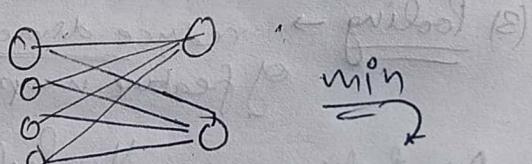
$$S_3 = S_5 \cdot w_3 \cdot H_3 \cdot (1 - H_3)$$

$w \rightarrow S$

Competitive Learning Algo

Self Organizing Maps  
to lower dimensions

$$D = \sqrt{\sum (w_i - w'_i)^2}$$
 euclidean



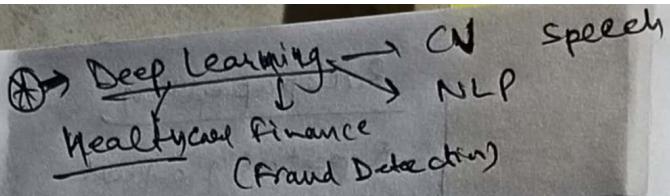
$$w_{\text{new}} = w_{\text{old}} + \alpha (x_i - w_i)$$

update

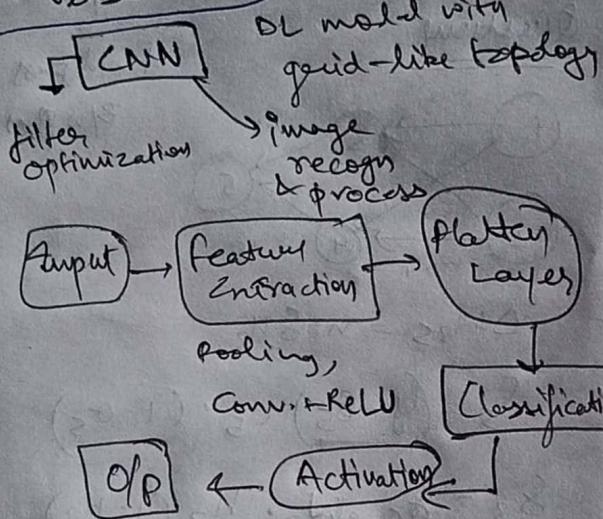
Stop when diff. changes = 0

Neural Network  $\rightarrow$  Deep Learning  
subset of ML that involves NN with many layers.

large amt. of Data



- Feature Extraction
- High performance
- Versatile



(1) Feature Extraction (convolution op.)

→ filters / kernel → matrix  
 (dot product b/w filter.)

(2) Padding → adding zeros to preserve dims

(3) Pooling → reduce dimensions of feature maps

Avg Pooling

avg value from each patch

Reduce Noise, Distortion

ReLU max(0, x)

Max Pooling

( $2 \times 2$ ) window

Reduce Noise, Distortion

ReLU max(0, x)

Policy ( $\pi$ )

↳ strategy / rule that defines action to be taken in each case.

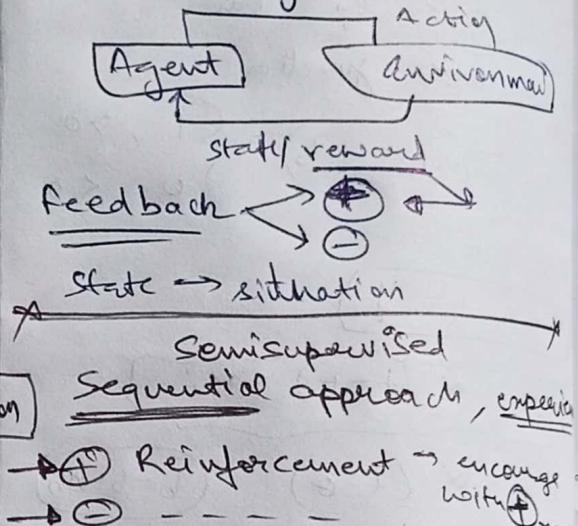
Adam → adjusts speed & learning rate

Gradient Descent

## Reinforcement Learning

feedback based ML.

An agent learns actions to perform by observing env. & results of action



- + Reinforcement → encourage with  $R_{new}$
- - - - - -

Maze Game, Grid game.

- Actions are interdependent
- Delayed feedback
- Time consuming

traffic light, industrial automation, driverless cars, recommendation system.

