i) Explain the mathematical formulation of linear Regression & logistic regression. dicuss how cost functions differ b/w them & the role of regularization.

A) **Linear regression**

linear regression predicts a continuous output as a linear combination of input features

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_n x_n = w^T x$$

$x = [x_1, x_2, \cdots, x_n]^T$ feature vector

$w = [w_0, w_1, \cdots, w_n]^T$ - model parameters (weights)

$\hat{y}$ = predicted value

**Cost function**

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

$m$ = number of training examples

$\hat{y}_i = w^T x_i$

This measures how far predictions are from actual

**Logistic regression**

logistic regression predicts a probability for a binary outcome (0 or 1):

$$\hat{y} = P(y=1 | x) = \sigma(w^T x)$$

where

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

**Cost function**

$$J(w) = -\frac{1}{m} \sum_{i=1}^{m} [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

this penalizes wrong predictions more heavily when the model is confident but incorrect.

# Role of Regularization

Regularization helps reduce overfitting by adding a penalty for large co-efficients in the cost function

## a) Ridge Regression

Adds a penalty equal to the square of the coefficients

$$J(w) = Loss + x \sum_{j=1}^{n} w_j^2$$

## b) Lasso Regression

Adds a penalty equal to the absolute value of coefficient

$$J(w) = Loss + x \sum_{j=1}^{n} |w_j|$$

⑤ Define the "curse of dimensionality". How does it affect algorithms like K-Nearest Neighbors (KNN)? Suggest two approaches to reduce its impact in practice.

## Ⓐ Curse of Dimensionality

The curse of dimensionality refers to the set of Problems that arise when data has too many features Campared to the number of samples

## How it Affects Algorithm like KNN

K-Nearest Neighbors depends heavily on distance metrics

1) All points appear equally distant

→the different between the nearest & farthest points becomes negligible

→ Distance loses it meaning so KNN can't distinguish blw close & far points efficiently

$$d \to \infty \qquad \frac{Dmax - Dmin}{Dmin} \to 0$$

a) overfitting risk increases

* KNN may fit noise rather than real pattern because each sample is isolated in high-dimensional space.

* Computational cost rises

more dimensions → higher time & memory needed to compute distances for all features.

* Approaches to Reduce its Impact

a) Dimensionality Reduction

→ PCA (principal Component Analysis)

projects data to a lower-dimensional space by combining correlated features

b) Feature Normalization / Regularization

* Helps balance the influence of features & reduce noise

③ A data base contains information about whether students pass or fail based on three attributes

| Study Hours | Attendence | Internal marks | Result |
|---|---|---|---|
| High | Good | High | Pass |
| High | Good | Low | Pass |
| Low | Good | High | Pass |
| Low | poor | Low | Fail |
| High | Poor | Low | Fail |
| Low | good | Low | Fail |

a) compute the Entropy of the target attribute
b) calculate the information gain for each attribute
c) identify the best attribute to spilt at the root node
d) draw the first level Decision Tree based on your calculations.

## Step 1

Pass = 3 (1, 2, 3)
Fail = 3 (4, 5, 6)

$$P(pass) = \frac{3}{6} = 0.5 \,, \quad P(fail) = 0.5$$

$$Entropy\,(Result) = -[0.5 \log_2(0.5) + 0.5 \log_2(0.5)] = 1.0$$

$$Entropy = 1.0$$

## Step 2

a)

| Study Hours | Pass | fail | Total | Entropy |
|---|---|---|---|---|
| High | 2 | 1 | 3 | $-[\frac{2}{3}\log_2\frac{2}{3} + \frac{1}{3}\log_2\frac{1}{3}] = 0.918$ |
| Low | 1 | 2 | 3 | $-[\frac{1}{3}\log_2\frac{1}{3} + \frac{2}{3}\log_2\frac{2}{3}] = 0.918$ |

$$Entropy_{studyhours} = \frac{3}{6}(0.918) + \frac{3}{6}(0.918) = 0.918$$

IG (study Hours) = 1.0 - 0.918 = 0.082

b)

| Attendence | pass | fail | Total | Entropy |
|---|---|---|---|---|
| Good | 3 | 1 | 4 | $-[\frac{3}{4}\log_2\frac{3}{4} + \frac{1}{4}\log_2\frac{1}{4}] = 0.811$ |
| Poor | 0 | 2 | 2 | $-[0\log_2(0) + 1\log_2(1)] = 0$ |

$$Entropy_{Attendence} = \frac{4}{6}(0.811) + \frac{2}{6}(0) = 0.541$$

IG attendence = 1.0 - 0.541 = 0.459

c) attribute

| Internal marks | Pass | fail | Total | Entropy |
|---|---|---|---|---|
| High | 2 | 0 | 2 | $-[1\log_2(1) + 0] = 0$ |
| Low | 1 | 3 | 4 | $-[\frac{1}{4}\log_2\frac{1}{4} + \frac{3}{4}\log_2\frac{3}{4}] = 0.811$ |

$$Entropy_{Internal\,Marks} = \frac{2}{6}(0) + \frac{4}{6}(0.811) = 0.541$$

IG (Internal Marks) = 1.0 - 0.541 = 0.456

**4) Discuss the max-margin principle of support vector Machines. Explain the concepts of hard margin, soft margin & the role of kernel functions in handling non-linear data**

