# [Machine Learning]

Machine learning (ML) is a subfield of artificial intelligence (AI) that focuses on developing algorithms and model that can automatically learn patterns and insights from data withought being explicitly programmed. The goal of ML is to enable computers to learn and improve from experience, just like a human.
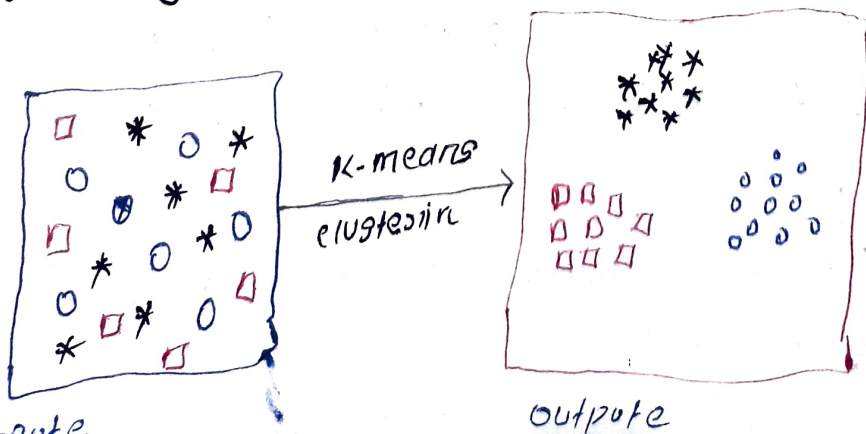
There are three main types of machine learning:

Supervised Learning: Supervised learning is a type of machine learning that uses labeled data to train machine learning model. A labelled dataset is one that has both inpute and output parameter. The model just need to map the inpute to the respective output or target value.

Supervised Learning algorithms:
- Linear Regression
- Logistic Regression
- SVM (support vector machine)
- KNN
- Decision Tree
- Naive Bayes.

Unsupervised Learning: Unsupervised learning is a type of machine learning that uses unlabeled data to train machine. unlabeled data doesn't have a fixed output variable. The model learn from the data, discovers the patterns and feature in the data, and return outpute.

- clustering
- PCA
- K-means clustering
- Hierarchical clustering



input → K-means clustering → outpute

**Reinforcement Learning:** In reinforcement learning, the algorithm learns by interacting with an environment and receive feedback in the form of rewards or punishment. The algorithms learn to take action that maximize the cumulative reward over time. This type of learning is often used in robotic, gaming and control system.

## Well posed Problems

Well-posed learning problem in machine learning is a problem that is well-defined, has a clear objective and has feasible solution.

1) Well-defined problem statement: The problem statement should be clear and unambiguous. The inputs, outputs and the objective of the problem should be well-defined.

2) Accessible and representative data: The data used to train the model should be representative of the real-world problem and should be accessible to the algorithms. Sufficient data required.

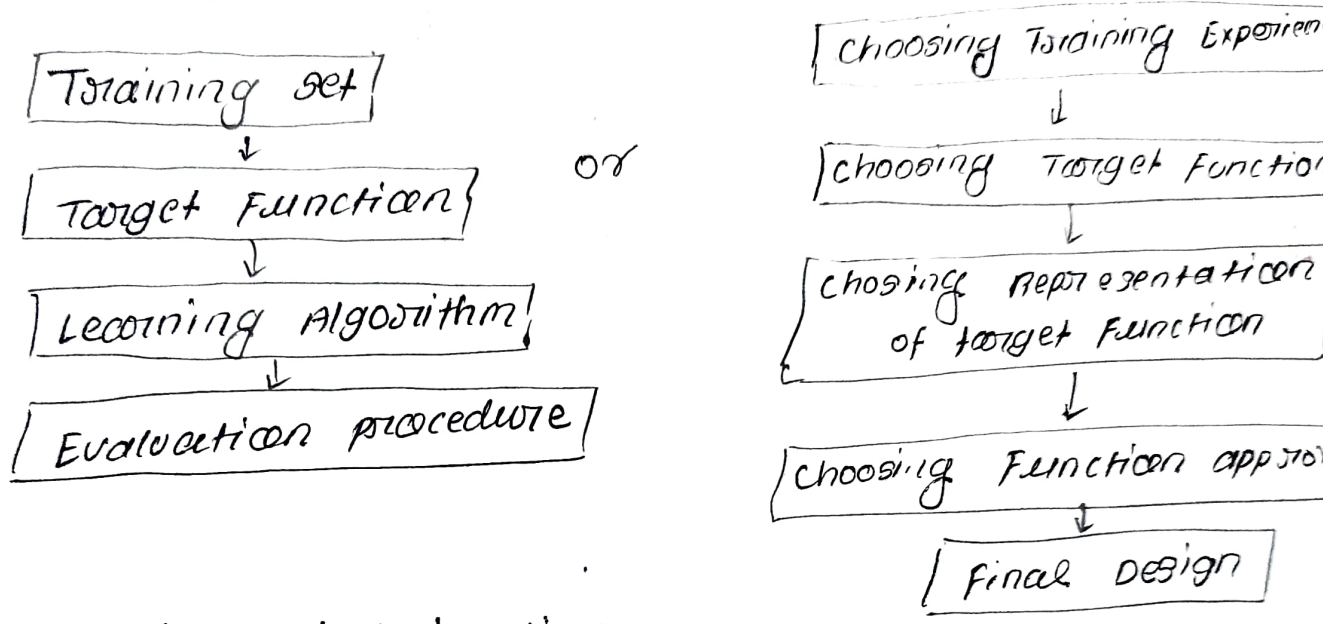3) Appropriate evaluation metrics: The evaluation matrix should be appropriate for problem at hand.

4) Feasible solution: The solution should be practical and should be achievable within a reasonable timeframe.

5) Generalizability: The solution obtain should be able to generalize well to unseen data.

- image classification
- sentiment analysis
- recomendation system

posed learning problem include:
- House price prediction
- image caption Generation

# Designing learning system

| Training set |

↓

| Target Function |

↓

| Learning Algorithm |

↓

| Evaluation procedure |

or

| choosing Training Experien |

↓

| choosing Target Function |

↓

| chosing Representation of target Function |

↓

| choosing Function appro |

↓

| Final Design |

# Empirical risk minimization

Empirical risk minimization (ERM) is a common approch i supervised learning, which involves minimizing the average los over a training dataset to find the best model parameters

on ERM, a model is trained on a training dataset, which con of input-output pairs. The model is then evaluated using a function, which measure the difference between the predic output to the truth output. The goal of ERM is to find model parameters that minimize the average loss over training dataset.