## Markov Decision Process

formal model for decision making where outcomes are partly random partly under control of agent

### Memoryless property

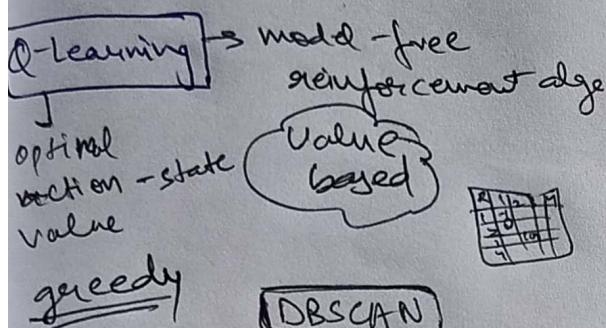
future state depends on current state  
 Action

Not on sequence of events

ROC → class distribution is balanced  
 $fp \approx fn$

PR curve → imbalanced, positive class predict dataset

State, Actions, Transition (P)  
Reward, Discount factor



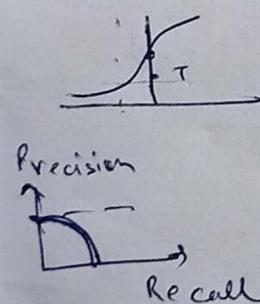
Density Based Spatial Clustering of Applications with Noise

$$\text{New Val} = \text{Old Val} + \alpha (\text{Reward} \times \frac{\text{Q}}{\text{Ngb}})$$

update table

Reward function  $Q(s, a)$   
assigns numerical value to an agent based on its current state & action. (process matter)

Value function  $V(x)$   
(Bellman eqn)  
a function that estimates the expected long term reward an agent will receive for a given state / action



$T \uparrow$ , Recall  $\uparrow$   
Precision  $\uparrow$   
 $FP \uparrow$ ,  $FN \uparrow$

Area under curve

Covariance matrix for each feature

$$S = \begin{bmatrix} X_1 X_1 & X_1 X_2 \\ X_2 X_1 & X_2 X_2 \end{bmatrix}$$

$$\frac{1}{N-1} \sum (X_{1k} - \bar{X}_1)(X_{2k} - \bar{X}_2)$$

$$(S - \lambda I) = 0$$

② Solve determinant

$$\text{get } \lambda \approx \lambda_1, \lambda_2$$

③ Eigen values

$$(S - \lambda I) \rightarrow \lambda = 0 \quad U = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$

$$\frac{u_1}{\alpha} = \frac{u_2}{\beta} = t$$

$$U = f(\lambda) = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$$

$$\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \quad \text{Eig}$$

④ unit eigen vector.

$$e_i = \begin{bmatrix} \frac{1}{\sqrt{1+|U|}} \\ \dots \end{bmatrix} \quad \text{unit vector}$$

$$e_1 = f(\lambda_1) \quad e_2 = f(\lambda_2)$$

⑤ First Principal Component

$$e_1^T = \begin{bmatrix} x_{1k} - \bar{x}_1 \\ x_{2k} - \bar{x}_2 \end{bmatrix}$$

for each example

$$p(y|x) = [p_{y|x}]^y (1-p_{y|x})^{1-y}$$

logistic maximum likelihood estimation

$$\frac{\partial}{\partial \theta} (\log LL) = \frac{1}{m} \sum_i (y_i - h_{\theta}(x)) x_i$$

$$y_i \log(p_{y|x}) + (1-y_i) \log(1-h_{\theta}(x))$$

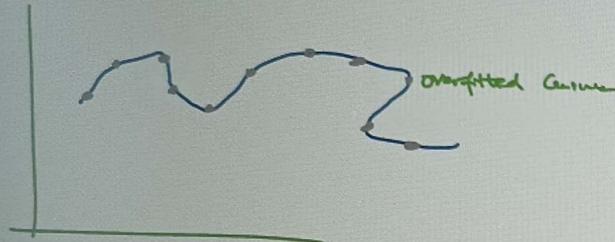
= Binary Cross Entropy

CNN

Convolutional Layer - Feature extraction  
Pooling layer → reduce size  
Fully Connected Layer → classification

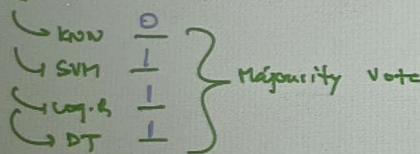
### Overfitting:

Model performs very well on training data but doesn't perform good on test data.



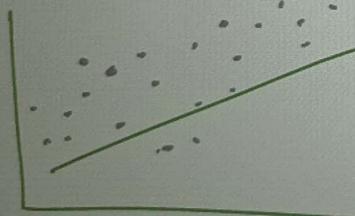
### Solutions:

- k fold cross validation
- Sufficient data → Ensembles
- Ensembling techniques.

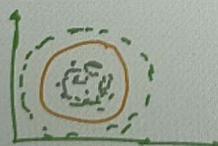


### Underfitting:

Model is not going to learn patterns from training data.



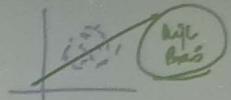
Underfit model will give poor performance on train as well as test data.



### Solutions:

- No. of features increase
- Model Complexity
- Reduce noise
- Increase the duration of training.

### BIAS AND VARIANCE:



Bias: wrong assumptions about data - like assuming data is linear in reality is follows a complex for<sup>n</sup>.

It is the inability of the model bcz of that there is diff in predicted value & actual value

Low Bias: lower assumptions make a model which closely matches the training (Simple model) dataset.

