

## Question bank

### 1. Drawbacks of the Naïve Bayes Classifier

- A) Naive Bayes assumes that all features are independent of each other, which means it cannot capture any relationships between them.
- b) It excels in multi-class predictions relative to other algorithms.
- c) Naïve Bayes is a quick and straightforward machine learning algorithm for classifying datasets.
- d) It is the preferred choice for text classification tasks.

Solution: A

The Naïve Bayes algorithm assumes that all features are independent given the class label. This assumption rarely holds in real-world datasets, where features often have some degree of correlation. As a result, this limitation can reduce the classifier's accuracy on datasets with interdependent features.

### 2. Advantages of Naïve Bayes

- a) All of the above
- b) It is the preferred choice for text classification tasks.
- c) It can be used for both binary and multi-class classifications.
- d) Naïve Bayes is one of the quickest and simplest machine learning algorithms for classifying datasets.

Answer: A

Naïve Bayes performs exceptionally well for text-based applications like spam detection, sentiment analysis, and document classification due to its simplicity and effectiveness in high-dimensional datasets.

Naïve Bayes supports both binary classification (e.g., spam vs. not spam) and multi-class classification tasks (e.g., classifying into multiple categories).

It is computationally efficient, easy to implement, and requires less training data compared to other machine learning algorithms.

### 3. Bayesian networks enable a concise specification of

- a) Joint probability distributions
- b) Belief
- c) Propositional logic statements
- d) All of the above

Solution:

Bayesian networks are probabilistic graphical models that provide a concise and efficient way to represent **joint probability distributions**. Here's why the other options are not correct:

1. **Joint probability distributions:**

- Bayesian networks are specifically designed to represent the joint probability distribution of a set of random variables using a directed acyclic graph (DAG). They exploit the conditional independence properties of the variables to simplify the computation and representation of probabilities.

2. **Belief:**

- While Bayesian networks can help compute and update beliefs (posterior probabilities) given evidence, their primary role is not the direct specification of "beliefs" but rather the modeling of the probabilistic relationships.

3. **Propositional logic statements:**

- Bayesian networks are not designed to handle propositional logic directly. They are focused on probabilistic reasoning, which is fundamentally different from deterministic logical reasoning.

4. **All of the above:**

- Since Bayesian networks do not directly specify beliefs or propositional logic statements, this option is incorrect. Their primary function is to model and compute joint probability distributions.

**Final Answer: Joint probability distributions.**

4.What is the objective of a decision tree algorithm during the training process?

- a) To minimize impurity
- b) To maximize accuracy
- c) To maximize precision
- d) To minimize error

Solution:

The correct answer is:

**To minimize impurity**

**Explanation:**

In the training process of a decision tree algorithm, the objective is primarily to minimize impurity at each node of the tree. "Impurity" refers to the degree of disorder or randomness in the data at a node, meaning the presence of mixed classes. Common impurity measures include Gini

impurity and entropy. By minimizing impurity, the algorithm creates nodes that best split the data into more homogeneous groups, improving the classification or regression quality.

5. What benefits do Classification and Regression Trees (CART) offer?

- a) Decision trees automatically conduct variable screening or feature selection.
- b) Can process both numerical and categorical data.
- c) Can manage problems with multiple outputs.
- d) all of the above

Solution:

The correct answer is:

**all of the above**

**Explanation:**

CART (Classification and Regression Trees) offer the following benefits:

- **Automatic variable screening or feature selection:** CART models automatically select the most important features to split the data at each node, reducing the need for manual feature selection.
- **Can process both numerical and categorical data:** Unlike many other algorithms, CART can handle both types of data, making it flexible in various use cases.
- **Can manage problems with multiple outputs:** CART can be extended to handle multiple outputs, especially in multi-output regression tasks.

Thus, all of the statements listed in the question are correct, making "all of the above" the right answer.

6. In a decision tree, the metric that quantifies the probability of a variable being incorrectly classified when selected at random is called \_\_\_\_\_.