Given a set of training data $\{(x_1, y_1), (x_2, y_2), (x_3, y_3) \ldots (x_n, y_n)\}$

where $x_i$ is the input and $y_i$ is the output, and model parameterized by $\theta$, we want to find the value of $\theta$ that minimizes the average loss over training data.

$$L_S(h) = \frac{1}{n} \sum_{i=1}^{n} L(h(x_i), y_i) = L_S(\theta) \sum_{i=1}^{n} L(\theta(x_i), y_i)$$

$$\text{minimize } \theta \frac{1}{n} \sum_{i=1}^{n} L(y_i, f(x_i; \theta))$$

$L$ = loss function
$f$ = model $\quad \theta$ = model parameter

## PAC (Probably Approximately correct).

It is a theoretical framework in machine learning that provides bounds on the number of training example need to learn a concept to a certain degree of accuracy whith high probability. The goal of PAC Learning is to find a hypothesis the true concept underlying a set of training example.

PAC is a way to measure how well a machine learning algorithm can learn a concept from a set of examples. ot like trying to learn new language by looking at example of sentences in that language.

PAC learning tell us that if we have set of example and a hypothesis (a guess about what the answer is). we can measure how over clouse over guess is to the truth answer. The goal is to find a hypothesis that is close to the true answer.

[Data preprocessing]
- Data cleaning
- Data integration
- Dimension reduction
- Feature extraction
- Data transformation

- Data spliting
- one hote encoding
- word embeding
- Tokenization
- Normalization

## Normalization

Normalization is a data preprocessing technique that involves scaling numerical feature in dataset to a standard range. The goal of normalization is to bring all the feature to a similar scale. so that no single feature domina to others, and to make it easier for machine learning alg to learn from the data.

# [Gradient Descent]

Gradient descent is a genetic optimization algorithm capable of finding optimal solution to a wide range of problem.

The general idea of gradient descent is to tweak parameters iteratively in order to minimize a cost function.

## ① Batch Gradient Descent

$$\frac{\partial}{\partial \theta_j} MSE(\theta) = \frac{2}{m} \sum_{i=1}^{m} (\theta^T x^i - y^i) x_j^i$$

Gradient vector of the cost function

$$\nabla_\theta MSE(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} MSE(\theta) \\ \frac{\partial}{\partial \theta_1} MSE(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n} MSE(\theta) \end{pmatrix} = \frac{2}{m} X^T (X\theta - y)$$

notice that this formula involves calculation over the full training set X at each gradient descent step! This is why the algorithm is called batch gradient descent: of uses the whole batch of training data at every step (actually. full gradient decent would probably be a better name)
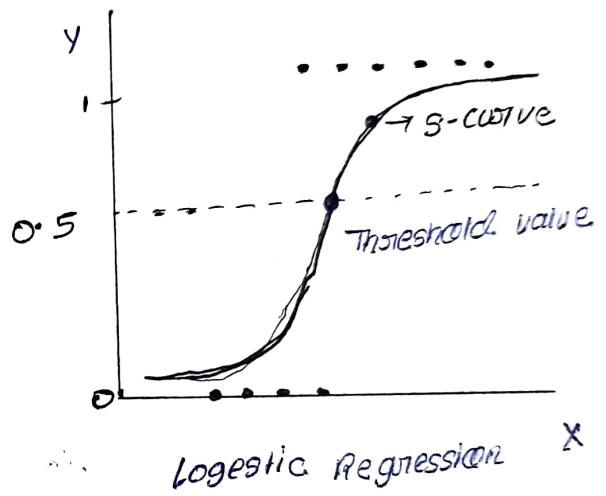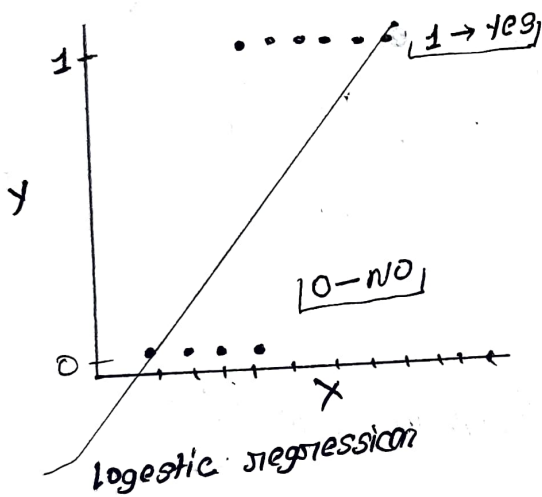
# [Logistic Regression]

Logistic regression is a statistical method used to predict the likelihood of a binary outcome like "yes" or "no".

It is a type of regression analysis that predicts the probability of occurrence of a categorical dependent variable based on one or more independent variables.

$$\text{logistic Regression} = \boxed{y = \frac{1}{1 + e^{-x}}}$$

Sigmoid



logestic regression

Here logestic regression is fail.



Logestic Regression

The value of the logistic must be between 0 and 1, which canno go beyond this limit, so it forms a curve like "S" form.

The S-form curve is called the sigmoid function or the logestic function.