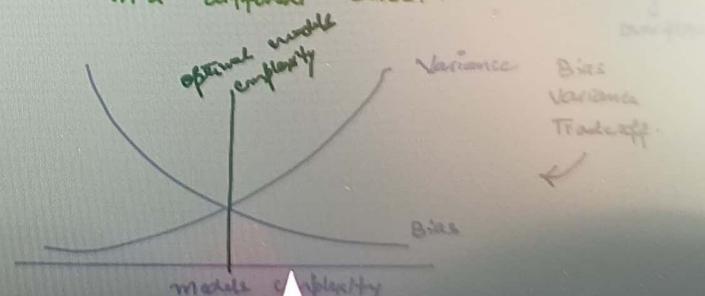
High Bias: more assumptions Model will not match training dataset closely. (Complex model)

### Variance:

→ measure of spread in data from its mean position  
→ sensitive to a subset which follows same distribution as your training dataset

Low Variance: model is very less sensitive to changes.

High Variance: model is very sensitive to changes and can result in significant changes if trained on a different subset.



# ROC Curve (Receiver Operator Characteristics Curve)

➤ An ROC curve (receiver operating characteristic curve) is a graph showing the performance of a classification model at all classification thresholds.

➤ This curve plots two parameters: TPR & FPR

➤ True Positive Rate (TPR)=  $\frac{\# \text{ true positives}}{\# \text{ actual positives}} = \frac{\# \text{ true positives}}{\# \text{ true positives} + \# \text{ false negatives}}$

➤ False Positive Rate (FPR)=  $\frac{\# \text{ false positives}}{\# \text{ actual negatives}} = \frac{\# \text{ false positives}}{\# \text{ true negatives} + \# \text{ false positives}}$

## Common Approaches for Reward Function Design

### Sparse vs. Dense Rewards

- **Sparse rewards** provide feedback only when a significant event occurs (e.g., the agent only receives a reward when it reaches the goal). Sparse rewards can make learning challenging but encourage exploration.
- **Dense rewards** provide feedback at every step (e.g., the agent receives a small reward for each correct action that moves it closer to the goal). Dense rewards facilitate faster learning but may lead to premature convergence on suboptimal strategies.

### Shaping Rewards

Semi-supervised learning is a machine learning paradigm where the dataset contains both labeled and unlabeled data. This approach is particularly useful when labeled data is scarce or expensive to obtain, but there is an abundance of unlabeled data available.

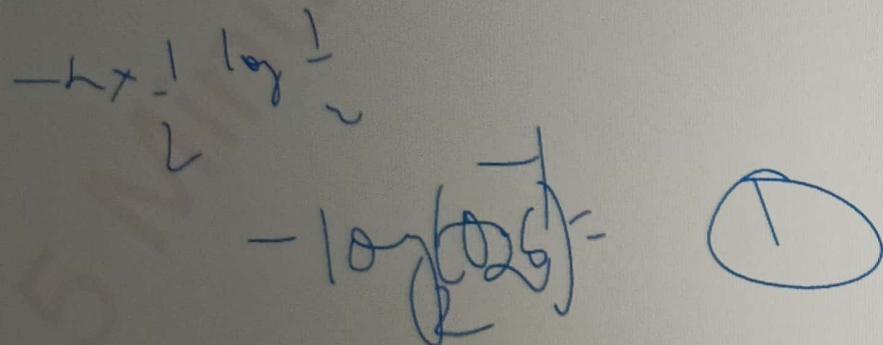
## Example of Semi-supervised Learning

Imagine you're working on a sentiment analysis task where you want to determine the sentiment (positive, negative, or neutral) of customer reviews for a product.

However, you only have a small dataset of labeled reviews, but there's a vast amount of unlabeled reviews available online.

$\lambda_1 \lambda_2$

So,  $\lambda_1 \rightarrow$  is our principal component.



③ value f<sup>n</sup>: value of state is total amount of reward an agent can expect to gain over the future starting from that particular state.

① Policy: defines agents behaviour for a given state / time. ✓  
state  $\xrightarrow{\text{mapping}} \text{Action}$

② Reward f<sup>n</sup>: It provides a numerical score based on state of environment

① Initialize Q-table

↳ It's a simple lookup table where we calculate max. expected future rewards for action at each state.

Actions:		↑	→	↓	←
S	0	0	0	0	
B	0	0	0	0	

Rows:- no. of states

# DBSCAN Clustering Algorithm Solved Example



P1: (3, 7)

P2: (4, 6)

P3: (5, 5)

P4: (6, 4)

P5: (7, 3)

P6: (6, 2)

P7: (7, 2)

P8: (8, 4)

P9: (3, 3)

P10: (2, 6)

P11: (3, 5)