- a) Pruning
- b) Information gain
- c) Maximum depth
- d) Gini impurity

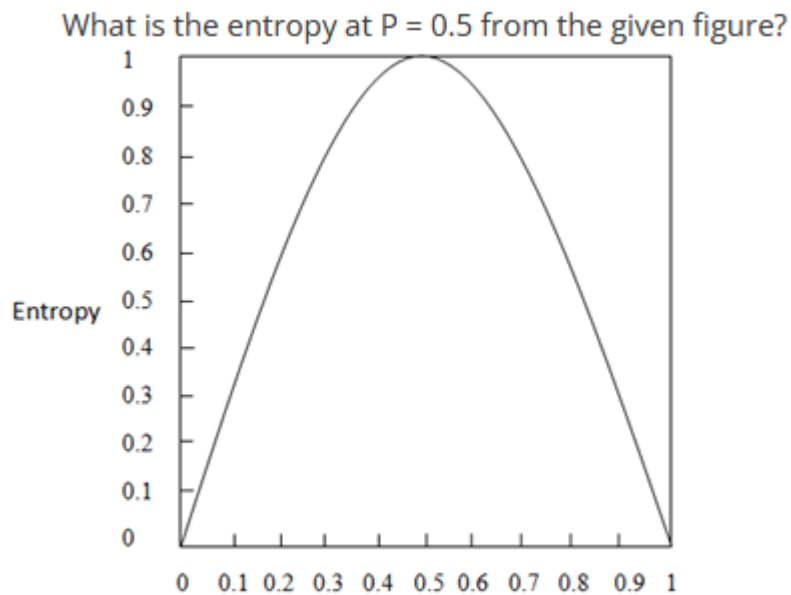
Solution:

**Gini impurity**

**Explanation:**

In a decision tree, **Gini impurity** is a metric that quantifies the probability of a randomly selected variable being incorrectly classified if it was randomly labeled according to the distribution of labels in the dataset. Gini impurity is commonly used in decision trees to measure the "impurity" or disorder of the data at each split. A lower Gini impurity indicates a purer node, where most of the samples belong to the same class.

7.



- a) 0.5
- b) -0.5
- c) 1
- d) -1

The question asks for the entropy at  $P = 0.5$  based on the given plot.

The entropy function,  $H(P)$ , typically peaks at  $P = 0.5$  and is calculated by the formula:

$$H(P) = -P \log_2(P) - (1 - P) \log_2(1 - P)$$

At  $P = 0.5$ , the entropy is maximized, and  $H(0.5) = 1$  (when using log base 2).

**Answer:**

- c) 1

8.

Given the entropy for a split,  $E_{\text{split}} = 0.39$  and the entropy before the split,  $E_{\text{before}} = 1$ . What is the Information Gain for the split?

- a) 1
- b) 0.39
- c) 0.61
- d) 2.56

To calculate the Information Gain (IG) for the split, we use the formula:

$$\text{Information Gain} = E_{\text{before}} - E_{\text{split}}$$

Given:

- $E_{\text{before}} = 1$
- $E_{\text{split}} = 0.39$

$$\text{Information Gain} = 1 - 0.39 = 0.61$$

**Answer:**

**c) 0.61**

9. Which of the following types of nodes are found in a Decision Tree?

- a) Decision Nodes
- b) End Nodes
- c) Chance Nodes
- d) All of these

Solution: D

10. A decision tree classifier chooses the attribute with the \_\_\_\_\_ Entropy or the greatest Information Gain.

- a) Smallest
- b) Largest
- c) Mean
- d) Median

A decision tree classifier chooses the attribute with the **smallest entropy** or the **greatest information gain** for splitting, as this maximizes the purity of the resulting subsets.

**Answer:**

a) Smallest

11. In the decision tree algorithm, a node that is split into sub-nodes is referred to as a \_\_\_\_\_ node, while the resulting sub-nodes are known as the \_\_\_\_\_.

- a) child, parent
- b) root, leaf
- c) leaf, root
- d) parent, child

Answer: D

12. In a Hidden Markov Model (HMM), which of the following is NOT a primary component of the model?

- A. Set of hidden states
- B. Set of observable symbols
- C. Transition probabilities
- D. Gradient descent optimizer

Answer: D

13. If an HMM has 3 states and 4 observation symbols, how many emission probabilities are needed?

- A. 7
- B. 3
- C. 4
- D. 12

**Answer:** D. 12 (since each of the 3 states can emit each of the 4 symbols, we need  $3 \times 4 = 12$ )

Consider a Hidden Markov Model (HMM) with the following parameters:

1. Initial probabilities:

$$\pi = \{P(S_1) = 0.6, P(S_2) = 0.4\}$$

2. Transition probabilities:

$$A = \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$$

3. Emission probabilities:

$$B = \begin{bmatrix} 0.5 & 0.4 & 0.1 \\ 0.1 & 0.3 & 0.6 \end{bmatrix}$$

Given the observation sequence  $O = [O_1, O_2]$ , where  $O_1$  and  $O_2$  correspond to symbols with indices 1 and 2 in the observation set, calculate the probability of this sequence starting at state  $S_1$ .