sigmoid function help to convert the linear combination of predictor variable into probability estimate that use to make prediction.
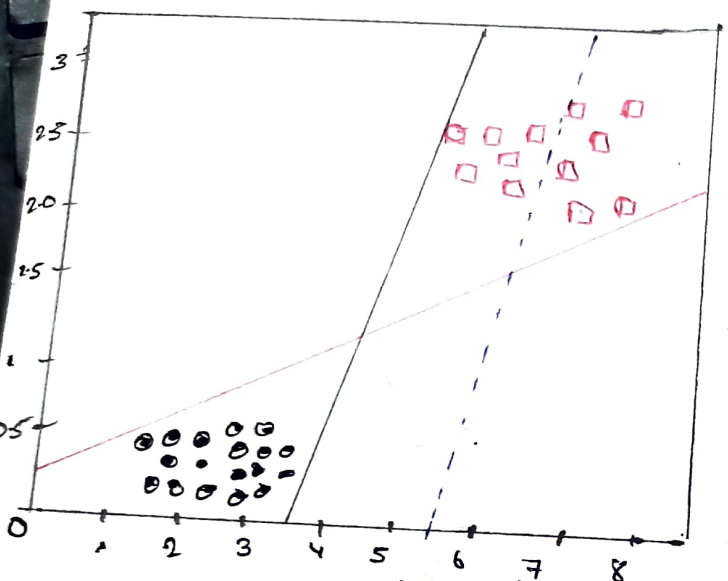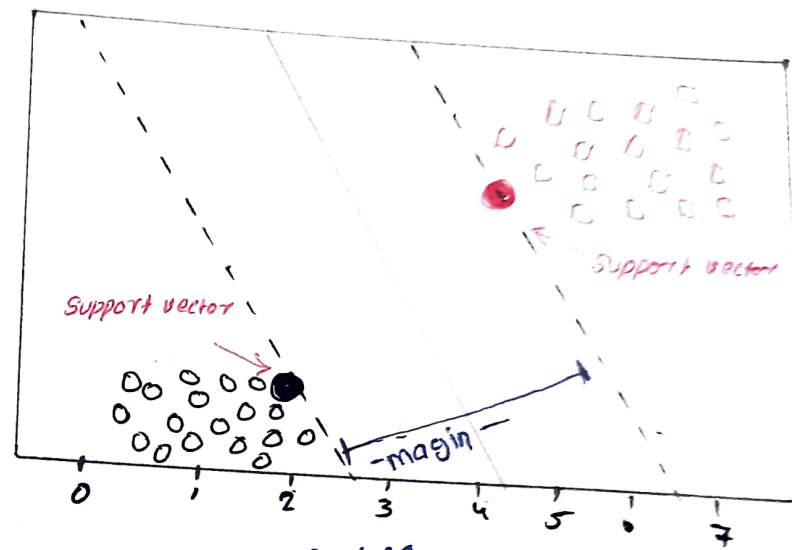
# [Support Vector Machines]

A support vector machine (SVM) is a powerful and versatile machine learning model capable of performing linear or nonlinear classification, regression and even novelty detection.

Snslitive
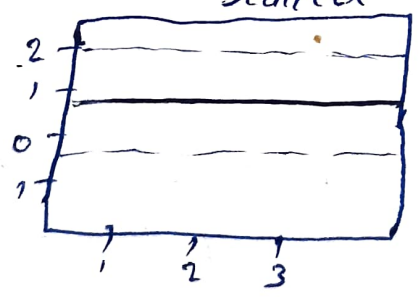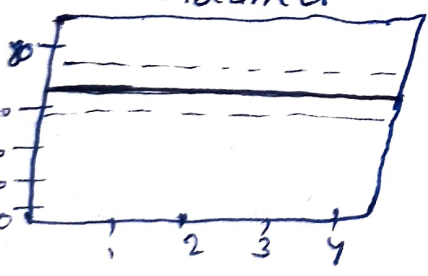it perform bad with outlier

outlier

$$\frac{2}{||w||}$$



linear classification
of 3 model

SVM

SVM is a maximum margin classifier. #[Go through book.]
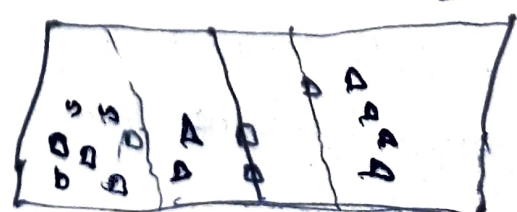
it is sensitive to the feature scales

unscalled          scalled



Sensitive to outlier henie we use soft margin instant of hard margin.

C₁                    C=100      [it underfit then reduce]

KNN, short for k-nearest neighbours, is a machine learning algorithm used for classification and regression problem. The KNN algorithm works by finding the k closest data point to the given inputs, and then using those data points to make prediction.

## Advantage of KNN

- Simple to understand and easy to implement
- No training period: unlike other machine learning algorithm that require training on a dataset.

KNN does not required a training period. It uses the the entire dataset for making prediction.

- Good performance on small dataset.
- works well with non-linear data.

· Euclidean Distance

$$d(p,q) = \sqrt{\sum_i (p_i - q_i)^2}$$

## Disadvantage

- computationally expansive
- sensitive to outliers
- Requires normalization
- The choice of k: The performance of KNN is highly depend on the choice of k, which can be a chanllenging task.

case. choosing a larger k-value can lead to smoother decision boundaries, but it also increase risk of ~~overfitting~~. underfitting.

k-less chance of ~~underfitting~~ overfitting.

# [Bayesian learning]

Bayesian learning is a machine learning approach that involves using bayesian statistics to make predictions and decisions

Bayesian learning is based on the Bayes theorem, which is a mathematical formula that describes the probability of an event based on prior knowledge and new evidence.

$P(x|y) \rightarrow$ probability of $y$ given $X$

Bayes theorem

$$P(x|y) = \frac{P(y|x)(P(x)}{P(y)}$$

$P(x||y) \rightarrow$ is called a **posterior**, which we need to calculate

$P(y|x) \rightarrow$ is called the **likelihood**. ot is the probability of evidence when hypothesis is true.

$P(x) \rightarrow$ called **prior probability**, probability of hypothesis before considering the evidence.

$P(y) \rightarrow$ called **marginal probability**.

(Q)

$P(cancer) = (0.008)$

$P(pos / cancer) = 0.98$

$P(pos / \sim cancer) = 0.03$

@ if a new patient comes in with a positive test result, what is the probability that he has cancer?

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

$$P(cancer / pos) = \frac{P(POS/cancer) * P(cancer)}{P(POS)}$$

$$P(POS) = P(pos/cancer) * P(cancer) + P(pos/\sim cancer) * P(\sim canc$$

$P(\sim cancer) = 1 - P(cancer) \rightarrow$ probability of having no cancer

substituting these value we will get --

$$P(cancer / POS) = \frac{(P(pos/cancer) * P(cancer))}{P(pos/cancer) * P(cancer) + P(pos/\sim cancer}{1 - P(cancer))}$$

$$= \frac{0.98 \times 0.008}{0.98 \times 0.008 + 0.03 * (1 - 0.008)}$$

$P(can/pos) = 0.208$

$$P(cancer \mid positive) = \frac{P(positive \mid cancer) \times P(cancer)}{P(positive)} \quad (1)$$

we have to find $P(positive)$

$$P(positive) = P(positive \mid cancer) \times P(cancer) + P(positive \mid \sim cancer) \times P(\sim cancer)$$

$$P(\sim cancer) = 1 - P(cancer)$$
$$\Rightarrow 1 - 0.008 = 0.992$$

$$P(positive) = P(positive \mid cancer) P(cancer) + P(positive \mid \sim cancer) P(cancer)$$

$$= 0.98 \times 0.008 + 0.03 \times 0.992$$

$$\Rightarrow 0.0376$$

$P(positive) = 0.0376$

we can put in equation (1)

$$= P(cancer \mid positive) = \frac{P(positive \mid cancer) \times P(cancer)}{P(positive)}$$

$$= \frac{0.98 \times 0.008}{0.0376}$$

$$\approx 0.208 \quad \sim$$

Decision tree are versatile machine learning algorithms that can perform both classification and regression tasks and even multioutput tasks.