P12: (2, 4)

minPts = 4 and epsilon ( $\epsilon$ ) = 1.9

	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	P11	P12
P1	0											
P2	1.41	0										
P3	2.83	1.41	0									
P4	4.24	2.83	1.41	0								
P5	5.66	4.24	2.83	1.41	0							
P6	5.83	4.47	3.16	2.00	1.41	0						
P7	6.40	5.00	3.61	2.24	1.00	1.00	0					
P8	5.83	4.47	3.16	2.00	1.41	2.83	2.24	0				
P9	4.00	3.16	2.83	3.16	4.00	3.16	4.12	5.10	0			
P10	1.41	2.00	3.16	4.47	5.83	5.66	6.40	6.32	3.16	0		
P11	2.00	1.41	2.00	3.16	4.47	4.24	5.00	5.10	2.00	1.41	0	
P12	3.16	2.83	3.16	4.00	5.10	4.47	5.39	6.00	1.41	2.00	1.41	0

## Ridge Regression

## Lasso Regression

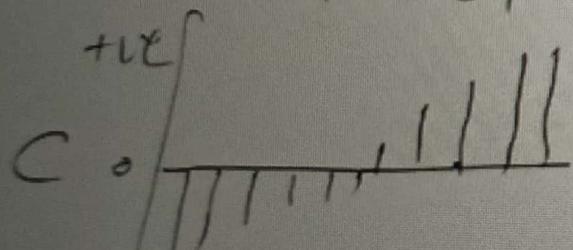
## Regularization

$$RR = \text{Loss} + \alpha \|w\|^2$$

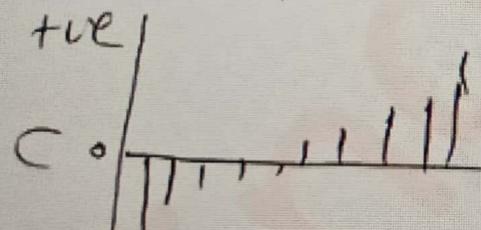
(penalty)