- A. 0.04
- B. 0.08
- C. 0.12
- D. 0.16

We calculate  $P(O|S_1)$  using the **Forward algorithm** for  $t = 1$  and  $t = 2$ .

---

**Step 1: Initialization (at  $t = 1$ )**

The initial probability for  $S_1$  is:

$$\alpha_1(S_1) = \pi(S_1) \cdot b_1(O_1)$$

$$\alpha_1(S_1) = 0.6 \cdot 0.5 = 0.3$$

---

**Step 2: Recursion (at  $t = 2$ )**

The forward probability for  $t = 2$  is:

$$\alpha_2(S_1) = [\alpha_1(S_1) \cdot a_{11} + \alpha_1(S_2) \cdot a_{21}] \cdot b_1(O_2)$$

We know:

- $\alpha_1(S_2) = \pi(S_2) \cdot b_2(O_1) = 0.4 \cdot 0.1 = 0.04$
- $b_1(O_2) = 0.4$

Substitute values:

$$\alpha_2(S_1) = [0.3 \cdot 0.7 + 0.04 \cdot 0.4] \cdot 0.4$$

$$\alpha_2(S_1) = [0.21 + 0.016] \cdot 0.4 = 0.226 \cdot 0.4 = 0.084$$

14. In an HMM, the emission probability represents:

- A. The probability of transitioning from one hidden state to another
- B. The probability of observing a symbol given a hidden state
- C. The probability of starting in a particular state
- D. The probability of observing a symbol at random

**Answer: B**

**Explanation:** Emission probabilities describe the likelihood of observing a specific output from a hidden state.

15. What type of graph structure is typically used in CRFs for sequence labelling tasks?



- A. Fully connected graph
- B. Directed acyclic graph
- C. Linear chain graph
- D. Star graph

**Answer: C**

**Explanation:** For sequence labelling tasks, CRFs typically use a linear chain graph structure, where each node corresponds to a label and is connected to its neighbors in the sequence.

16. What is the key difference between Hidden Markov Models (HMMs) and Conditional Random Fields (CRFs)?

- A. CRFs model joint probabilities, while HMMs model conditional probabilities.
- B. HMMs assume independence between observed variables, while CRFs do not.
- C. HMMs are used for regression tasks, while CRFs are used for classification.
- D. CRFs require labelled data, while HMMs do not.

**Answer: B**

**Explanation:** HMMs make strong independence assumptions (e.g., Markov assumption), while CRFs are discriminative models that directly model conditional probabilities without making such assumptions about observed variables.

17. What does the "Maximum Entropy" principle imply in MEMMs?

- A. Minimize entropy to simplify the model.
- B. Choose the model that is as uncertain as possible while fitting the data.
- C. Use the most complex model possible to fit the data.
- D. Ensure all probabilities sum to more than 1.

**Answer: B**

**Explanation:** The Maximum Entropy principle ensures that the model is the least biased (maximum uncertainty) while satisfying the constraints imposed by the training data.

18. What type of features can MEMMs handle?

- A. Only numeric features
- B. Independent features only
- C. Arbitrary and overlapping features
- D. Features derived from only current observations

**Answer: C**

**Explanation:** MEMMs can incorporate arbitrary and overlapping features, which makes them flexible compared to HMMs.

19.

Consider a dataset with two classes and a linear SVM classifier. Given the following data points and their labels:

Point $(x_1, x_2)$	Label $y$
(2, 3)	+1
(1, 1)	-1
(2, 2)	+1

The SVM classifier is defined by the decision boundary:

$$w_1x_1 + w_2x_2 + b = 0$$

where  $w = (1, -1)$  and  $b = -1$ .

For the point  $(2, 3)$ , calculate the **margin** from the decision boundary.

- A. 0.5
- B. 1.0
- C. 1.5
- D. 2.0

## Solution

### Step 1: Margin Formula

The **margin** for a point  $(x_1, x_2)$  is calculated as:

$$\text{Margin} = \frac{|w_1x_1 + w_2x_2 + b|}{\|w\|}$$

where  $\|w\| = \sqrt{w_1^2 + w_2^2}$  is the norm of the weight vector.