## CART (classification and Regression Tree)

* CART algorithm is a greedy algorithm.

cost function for classification

$$J(k, t_k) = \frac{m_{left}}{m} G_{left} + \frac{m_{right}}{m} G_{right}$$

where $G_{left/right}$ → measures the impurity of the left/right subset

$M_{left/right}$ — is the number of instances in the left/right sub

The algorithm work by first splitting the training set into two subset using a single feature $k$ and a threshold $t_k$

it search for the pair $(k, t_k)$ That produces the purest subset weight by their size.

once the CART algorithm has successfully split the training
set in two, it splits the subset using the same logic, then the
sub-subset, and so on, recursively. It stops recursion once
it reaches the maximum depth (defined by the max-depth hyperparameter)
or if it cannot find the split that will reduce impurity.

## ID3 (Iterative Dichotomiser 3)

It as decision tree algorithm used in ML for classification task.
The algorithm builds a decision tree by recursively selecting
the input feature that provides the most information gain
about the target variable, based on the entropy or impurity
of the dataset.

### Disadvantage
- Attributes must be in nominal values
- Dataset must not include missing data
- The algorithm tend to fall into over-fitting.

─── ✗ ─── ✗ ─── ✗ ─── ✗ ─── ✗ ───

$$\left[ Entropy(Decision) = \sum P(I) \cdot \log_2 P(I) \right]$$

$$\left[ Gain(S,A) = Entropy(S) - \sum P(S|A) \cdot Entropy(S|A) \right]$$

**(Q)** Data set given

- Firstly, we need to calculate global entropy.

There are 14 example; 9 instances says → yes
5 instances says → no

$$\text{Entropy (Decision)} = \sum - p(i) \cdot \log_2 p(i)$$

$$= - p(yes) \cdot \log_2(yes) - p(no) \cdot \log_2 p(no)$$

$$= -(9/14) \cdot \log_2(9/14) - (5/14) \cdot \log_2(5/14)$$

$$= 0.9401$$

**1) # Wind factor on decision**

$$\text{Gain (Decision, wind)} = \text{Entropy(Decision)} - \sum [P(\text{Decision}|w] \cdot \text{Entropy(Decision)}$$

⇒ Wind attribute has two labels: weak and strong.

calculate (Decision | wind = weak) and (Decision | wind = strong)

| Day | wind | Decision |
|---|---|---|
| 1 | week | no |
| 3 | week | yes |
| 4 | week | yes |
| 5 | week | yes |
| 8 | week | no |
| 9 | week | yes |
| 10 | week | yes |
| 13 | week | |

there are instances for weak wind. Decision of 2 items are no and 6 items are yes.

Entropy (Decision | wind = weak) = $-P(no) \cdot \log_2 P(no) - P(yes) \cdot \log_2 P(yes)$

$$= -\left(\frac{2}{8}\right) \cdot \log_2 \left(\frac{2}{8}\right) - \left(\frac{6}{8}\right) \cdot \log_2 \left(\frac{6}{8}\right) = 0.811$$

[Strong wind factor on decision]

Entropy (Decision | wind = strong) = $-P(no) \cdot \log_2 P(no) - P(yes) \cdot \log_2 P(yes)$

$$= -\left(\frac{3}{6}\right) \cdot \log_2 \left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \cdot \log_2 \left(\frac{3}{6}\right) = 1$$

Gain (Decision, wind) ▬ [P(Decision | wind = weak) . Entropy (Decision

wind = weak)] − [P(Decision | wind = strong) . Entropy (Decision | wind

= strong)]

Entropy of Strong wind

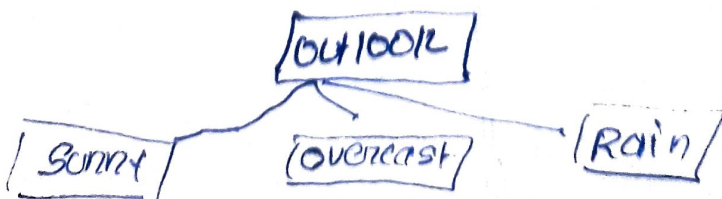$$= 0.940 - \left[\left(\frac{8}{14}\right) \cdot 0.811\right] - \left[\left(\frac{6}{14}\right) \cdot 1\right] = 0.048$$

weak wind    weak + strong    strong    Entropy of strong wind
              wind

[Information Gain of all attributes]

1) Gain (Decision, outlook) = 0.246

2) Gain (Decision, Temperature) = 0.029

3) Gain (Decision, Humidity) = 0.151

4) Gain (Decision, wind) = 0.048

[outlook]

[Sunny]    [overcast]    [Rain]

# # [CART Algorithm]

it can handle both classification and regression and regression tasks.

This algorithm uses a new metric named gini index to create decision points for classification tasks.