$$\|w\|^2 = w_1^2 + w_2^2 + w_3^2 + w_4^2 + \dots + w_n^2$$

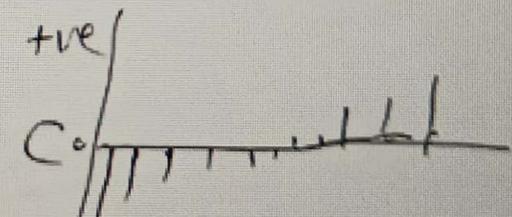
$$\alpha = 0.01$$



$$\alpha = 1$$

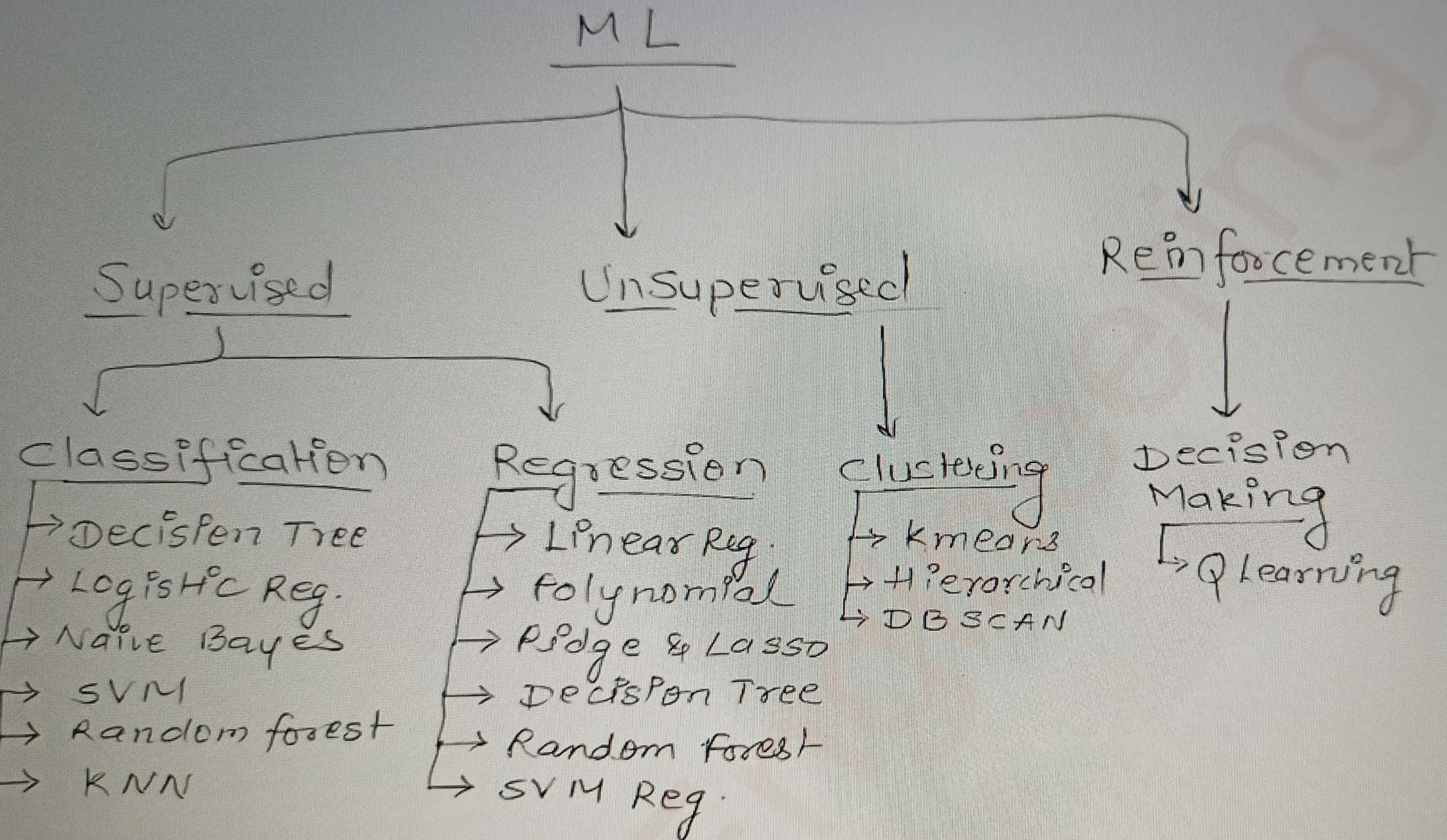


$$\alpha = 2$$



$$LR = \text{Loss} + \alpha \|w\|$$

$\nearrow$   
absolute  
value.



## Markov Decision Process (MDP)

- A Markov Decision Process (MDP) is a mathematical framework used to describe an environment in reinforcement learning. It provides a formal model for decision-making in situations where outcomes are partly random and partly under the control of a decision-maker (agent).

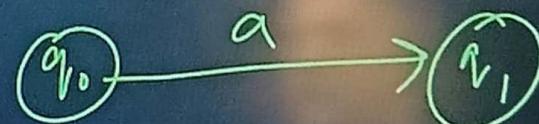
- Characteristics of MDP**

- Markov Property:**

- The future state depends only on the current state and action, not on the sequence of events that preceded it.
- This is known as the "memoryless" property.

- Policy ( $\pi$ ):**

- A strategy or rule that defines the action to be taken in each state.
- Example: A policy that always chooses the action that leads to the highest reward.



**1. States (S):**

- The set of all possible states in which an agent can exist.
- Example: Different rooms in a building.

**2. Actions (A):**

- The set of all possible actions an agent can take.
- Example: Moving from one room to another.

**3. Transition Function (P):**

- A probability function  $P(s'|s, a)$  that defines the probability of transitioning to state  $s'$  from state  $s$  when action  $a$  is taken.
- Example: The probability of successfully moving from one room to another.

**4. Reward Function (R):**

- A function  $R(s, a, s')$  that provides the immediate reward received after transitioning from state  $s$  to state  $s'$  due to action  $a$ .
- Example: Receiving a reward for successfully reaching a target room.

**5. Discount Factor ( $\gamma$ ):**

• A value between 0 and 1 that indicates the importance of future rewards.

## Challenges of Reinforcement Learning

### Reward Design:

- Designing appropriate rewards and determining their value can be complex. Example: In many games, setting suitable rewards is a significant challenge.

### Absence of a Model:

- Some environments lack a fixed structure or rules, requiring simulations to gather experience. Example: Unlike chess, many games need simulations because they have no underlying model.

### Partial Observability of States:

- Not all states are fully observable, leading to uncertainty. Example: In weather forecasting, complete information about the state is often unavailable.

### Complexity:

- Large board configurations and numerous possible actions increase complexity. Example: Games like Go have vast possibilities, making labeled data unavailable and increasing algorithmic complexity.

### Time-Consuming Operations:

- The need to explore extensive state spaces and possible actions increases the time required for learning. Example: Complex scenarios result in longer computational times, making the training process slow.

## Chapter 5 | Reinforcement Learning

s1	s2	s3	
s5	s6	s7	
	s10	s11	s12
s9			



Criteria	Supervised Learning	Reinforcement Learning	Unsupervised Learning
Mapping	Present	Present	Not present
Feedback	Instantaneous feedback once the model is created	Constant feedback from environment	No feedback from environment
Supervisor	Presence of supervisor and labeled data	No supervisor and labeled dataset is not available	No supervisor
Decisions	Independent and based on training input	Dependent and made sequentially	Independent
Examples	Classifiers (e.g., image classification)	Chess, Go Games, Robotics	Clustering (e.g., customer segmentation)