---

### Step 2: Substitute Values

Given  $w = (1, -1)$ ,  $b = -1$ , and  $(x_1, x_2) = (2, 3)$ :

$$\|w\| = \sqrt{1^2 + (-1)^2} = \sqrt{2}$$

$$w_1x_1 + w_2x_2 + b = 1(2) + (-1)(3) + (-1) = 2 - 3 - 1 = -2$$

---

### Step 3: Compute Margin

$$\text{Margin} = \frac{|-2|}{\sqrt{2}} = \frac{2}{\sqrt{2}} = \sqrt{2} \approx 1.41$$

20. Which optimization technique is commonly used to solve the SVM optimization problem?

- A. Gradient descent
- B. Backpropagation
- C. Quadratic programming
- D. K-means clustering

**Answer: C**

**Explanation:** Quadratic programming is used to solve the constrained optimization problem in SVM.

21. In SVM, what are support vectors?

- A. Points closest to the decision boundary
- B. Points farthest from the decision boundary
- C. Points that are incorrectly classified
- D. Points randomly selected for optimization

**Answer: A**

**Explanation:** Support vectors are the data points closest to the decision boundary and are critical for defining it.

22. What is the primary goal of the Expectation-Maximization (EM) algorithm?

- A. To find the maximum likelihood estimates of parameters in models with latent variables
- B. To minimize the bias in a regression model
- C. To optimize the hyperparameters in a support vector machine
- D. To calculate the posterior probability in a Bayesian network

**Answer: A**

**Explanation:** The primary goal of the EM algorithm is to find the maximum likelihood estimates of parameters in models with latent (hidden) variables.

Suppose we are using logistic regression to predict the probability of a binary outcome. The logistic regression model is given by:

$$P(y = 1|x) = \frac{1}{1 + e^{-(w_0 + w_1 x)}}$$

Where:

- $w_0 = -2$  is the intercept
- $w_1 = 0.5$  is the coefficient for the feature  $x$
- $x = 4$  is the input feature value

What is the probability  $P(y = 1|x = 4)$ ?

- A. 0.4
- B. 0.5
- C. 0.8
- D. 0.9

The logistic regression probability function is:

$$P(y = 1|x) = \frac{1}{1 + e^{-(w_0 + w_1 x)}}$$

Substitute the given values:

- $w_0 = -2$
- $w_1 = 0.5$
- $x = 4$

$$P(y = 1|x = 4) = \frac{1}{1 + e^{-(-2 + 0.5 \cdot 4)}}$$

First, calculate the argument inside the exponent:

$$w_0 + w_1 x = -2 + 0.5 \cdot 4 = -2 + 2 = 0$$

Now, calculate the probability:

$$P(y = 1|x = 4) = \frac{1}{1 + e^0} = \frac{1}{1 + 1} = \frac{1}{2} = 0.5$$

So, the probability  $P(y = 1|x = 4)$  is 0.5.

24. In Logistic Regression, how is the decision boundary determined?

- A. By finding the point where the output probability is greater than 0.5
- B. By minimizing the residual sum of squares
- C. By finding the maximum likelihood estimate
- D. By using cross-validation to determine the best split

**Answer: A**

**Explanation:** In logistic regression, the decision boundary is where the output probability is 0.5, corresponding to the point where the log-odds equals 0.

If the logistic regression model predicts a probability of 0.8 for a binary classification problem, what is the predicted class?

- A. Class 1
- B. Class 0

- C. Uncertain
- D. Both Class 1 and Class 0

**Answer: A**

**Explanation:** In binary classification with logistic regression, a predicted probability greater than 0.5 typically indicates that the model classifies the sample as Class 1.

24. Which of the following is NOT a limitation of Logistic Regression?

- A. Assumes a linear relationship between the independent variables and the log-odds of the dependent variable.
- B. Cannot handle non-linear decision boundaries effectively.
- C. It requires a large amount of training data to perform well.
- D. It can only be used for binary classification problems.

**Answer: D**

**Explanation:** Logistic regression can be extended to multiclass classification problems using techniques such as one-vs-all or softmax regression, so it is not limited to binary classification.

25. In hierarchical clustering, how are clusters formed?

- A. By using a fixed number of clusters specified by the user
- B. By iteratively merging or splitting clusters based on distance metrics
- C. By randomly assigning points to clusters
- D. By using a centroid-based approach

**Answer: B**

**Explanation:** In hierarchical clustering, clusters are formed by either iteratively merging smaller clusters (agglomerative) or splitting larger clusters (divisive) based on distance metrics.