## Gini index

$$Gini = 1 - \sum (P_i)^2 \text{ for } i=1 \text{ to number of class.}$$

| Outlook | yes | no | number of instances | Gini |
|---|---|---|---|---|
| sunny | 2 | 3 | 5 | 0.48 |
| overcast | 4 | 0 | 4 | 0 |
| Rain | 3 | 2 | 5 | 0.48 |

$Gini \, (outlook \neq sunny) = 1 - (2/5)^2 - (3/5)^2 = 1 - 0.16 - 0.36 = 0.48$

$Gini \, (outlook = overcast) = 1 - (4/4)^2 - (0/4)^2 = 0$

$Gini \, (outlook = Rain) = 1 - (3/5)^2 - (2/5)^2 = 1 - 0.36 - 0.16 = 0.48$

$Gini \, (outlook) = (5/14) \times 0.48 + (4/14) \times 0 + (5/14) \times 0.48$

$= 0.171 + 0 + 0.171$

$= 0.342$

| Temperature | Yes | No | No of instances |
|---|---|---|---|
| Hot | 2 | 2 | 4 |
| Cool | 3 | 1 | 4 |
| Mild | 4 | 2 | 6 |

$Gini\,(Tem = Hot) = 1 - (2/4)^2 - (2/4)^2 = 0.5$

$Gini\,(Tem = cool) = 1 - (3/4)^2 - (1/4)^2 = 1 - 0.562 - 0.0625 = 0.375$

$Gini\,(tem = mild) = 1 - (4/6)^2 - (2/6)^2 = 0.0439$

$Gini\,(tem) = (4/14) \times 0.5 + (4/14) \times 0.375 + (6/14) \times 0.439$

$= 0.439$

Similarly we calculate for all coloums.

| Feature | Gini Index |
|---|---|
| Outlook | 0.342 |
| temp | 0.439 |
| Humidity | 0.367 |
| wind | 0.428 |

The winner is outlook because it cost is minim
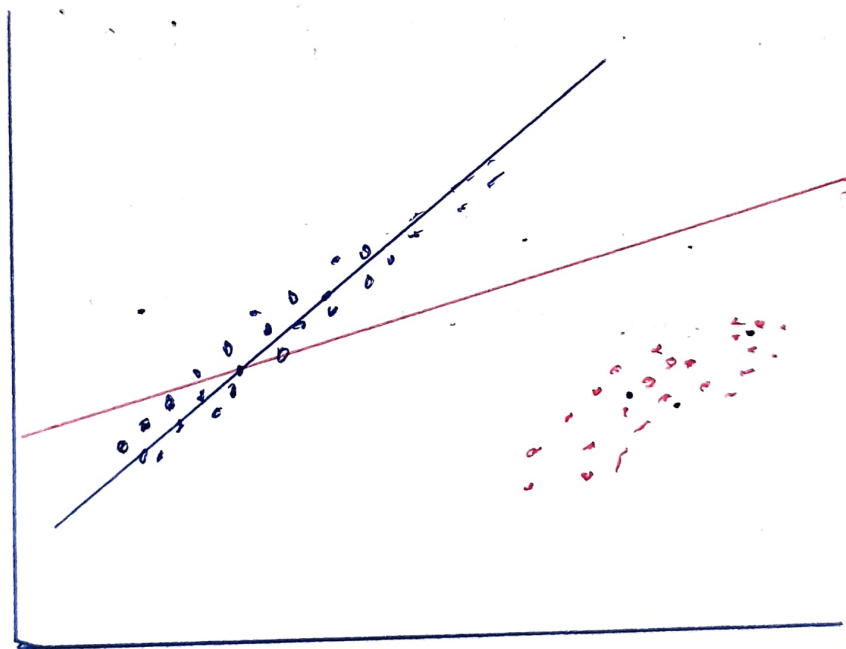
we can see how robustly fit a linear model to faulty data using RANSAC algorithm.

The ordinary linear regressor is sensitive to outliers, and the fitted line can easily be skewed away from the truth underlying relationship of data.

The RANSAC regressor automatically split the data into inliers and outliers, and the fitted line is determined only by the identified inliers.



**... → inliers**
**... → outliers**
**— RANSAC regressor**
**--- Linear regressor**

RANSAC = (Random Sample consensus)
fit a model from random subsets
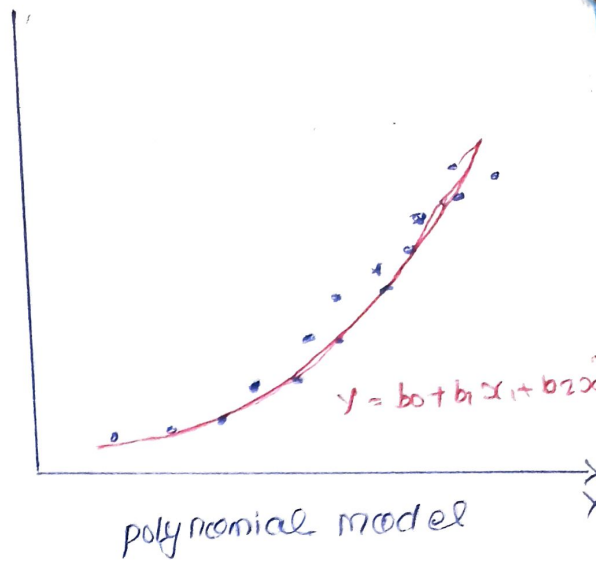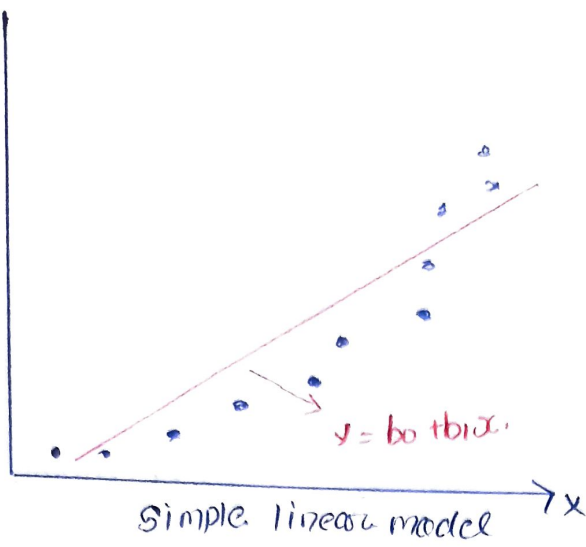of inliers from the complex data set.

a reasonable result with a certain probability, which is dependent on the number of iterations (max-trials parameter). It is typically used for linear found non-linear regression problems and is especially popular in the field of photogrammetric computer vision.

# Polynomial Regression

Polynomial Regression is a regression algorithm that models the relationship between a dependent (y) and independent variable as $n^{th}$ degree polynomial.

$$\text{Equation} = b_0 + b_1 x_1 + b_2 x_1^2 + b_3 x_1^3 - - b_n x_1^n$$

- It is also called the special case of multiple linear regression
- It is a linear model with some modification in order to increase the accuracy.

- It work well on non-linear data

- In polynomial regression, the original features are converted into polynomial feature of required degree (2, and then modeled using a linear model.

$y = bo + b_1 x$.