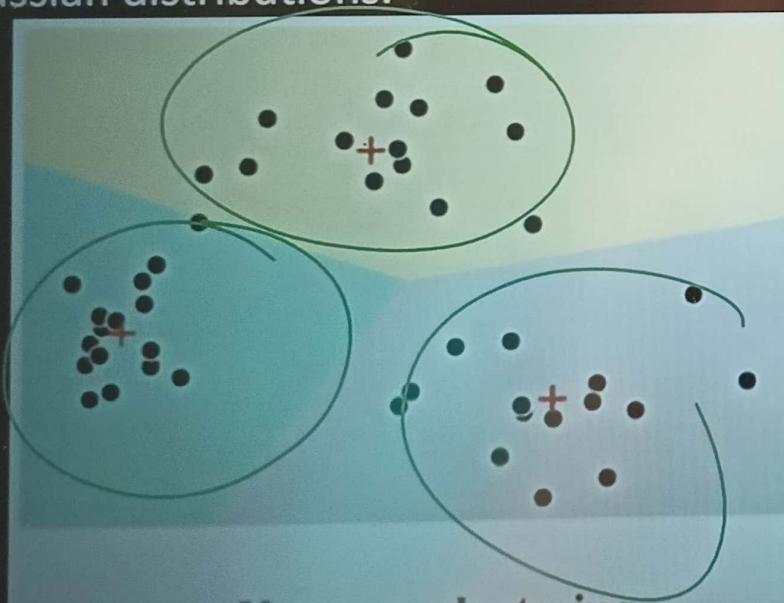
## Linear SVM Vs Non-Liner SVM

Conclusion Point	Linear SVM	Non-Linear SVM
Data Suitability	Best for linearly separable data	Suitable for complex, non-linearly separable data
Decision Boundary	Straight line or hyperplane	Non-linear boundary using kernel functions
Training Complexity	Simpler and faster	More computationally intensive
Kernel Usage	Not required	Requires kernel functions (e.g., RBF, polynomial)
Flexibility and Applications	Less flexible, used for simpler problems	Highly flexible, used for complex data patterns

Aspect	Linear Regression	Logistic Regression
Type of Model	Supervised regression model	Supervised classification model
Prediction Outcome	Predicts continuous values	Predicts binary outcomes (0 or 1)
Mathematical Model	Uses linear functions	Uses logistic functions with an activation function
Purpose	Estimates values of a dependent variable	Estimates probability of an event
Example	Predicting house prices based on size	Predicting whether a patient has a disease or not

## Partitional Clustering:

- **Centroid-based Clustering (e.g., K-means):** Partitions data into clusters, each represented by a centroid. Clusters minimize distance between data points and centroid, optimizing intra-cluster similarity and inter-cluster dissimilarity. For example, retail customers can be clustered by buying patterns, with each cluster's centroid reflecting average behavior.
- **Model-based Clustering:** Uses a statistical model for each cluster, finding the best data fit. For instance, Gaussian mixture models assume data points in each cluster are Gaussian distributed. This method is used in image processing to model different textures as coming from different Gaussian distributions.



K-means clustering



Mixture model (Gaussian)

**Example 3.75.** A coin is tossed  $(m + n)$  times, ( $m > n$ ). Show that the probability of at least  $m$  consecutive heads is  $\frac{n+2}{2^{m+1}}$ .

**Solution.** Since  $m > n$ , only one sequence of  $m$  consecutive heads is possible. This sequence may start either with the first toss or second toss or third toss, and so on, the last one will be starting with  $(n + 1)$ th toss.

Let  $E_i$  denote the event that the sequence of  $m$  consecutive heads starts with  $i$ th toss. Then the required probability is :  $P(E_1) + P(E_2) + \dots + P(E_{n+1})$ . ... (\*)

$$P(E_1) = P[\text{Consecutive heads in first } m \text{ tosses and head or tail in the rest}] = \left(\frac{1}{2}\right)^m$$

$$P(E_2) = P[\text{Tail in the first toss, followed by } m \text{ consecutive heads and head or tail in the next}] = \frac{1}{2} \left(\frac{1}{2}\right)^m = \frac{1}{2^{m+1}}$$

In general,

$$P(E_r) = P[\text{tail in the } (r-1)\text{th trial followed by } m \text{ consecutive heads and head or tail in the next}]$$

$$= \frac{1}{2} \left(\frac{1}{2}\right)^m = \frac{1}{2^{m+1}}, \forall r = 2, 3, \dots, n+1.$$

$$\text{Substituting in (*), required probability} = \frac{1}{2^m} + \frac{n}{2^{m+1}} = \frac{2+n}{2^{m+1}}.$$

## Central Limit Theorem

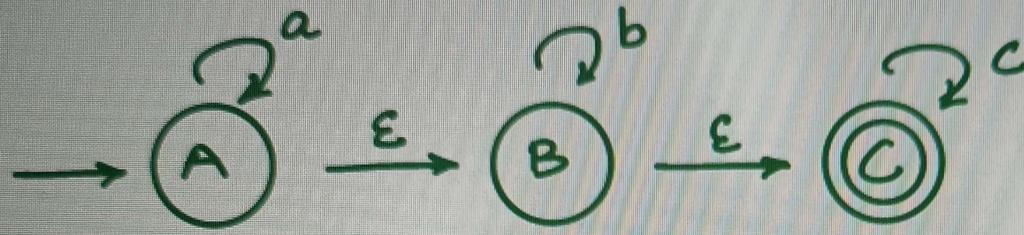
Theorem :- If  $\bar{X}$  is the mean of a random sample of size  $n$  taken from a population with mean  $\mu$  and variance  $\sigma^2$ , then the sampling distribution of the sample mean tends to be normal distribution with mean  $\mu$  and variance  $\frac{\sigma^2}{n}$  as the sample size tends to be large ( $n \geq 30$ ), regardless the form of the parent population,  
i.e  $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$ .