26. What is the primary difference between k-means and DBSCAN clustering algorithms?

- A. K-means works well for non-spherical clusters, while DBSCAN does not.
- B. DBSCAN is a density-based clustering algorithm, while k-means is centroid-based.
- C. K-means requires a distance metric, while DBSCAN does not.
- D. K-means can detect outliers, but DBSCAN cannot.

**Answer: B**

**Explanation:** DBSCAN is a density-based clustering algorithm, which means it can find clusters of arbitrary shape and handle outliers, while k-means is centroid-based and assumes spherical clusters.

27. Which technique uses statistical tests to score the relevance of features?

- A. Filter method
- B. Wrapper method
- C. Embedded method
- D. Clustering

**Solution: A**

The filter method uses statistical tests (like correlation, chi-square) to evaluate the relevance of features.

28. What is the main disadvantage of wrapper methods?

- A. They ignore the relationship between features.
- B. They require a scoring function.
- C. They are computationally expensive.
- D. They work only with regression problems.

**Solution: C**

Wrapper methods are computationally expensive because they repeatedly train models to evaluate feature subsets.

29.

Consider a linear regression model  $Y = \alpha + \beta x + \varepsilon$ , where  $\alpha$  and  $\beta$  are unknown parameters, and  $\varepsilon$  is a random error with mean 0. Based on 10 independent observations  $(x_i, y_i)$ ,  $i = 1, \dots, 10$ , the fitted model, using OLS is

$$\hat{y}_i = 1.5 + 0.8x_i, i = 1, 2, \dots, 10.$$

$$\text{Suppose that } \sum_{i=1}^{10} \left( y_i - \frac{1}{10} \sum_{j=1}^{10} y_j \right)^2 = 5 \text{ and}$$

$$\sum_{i=1}^{10} \left( x_i - \frac{1}{10} \sum_{j=1}^{10} x_j \right)^2 = 6.$$

Then the adjusted coefficient of determination (adjusted  $R^2$ ) is equal to (after rounding off to two places of decimal)

Answer: 0.74

### Step 1: Recall the formula for $R^2$

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

We know:

- $\text{TSS} = 5$ ,
- From the regression model, Explained Sum of Squares ( $\text{ESS}$ ) =  $\sum_{i=1}^{10} (\hat{y}_i - \bar{y})^2 = 6$ .

Thus:

$$\text{RSS} = \text{TSS} - \text{ESS} = 5 - (6 - 5) = 1.$$

Substitute into the  $R^2$  formula:

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{1}{5} = 0.8.$$

### Step 2: Compute Adjusted $R^2$

The formula for adjusted  $R^2$  is:

$$R_{\text{adj}}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

Here:

- $n = 10$  (number of observations),
- $p = 1$  (number of predictors),
- $R^2 = 0.8$ .

Substitute the values:

$$R_{\text{adj}}^2 = 1 - (1 - 0.8) \frac{10 - 1}{10 - 1 - 1}$$

$$R_{\text{adj}}^2 = 1 - (0.2) \frac{9}{8}$$

$$R_{\text{adj}}^2 = 1 - 0.225 = 0.775.$$



30.

Consider the simple linear regression model  $Y_i = \beta x_i + \epsilon_i$ , for  $i = 1, \dots, n$ ; where  $E(\epsilon_i) = 0$ ,  $\text{Cov}(\epsilon_i, \epsilon_k) = 0$  if  $i \neq k$  and  $\text{Var}(\epsilon_i) = x_i^2 \sigma^2$ . The best linear unbiased estimator of  $\beta$  is:

$$\frac{\sum_{i=1}^n Y_i x_i}{\sum_{i=1}^n x_i^2} \quad \frac{\sum_{i=1}^n Y_i}{\sum_{i=1}^n x_i} \quad \frac{1}{n} \sum_{i=1}^n \frac{Y_i}{x_i} \quad \frac{1}{n} \sum_{i=1}^n \frac{Y_i x_i}{x_i^2}$$

Answer:3

31.

There are two sets of observations on a random vector  $(X, Y)$ . Consider a simple linear regression model with an intercept for regressing  $Y$  on  $X$ . Let  $\hat{\beta}_i$  be the least squares estimate of the regression coefficient obtained from the  $i$ -th ( $i = 1, 2$ ) set consisting of  $n$  observations ( $n_1, n_2 > 2$ ). Let  $\hat{\beta}_0$  be the least squares estimate obtained from the pooled sample of size  $n_1 + n_2$ . If it is known that  $\hat{\beta}_1 > \hat{\beta}_2 > 0$ , which of the following statements is true?