Simple linear model

$y = bo + b_1 x_1 + b_2 x$

polynomial model

## Equation of the polynomial Regression model

Simple linear equation : $y = bo + b_1 x$

multiple — — — — : $y = bo + b_1 x + b_2 x_2 + b_3 x_3 - -$

polynomial — — — — : $y = bo + b_1 x_1 + b_2 x_1^2 + b_3 x_1^3 -$

# # ARIMA

Autoregressive integrated moving Average is a model
use to analyze and forcast time series data.

Time series data is a sequence of observations record
over time, such as stock prices, temperature reading, or sal

**i) Autoregressive (AR) component:** This component looks at the relationship between an observation and a certain number of past observation.

**ii) Integrate (I):** it deals with transforming the data to make it stationary

**iii) Moving Average (MA)** it is the influence of past error term in predecting future values.

it consider the average of the error and their relationship to previous error.

By combining these components, ARIMA helps cature the patterns, trends and dependencies present in the time series data.

**SARIMA:** Seasonal (ARIMA) is a extension of ARIMA that incorporates seasonality into the analysis.

# Regularized methods for Regression

Regularized method for regression are techniques used to mitigate (reduce) overfitting and improve the generalization performance of regression model. These methods add a regularization term to the loss function during model training, which help control the complexity of the model.

There are two commonly used regularized method for regression

i) Ridge Regression (L2 Regularization)

ii) Lasso Regression (L1 Regularization)

iii) Elastic net (Ridge + Lasso)

i) Ridge Regularization: Ridge Regression adds a penalty term to the loss function, which is proportional to the sum of square coefficient of the regression model.

This penalty term discourages large coefficient values, making the model less sensitive to individual datapoint and reducing the likelihood of overfitting.
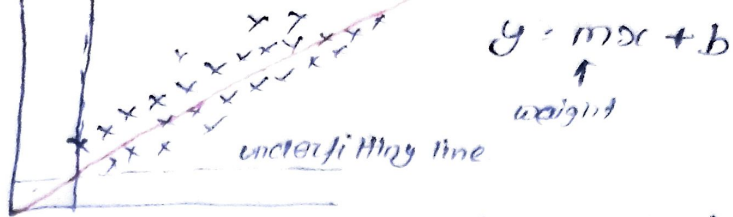
Loss function of ridge = Loss + alpha * (sum of square of
$$\underset{\text{lemda}}{\text{or}}$$ cofficients)

$$LOSS = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

Loss function of ridge = $\sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda (m^2)$

alpha is a hyperparameter which control the strength of regulari2 alpha → more shrinking of cofficient to zero, reducing model complexity.

y = mx + b

$$y = mx + b$$

weight ↑ (pointing to m)

underfitting line

when $m >>> b$ then overfitting

when $b >>> m$ then underfitting

when $m ≈ b$ then best fit.

training data

→ best fit line according to trailing data but overfitting, not perform good on testing data

→ Ridge regularization line reduce overfitting problem

• → training data
x → testing data

→ Shrinkage cofficient

$$\text{Ridge LOSS} = \sum_{i=1}^{n} (Y_i - \hat{Y_i})^2 + \boxed{\lambda ||W||^2}$$

$$\lambda (\omega_1^2 + \omega_2^2 + \omega_3^2 ---)^2$$

$$\lambda → 0 \text{ to } \infty$$

Lasso regression also add a penalty term to the loss function, but instead of the sum of squared coefficient, it uses the sum of the absoulte values of the coeffients.

This encourages sparsity in the model, meaning that some coeffients may become exactly zero.

Lasso regression is useful for feature selection, as it tends to set less important feature to zero, effectively removing them from model.

$$\sum_{i=1}^{n} (y_i - \hat{y_i})^2 + \lambda ||w||$$

### 3) Elastic Net

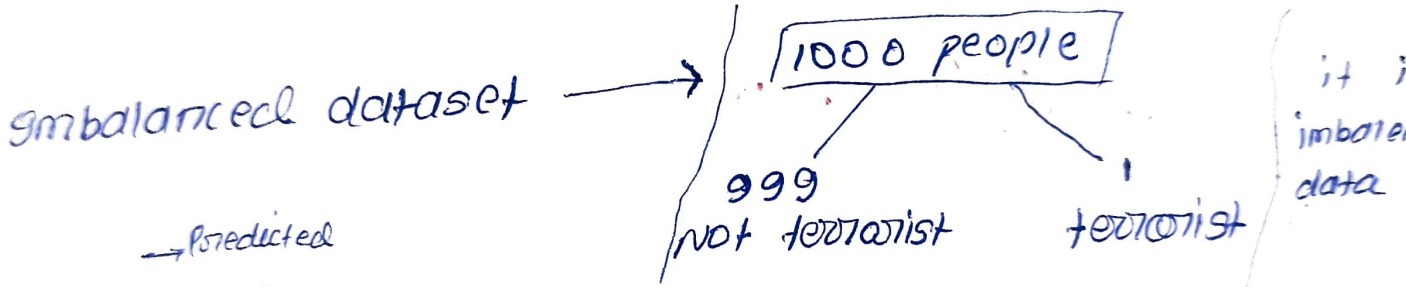$$\sum_{i=1}^{n} (y_i - \hat{y_i})^2 + \lambda ||w||^2 + \lambda ||w||$$

it is used to calculate the relationship between different features in the dataset.

$$-1 \text{ to } 1$$

1) positive correlation $(>0)$

2) negative correlation $(<0)$

3) no correlation $(0)$

## # confusion Matrix

imbalanced dataset $\longrightarrow$ | 1000 people |

it is imbalance data

999 not terrorist          terrorist

→ Predicted



Accuracy $= \dfrac{999}{999+1} \approx 99.9\%$ accuracy

here accuracy is 99.9% according to it but it is not g

way to mesure accuracy because data is imbalance

**Precision:** what proportion of predicted positive is really positive

(jitne chizho ko positive bola hai usme se kitne chijhe sach me positive hai)

or

predicted positive value is what percent really positive



Predicted

|        |          | negative       | positive        |
|--------|----------|----------------|-----------------|
| actual | negative | True negative  | False positive  |
|        | positive | False negative | True positive   | recall

precision

$$Precision = \frac{True\ positive}{True\ positive + False\ positive}$$

$$= \frac{True\ positive}{Total\ predicted\ positive}$$

**II) Recall:** what proportion of actual positive is correctly classified

jitne logo ko sach me cancer that usme se kitno logo ko detect kiya gata.

$$Recall = \frac{True\ positive}{True\ positive + False\ negative}$$

$$= \frac{True\ positive}{Total\ actual\ positive}$$

$$F_1 = 2 \times \frac{\text{precision} \times \text{Recall}}{\text{precision} + \text{Recall}}$$

# No Free Lunch (NFL) theorem is a fundamental

oresult in the field of machine learning and optimization
or essentially states that there is no on' universal algorithm
that performs optimally. across all possible problem domain
oor data distributions.

it argues that when averaging over all problem instances,
every algorithm will have the same performance on average .
on other word , no algorithms can outperform any other
algorithm when considering the entire space of possible probl

The NFL theorem also highlights the important of
understanding the problem at hand , exploring different
algorithm and adapting them to specific problem.