Note : ① In term of standard normal variate, we can write it as

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

- ② In case of  $n < 30$ , the approximation is good only if the population is not too different from a normal distribution.
- ③ If the population is to be normal, then the sampling distribution of  $\bar{X}$  is normal for any size of  $n$ .
- ④ CLT is really useful because it characterizes large samples from any distribution.

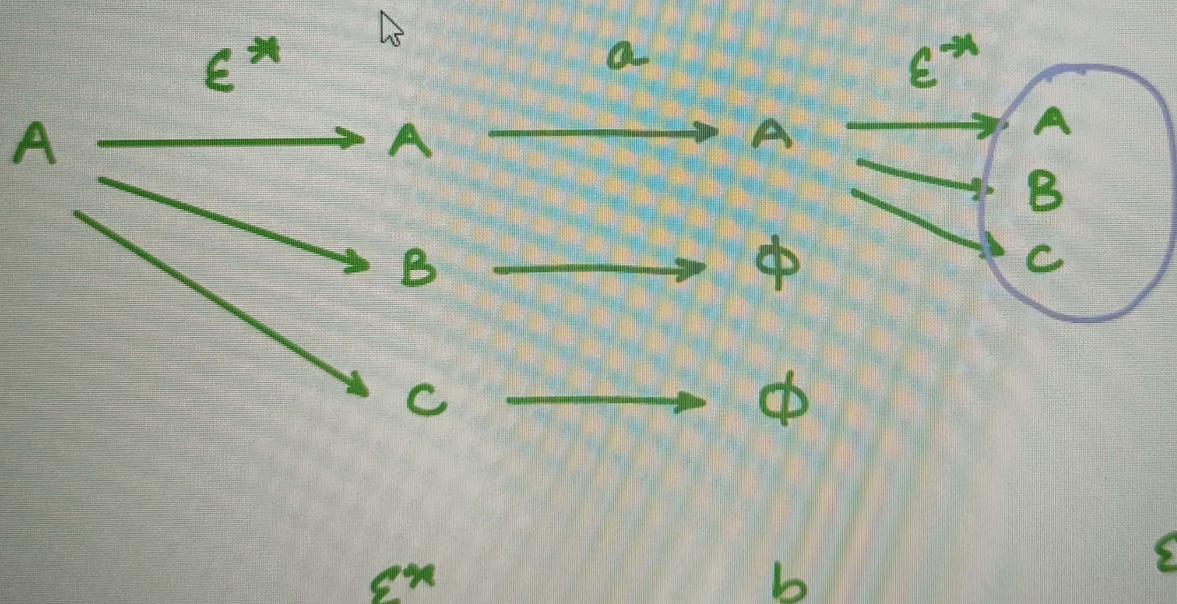
eq:



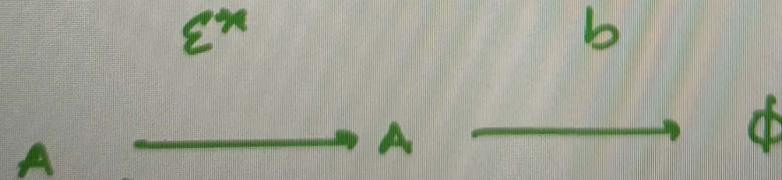
$$\Sigma = \{a, b, c\}$$

$\epsilon$ -closure(A) = AB  
 $\epsilon$ -closure(B) = BC  
 $\epsilon$ -closure(C) = C

A, a



A, b



$\epsilon^*$

## Machine learning

gradient descent - optimization algo

used to train model by

minimizing error b/w predicted & actual value.

one hot encoding → representation of categorical values as numerical values.

Regression → feature prediction from another feature.

Tree

→ Precision & Recall (phone)

new many predicted

To true pos.

True positive  
from emp!

TP	FP
FN	TN

$$F_1 = \frac{2 \cdot TP}{2 \cdot TP + FP + FN}$$

Precision → all predicted +ve,  
how many TP

Recall → all actual +ve,  
how many TP.