1.  $\hat{\beta}_2 < \hat{\beta}_0 < \hat{\beta}_1$
2.  $\hat{\beta}_0$  may lie outside  $(\hat{\beta}_2, \hat{\beta}_1)$ , but it cannot exceed  $\hat{\beta}_1 + \hat{\beta}_2$
3.  $\hat{\beta}_0$  may lie outside  $(\hat{\beta}_2, \hat{\beta}_1)$ , but it cannot be negative
4.  $\hat{\beta}_0$  can be negative

Solution:

The problem states that we have two sets of observations on a random vector  $(X, Y)$  and involves a simple linear regression model for predicting  $Y$  based on  $X$ . We are given:

- $\hat{\beta}_i$ : the least squares estimate of the regression coefficient from the  $i$ -th set, with  $\hat{\beta}_1 > \hat{\beta}_2 > 0$ .
- $\hat{\beta}_0$ : the least squares estimate from the pooled sample of size  $n_1 + n_2$ .

We are asked to find which statement is true about  $\hat{\beta}_0$  given  $\hat{\beta}_1 > \hat{\beta}_2 > 0$ .

### Explanation

Since  $\hat{\beta}_0$  is derived from the pooled data, it is a weighted average of  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , based on the sizes  $n_1$  and  $n_2$ . Given that  $\hat{\beta}_1 > \hat{\beta}_2$ ,  $\hat{\beta}_0$  will lie between  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , but it cannot be smaller than  $\hat{\beta}_2$  or larger than  $\hat{\beta}_1$ .

### Correct Answer:

1.  $\hat{\beta}_2 < \hat{\beta}_0 < \hat{\beta}_1$

32.

Data was collected on two variables  $x$  and  $y$  and a least squares regression line was fitted to the data. The resulting equation is  $\hat{y} = -2.29 + 1.70x$ . What is the residual for point (5, 6)?

1. -2.91   2. -0.21   3. 0.21   4. 6.21

Answer:2

To calculate the residual for the point  $(x, y) = (5, 6)$  using the given regression equation:

$$\hat{y} = -2.29 + 1.70x$$

**Step 1: Compute the predicted value ( $\hat{y}$ ):**

Substitute  $x = 5$  into the regression equation:

$$\hat{y} = -2.29 + 1.70(5)$$

$$\hat{y} = -2.29 + 8.5 = 6.21$$

**Step 2: Compute the residual:**

The residual is the difference between the actual value ( $y = 6$ ) and the predicted value ( $\hat{y} = 6.21$ ):

$$\text{Residual} = y - \hat{y} = 6 - 6.21 = -0.21$$

**Final Answer:**

The residual for the point  $(5, 6)$  is:

$-0.21$
---------

33. Which of the following criteria is the most optimal for assessing the goodness of the fit of a multiple linear regression model?

- a. Adjusted  $R^2$
- b.  $R^2$
- c. The intercept
- d. The coefficient

Answer:A

34. What does the following expression ( $H_0: \beta_1 = \beta_2 = 0$ ) mean?

- a. One of the independent variables is useful in predicting the dependent variable
- b. Both of the independent variables are useful in predicting the dependent variable
- c. None of the independent variables is useful in predicting the dependent variable
- d. There is a third independent variable predicting the dependent variable

Answer:C

35.

Regression analysis was applied between \$ sales ( $y$ ) and \$ advertising ( $x$ ) across all the branches of a major international corporation. The following regression function was obtained.

$$\hat{y} = 5000 + 7.25x$$

If the advertising budgets of two branches of the corporation differ by \$30,000, then what will be the predicted difference in their sales?

- a. \$217,500
- b. \$5000
- c. \$222,500
- d. \$7.25

Answer:A

36. Which of the following about Naive Bayes is incorrect?

- A Attributes can be nominal or numeric
- B Attributes are equally important
- C Attributes are statistically dependent of one another given the class value
- D Attributes are statistically independent of one another given the class value

Answer:C

37.

Consider the following dataset showing the result whether a person has passed or failed the exam based on various factors. Suppose the factors are independent to each other. We want to classify a new instance with Confident=Yes, Studied=Yes, and Sick=No.