# Error decomposition

Error decomposition refers to the process of breaking down the total error in a prediction or estimation task into different component of understanding their individual contribution.

1) Bais :

2) variance :

3) irreducible Error : The irreducible error represents the inherent noise or uncertainty in the data that cannot be reduced, even with perfect model.

# Performance measure

The ultimate goal of performance measure is to reduce errors.

1) MAE (Mean absolute error (MAE)

$y_i$ - target value     $\hat{y_i}$ = predicted value

$$\frac{\sum\limits_{i=1}^{n} |y_i - \hat{y_i}|}{n}$$

• Disadvantage
it is modulous it is not differencia at 0

• Advantage
1) Same unit as y-axis
2) Robust to outliers

2) MSE (mean square error)

$$\sum_{i=1}^{n} \frac{(y_i - \hat{y_i})^2}{n}$$

adv



Differentiable at
all points

dis

• not good for outliers



covering more
area outliers

• unit of error = square of y- axis unit

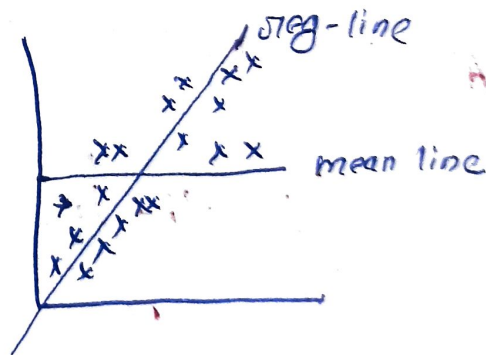RMSE $= \sqrt{\sum_{i=1}^{n} \frac{(y_i - \hat{y_i})^2}{n}} = \sqrt{MSE}$

•Same unit as y.

→ coefficient of determination ($R^2$) or goodness of fit

$$R^2 = 1 - \frac{SSR}{SSm}$$

$SS_R$ = Sum of square error
in the regression line

$SSm$ = Sum of square error
in mean line



reg-line

mean line

$$R^2 = 1 - \frac{\left[\sum_{i=1}^{n} (y_i - \hat{y_i})^2\right]_{reg}}{\left[\sum_{i=1}^{n} (y_i - \bar{y_i})^2\right]_{mean}}$$

5> MSLE (mean square logerithimic error)

$$MSLE = \frac{1}{n} \sum_{i=1}^{n} \left( \log(1 + \hat{t}_i) - \log(1 + t_i) \right)^2$$

6> MAPE (mean absolute percentage error) also known as MAPE

$$MAPE = \frac{100\%}{n} \sum_{t=1}^{n} \left| \frac{A_t - F_t}{A_t} \right|$$

$A_t \rightarrow$ Actual value

$F_t \rightarrow$ Forecast value

# UC Dimension

VC $\rightarrow$ Vapnik - chervonenkin

VC- dimension represent the maximum number of data point that can be shattered or perfectly separated by hypothesic class
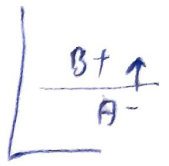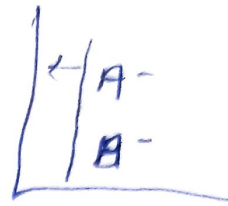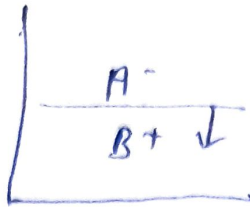
For two dat point = n = 2,     total 4 label $(2^2)$

" three "    " (n=3)    total 8 (label) $(2^3)$

So, $(2^n)$

Let talle two class of posetiue and negatiue $n=2$



A+
→B+

↓ DESIGN Boundary

A-
B+ ↓

←|A-
|B-

B+ ↑
A-

Four lablal

$n=3$



A+
B+
C+ →

A-
B+ ↓
C+    ____, total 8 posibility

VC concept apply in varicour ml, like, Svm, neural netwoork decision tree etc

# # <u>Bagging and boosting</u> @

Bagging and boosting are both ensemble learning technique that aim to improue the performance of individual base classifiers by combining predictions.
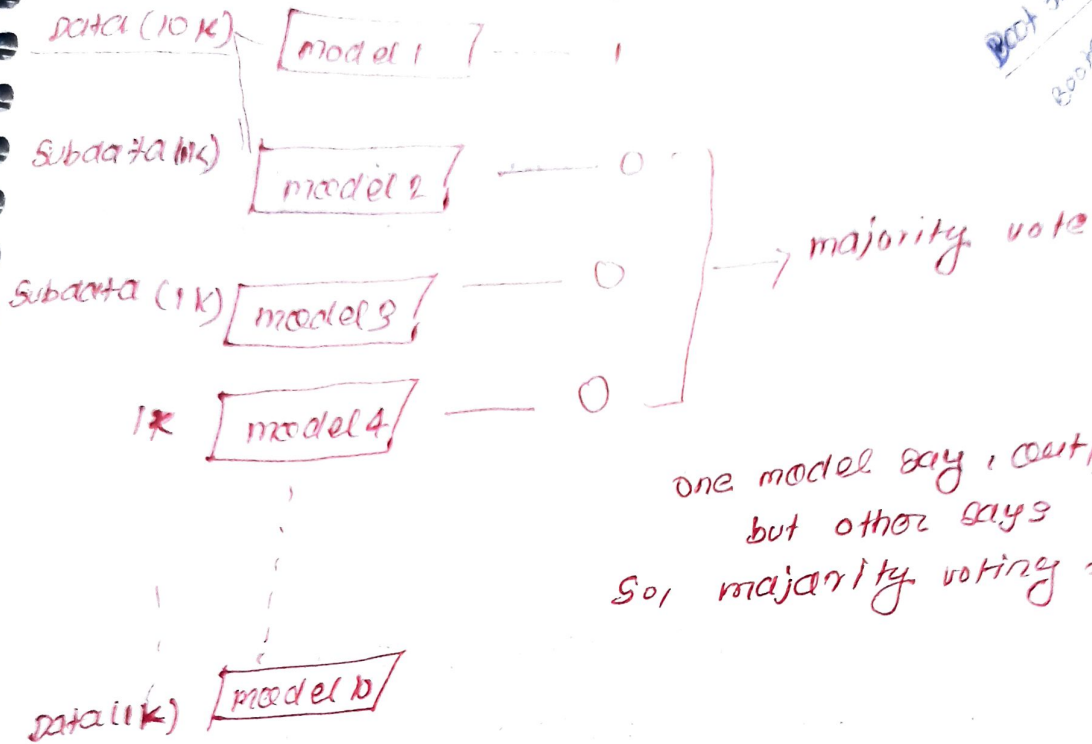
# Bagging