### 1) support vector machine

Its main goal is to find the optimal separating hyperplane that best the margin b/w two classes

#### Hard margin SVM

* Assumes linearly separable data
* Maximizes distance between Support vectors & hyperplance

#### Soft Margin Svm

* Allows Some misclassification using penalty terme
* Bances margin maximization with error
  object $= min \frac{1}{2} ||w||^2 + C \varepsilon c_i$

#### Kennel tricks

maps data into higher-dimensional Space to find a linear baundary

#### Common Kennels

RBF :- Handles circular/non linear boundaries
Polynomial : captures curved relationships
sigmoid : similar to neural networks

### 5) campare bagging, Boosting, Random forest, & strucking in terms of methodology, base leavenergy & bias variance trade off. provide one read world use case for each

### 1) Bagging

#### Methodology :-

Multiple models are trained on different bootstrap samples of the training data

#### Base Learners

Typically high-variance, low-bias models like Decision Trees

## Bias-Variance Trade-off

Reduces Variance without increasing bias much - helps avoid overfitting

## Real-world use case

Credit Risk predicition - combining multiple decision trees to stabilize predictions about loan default risk.

## Boosting :-

### Methodology

Models are built sequentially, where each model is to correct the errors of the previous one.

### Base Learners

Usally weak fearns, like shallow decision trees

### Bias-variance Trade-off

Reduces bias but may increase variance if over fitted

### Real-world use case

Spam Email Detection - boosting refines predictions to distinguish spam from legitimate mail accurately.

## Random forest :-

### Methodology :

An improved version of bagging where each tree is trained on a bootstrap sample & uses a random subset of features at each spilt - increasing model diversity

### Base learners

Decision Tress

### Bias-variance Trade-off

Greatly reduces variance while mataing low bias - more stable that bagging alone

### Real world use case

Medical Diagnosis - predicting diseases using patient features with high accuracy & robustness.

# Stacking

## Methodology

Combines different of models & uses a meta - model to leasn how to best combine their outputs

## Base Leasns

## Heterogeneous

## Bias-Variance Trade-off

Tries to balance both bias & variance by leveraging strengths of diverse methods

## Real-world use case

House price prediction stacking regression models improves accusancy over any single model

6) explain the working principle of gaussian Naive Bayes and Gaussian processes. How would model calibration improve the reliability of predicted probabilities?

## Gaussian Naive Bayes

* It is uses Bayes rule to predict the class
* It assumes all features are independent
* Each feature follows a normal curve.
* It picks the class with the highest probability

## Gaussian processes

* used for regression
* It assumes data points come from a Gaussian pattern.
* It uses a kernel to find how similar points are.
* Gives both a prediction & how certain it is

## Model calibration

* Makes the model's probilities more accounts
* Example: if it says 80% change ,it really happens about 8 out of 10 times
* It makes prediction more trustwostly.

7) Differentiate between accuracy, precision, recall, f1-core, Roc and AUC. why are these metrics important for imbalance datasets? Also, briefly discuss the importance of model explainability in modern AI systems.

## Accuracy

* Tells how many predictions are correct overall.
* Formula: (correct predictions / total predictions)
* problem: Can be misleading if one class has many more samples

## precision

* out of all predicted positives, how many are actually Positions
* focus : Avoiding false positives
* example: Good for spam detection

## Recall:

* out of all actual postives, how many are correctly found
* focus: Avoiding false negatives
* Example: Good for disease detection

## F1-Score

* The balance b/w precision & recall
* Formula: $2 \times ($ precision $\times$ Recall$) / ($ precision $+$ Recall$)$
* Good when data is imbalance

## Roc (Receiver operating charateristic)

* shows how to model performs at different thresholds
* plots true positive Rate vs False positive Rate,

## AUC (Area under curve)

* Measures the overall performance at different theSold
* Higher AUC = better model.

## why important for imbalanced data

* Accuracy may look high even if the model ignores the smaller dataclass
* precision, Recall, $F_1$, & AUC give a better picture of real performance.