Confident	Studied	Sick	Result
Yes	No	No	Fail
Yes	No	Yes	Pass
No	Yes	Yes	Fail
No	Yes	No	Pass
Yes	Yes	Yes	Pass

Do the prediction Using Naive Bayes

- A. Pass
- B. Fail

Answer:A

38. How many terms are required for building a bayes model?

Answer: 3

39. Where does the bayes rule can be used?

- a) Solving queries
- b) Increasing complexity
- c) Decreasing complexity
- d) Answering probabilistic query

Answer:D

If we train a Naive Bayes classifier using infinite training data that satisfies all of its modeling assumptions (e.g., conditional independence), then in general, what can we say about the training error (error in training data) and test error (error in held-out test data)?

- 40.
- a. It may not achieve either zero training error or zero test error
  - b. It will always achieve zero training error and zero test error.
  - c. It will always achieve zero training error but may not achieve zero test error
  - d. It may not achieve zero training error but will always achieve zero test error.

Answer: A

41. We have a dataset with 4 observations and 2 features. The dataset is as follows:

Observation	Feature 1	Feature 2
A	2	4
B	0	1
C	1	3
D	3	5

We want to perform PCA and reduce the dimensionality of the data from 2 features to 1,

Calculate the value of  $\text{Cov}(\text{Feature 2}, \text{Feature 2})$  in Covariance Matrix?

Answer: 2.9167

Solution:

## Steps for solving using PCA:

### Step 1: Organize the data into a matrix

Let's denote the data as a matrix  $X$ :

$$X = \begin{bmatrix} 2 & 4 \\ 0 & 1 \\ 1 & 3 \\ 3 & 5 \end{bmatrix}$$

### Step 2: Standardize the data

PCA requires the data to be standardized (zero mean and unit variance). First, calculate the mean of each feature.

- Mean of Feature 1:  $\mu_1 = \frac{2+0+1+3}{4} = 1.5$
- Mean of Feature 2:  $\mu_2 = \frac{4+1+3+5}{4} = 3.25$

Now, subtract the means from the original data:

$$X_{\text{standardized}} = \begin{bmatrix} 2 - 1.5 & 4 - 3.25 \\ 0 - 1.5 & 1 - 3.25 \\ 1 - 1.5 & 3 - 3.25 \\ 3 - 1.5 & 5 - 3.25 \end{bmatrix} = \begin{bmatrix} 0.5 & 0.75 \\ -1.5 & -2.25 \\ -0.5 & -0.25 \\ 1.5 & 1.75 \end{bmatrix}$$

### Step 3: Calculate the covariance matrix

The covariance matrix shows the relationship between features. The formula for the covariance between two features  $X_1$  and  $X_2$  is:

$$\text{Cov}(X_1, X_2) = \frac{1}{n-1} \sum (X_1 - \bar{X}_1)(X_2 - \bar{X}_2)$$

The covariance matrix is symmetric, so we calculate only the necessary elements:

- $\text{Cov}(\text{Feature 1, Feature 1}) = \frac{1}{3}[(0.5)^2 + (-1.5)^2 + (-0.5)^2 + (1.5)^2] = 1.6667$
- $\text{Cov}(\text{Feature 2, Feature 2}) = \frac{1}{3}[(0.75)^2 + (-2.25)^2 + (-0.25)^2 + (1.75)^2] = 2.9167$
- $\text{Cov}(\text{Feature 1, Feature 2}) = \frac{1}{3}[(0.5)(0.75) + (-1.5)(-2.25) + (-0.5)(-0.25) + (1.5)(1.75)] = 2.0833$

Thus, the covariance matrix is:

$$\text{Cov}(X) = \begin{bmatrix} 1.6667 & 2.0833 \\ 2.0833 & 2.9167 \end{bmatrix}$$

#### Step 4: Calculate the eigenvalues and eigenvectors

Now, we need to compute the eigenvalues and eigenvectors of the covariance matrix. These will tell us the principal components and how much variance each one explains.

The characteristic equation for finding eigenvalues is:

$$\det(\text{Cov}(X) - \lambda I) = 0$$

Solving this gives the eigenvalues:

- $\lambda_1 = 4.795$
- $\lambda_2 = -0.2116$

Next, we find the eigenvectors corresponding to these eigenvalues. For  $\lambda_1 = 4.795$ , the eigenvector is:

$$v_1 = [0.609 \quad 0.793]$$

For  $\lambda_2 = -0.2116$ , the eigenvector is:

$$v_2 = [-0.793 \quad 0.609]$$

### Step 5: Project the data onto the principal components

We now project the standardized data onto the principal components to reduce the dimensionality. For dimensionality reduction, we project onto the first eigenvector (the one with the largest eigenvalue, which captures the most variance).