Boot Strapping          Aggregation
                        output on test data (x)

Data (10K) — [model 1] — 1

Subdata (1k) — [model 2] — 0

Subdata (1k) [model 3] — 0   } → majority vote

1K [model 4] — 0

Data (1k) [model b]

One model say, output should 1
but other says 0
So, majority voting = 0   So, output = 0

we give our data and give to each model individualy for training.

Bagging involves training multiple base classifiers independently on different subset of the training data. Each base classifier is trained on a random selected subset of the original training data.

The final prediction is then made by aggregating the predictions of all base classifiers typically through majority voting.

Ex, Random Forest
     Bagged Decision tree

# Boosting

Boosting is an ensemble learning technique that combines multiple weak learner or base model to create a strong predictive model.

It is a sequential learning process where each subsequent model in the ensemble is trained to correct mistake made by previous model.

Et,
- AdaBoost , Gradient Boosting

# Ensemble learning is a machine learning technique that combines the prediction of multiple individual model to make final prediction
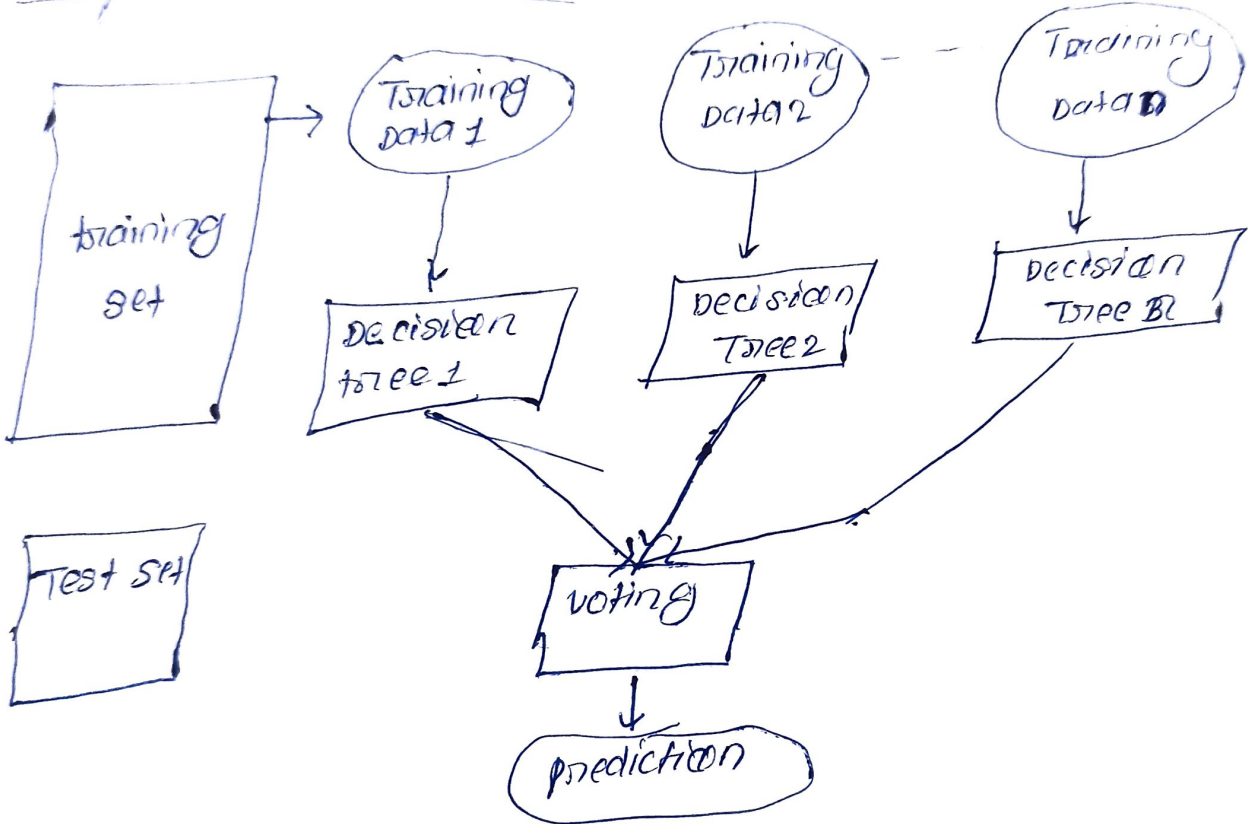
Ex Bagging , Boosting , Random Forest

Random Forest = Bagging + Decision tree .

in Bagging we use any algo-- if we use Decision tree in bagging then it called random forest.

The majority voting classifier is a simple ensemble learning technique that combine the predictions of multiple individual classifiers or model by majority voting to make the final prediction.

# Random Forest



$$cormax = \left[ function \; ka \; differentian = 0 \right]$$

(0) well-posed learning problem consist of
→ input data, output data, problem statement.

2) primary goal of empirical risk minimization

   Ans= To minimize risk associated with m model

3) what is the role of inductive bias in empirical risk minimization

   Ans - To bias the model towards certain type of data

4) PAC → probably Approximately correct learning

5) inductive bias helps models generalize well to unseen data

   ⇒ True

6) what does the "approximately correct" aspect in PAC learning

   Ans The model is approximately accurate in its
   Generalization to unseen data.

7) VC dimension in PAC refer to ?

   Ans The capacity of hypothesis space

8) why is data preprocessing an essential step
   ⇒ it helps improve the performance and reliability of