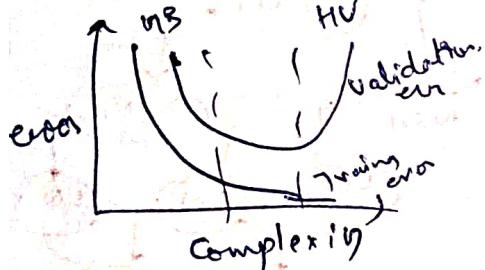
covariance

$$C(x,y) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

+ve, -ve, 0  
direction of linear relations  
+ve two variables

High bias underfit

Bias measured by how close the predictions match actual values



Bias  
↑ model  
complexity

Variance  
↓ model  
complexity  
↑ training  
data

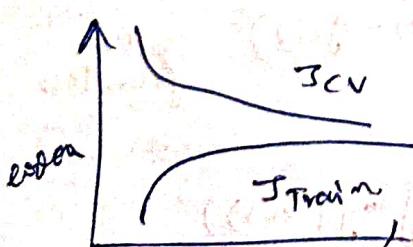
(Regularization)  $\lambda$   $\rightarrow$  Bias-Variance Trade-off  
 $J_{CV} = \text{Bias}^2 + \lambda \text{Var}$

Compute

$J_{CV}$  for

all  $\lambda$  and choose

$$\min J_{CV}$$



training  
example

high P, R, minimizes  
(FP, TP)

$$TPR = \frac{\#TP}{\# \text{actual} +ve}$$

$$FPR = \frac{\#FP}{\# \text{predicted} +ve}$$



Entropy  $\rightarrow$  degree of purity.  $H(X) = -\sum p(x_i) \log_2(p(x_i))$

(1G) Information gain measures changes in entropy due to others feature  $I(X;Y) = H(X) - H(X|Y)$

Mutual information  $MI = H(X) + H(Y) - H(X,Y)$   
information shared b/w  $X \& Y$ .

### Decision Trees

- ① SLT. for classification. Tree structured classifier, internal nodes represents features, branch represents decision rules & leaf nodes represent the outcome.

CART algos Classification & Reg. Tree Alg.

### Attribution Selection Measure

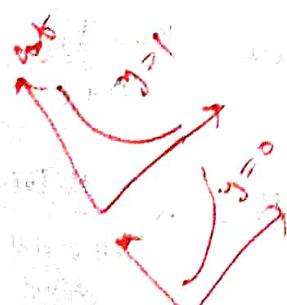
Information Gain

Tree Structure, splitting,

- ② pruning  $\rightarrow$  remove nodes which add less value.

Impurity

Gini index  $\rightarrow$  probability of misclassification in CART



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

cost function

$$\begin{aligned} \text{Cost for logistic} &= -y \log(h(x)) \\ \text{for reg} &= -(1-y) \log(1-h(x)) \end{aligned}$$

Convergence

### Naive Bayes

based on Bayes Theorem

## # classification

$$P(AB) = P(B/A) \cdot P(A)$$

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

## K-nearest neighbour

- ① find K nearest points and  
find its class category but euclidean distance

Identify category of KNN.

For Regression: find avg.

## Limitations

- not for huge dataset
- outliers sensitive
- sensitive to missing value.

DECISION non parametric

## TREE

Information =  $I(Total)$

Gain

$$= \sum n_i \text{Ent}(E_i)$$

$$= \frac{1}{H_{\text{Total}}} \sum H_i$$

Impurity

pure → only one class per split

split with 1/N.

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

$$F_1 = \frac{(1+\beta^2)}{\beta^2(P+R)} \frac{P \cdot R}{P+R}$$

$$\rightarrow FP \gg F.N, \quad \beta = 1/2$$

$$\rightarrow FN \gg FP, \quad \beta = 2$$

$$\rightarrow FP, FN \quad \beta = 1$$

- ① find entropy & information
- ② gain & pick maximum or
- ③ Do further, make tree.

$$(IG) \rightarrow \sum p_i \log p_i$$

$$(G_L)$$

$$[0, 0.5]$$

$$n(S) \rightarrow (0, 1)$$

No missing data, may have overfitting.

Naive Bayes → all independent

and & test weaker  
spam class pred.

## Pruning

→ maximum levels / depth

→ minimum threshold.

→ node becomes pure

naive bayes

probabilistic

classifier based  
on Bayes Thm.

(i) ID3

used  
entropy

& (IG)

stopping  
criterion

Tree Pruning

## Pruning

PrePruning

→ Early stopping

→ faster

Post Pruning

pruning  
after full  
growth

Better  
Generalization

→ no  
overfit

## Linear Discriminant Analysis (LDA)

→ dimensionality reduction technique

→ high dimensional space to lower Dim. and project  
into 1D line → separating line → separating lines of features

### Assumptions

→ Gaussian Distribution → Covariance are equal matrices

→ Linearly Separable

- Maximize Distance b/w means
- Minimize Variation b/w classes

### Methods to avoid

#### Overfitting in Decision Trees

1. Pruning Techniques → remove parts not necessary
2. Limit tree depth → reduce complexity
3. Cross Validation → better generalization
4. Feature Selection → remove noisy ones, select appropriate

find a line for best fit for set of points (Least Square Techniques)

$$\# a^n \cdot b^m \cdot c^n \mid n, m \geq 1$$

alphabet, stack top i move

$$S \rightarrow aR \cup I \cup aSr$$