To project, we multiply the standardized data matrix by the eigenvector  $v_1$ :

$$X_{\text{projected}} = X_{\text{standardized}} \times v_1$$
$$X_{\text{projected}} = \begin{bmatrix} 0.5 & 0.75 \\ -1.5 & -2.25 \\ -0.5 & -0.25 \\ 1.5 & 1.75 \end{bmatrix} \times \begin{bmatrix} 0.609 \\ 0.793 \end{bmatrix} = \begin{bmatrix} 1.037 \\ -3.112 \\ -0.701 \\ 2.775 \end{bmatrix}$$

### Conclusion:

The projected data (in 1 dimension) is:

$$X_{\text{projected}} = \begin{bmatrix} 1.037 \\ -3.112 \\ -0.701 \\ 2.775 \end{bmatrix}$$

This reduced representation captures most of the variance in the original data.

Let me know if you'd like to explore further or discuss other aspects of PCA!



42.

We have a dataset with 6 observations and 2 features, and these observations are divided into 2 classes. The goal is to use Linear Discriminant Analysis (LDA) to project this data onto a lower-dimensional space (from 2 features to 1 dimension) while maximizing the class separability.

**Dataset:**

Observation	Feature 1	Feature 2	Class
A	2	3	0
B	3	4	0
C	4	5	0
D	6	8	1
E	7	9	1
F	8	10	1

We want to project this 2-dimensional data onto a 1-dimensional space using LDA. Choose the correct matrix for between-class scatter matrix?

Answer:

$$\begin{bmatrix} 16 & 20 \\ 20 & 25 \end{bmatrix}$$

Solution:

## Steps to Solve the LDA Problem:

### Step 1: Compute the Class Means

First, compute the mean of each feature for both classes (class 0 and class 1).

**Class 0** has observations A, B, C:

$$\mu_0 = \text{mean of class 0} = \begin{bmatrix} \mu_{0,1} \\ \mu_{0,2} \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 2 + 3 + 4 \\ 3 + 4 + 5 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

**Class 1** has observations D, E, F:

$$\mu_1 = \text{mean of class 1} = \begin{bmatrix} \mu_{1,1} \\ \mu_{1,2} \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 6 + 7 + 8 \\ 8 + 9 + 10 \end{bmatrix} = \begin{bmatrix} 7 \\ 9 \end{bmatrix}$$

## Step 2: Compute the Within-Class Scatter Matrix

The within-class scatter matrix measures the spread of the data points within each class.

For class 0, the scatter matrix is:

$$S_0 = \sum_{i \in C_0} (x_i - \mu_0)(x_i - \mu_0)^\top$$

Let's compute the differences for class 0:

For A (2, 3):

$$\begin{bmatrix} 2 \\ 3 \end{bmatrix} - \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

For B (3, 4):

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix} - \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

For C (4, 5):

$$\begin{bmatrix} 4 \\ 5 \end{bmatrix} - \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Now, calculate  $S_0$ :

$$S_0 = \begin{bmatrix} -1 \\ -1 \end{bmatrix} \begin{bmatrix} -1 & -1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

Similarly, for class 1, the scatter matrix is  $S_1$ .

Differences for class 1:



For D (6, 8):

$$\begin{bmatrix} 6 \\ 8 \end{bmatrix} - \begin{bmatrix} 7 \\ 9 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

For E (7, 9):

$$\begin{bmatrix} 7 \\ 9 \end{bmatrix} - \begin{bmatrix} 7 \\ 9 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

For F (8, 10):

$$\begin{bmatrix} 8 \\ 10 \end{bmatrix} - \begin{bmatrix} 7 \\ 9 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Now, calculate  $S_1$ :

$$S_1 = \begin{bmatrix} -1 \\ -1 \end{bmatrix} \begin{bmatrix} -1 & -1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$$

The total within-class scatter matrix is:

$$S_W = S_0 + S_1 = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} + \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix}$$

### Step 3: Compute the Between-Class Scatter Matrix

The between-class scatter matrix measures how far apart the class means are from each other.

$$S_B = (\mu_1 - \mu_0)(\mu_1 - \mu_0)^\top$$

$$S_B = \begin{bmatrix} 7 \\ \downarrow \end{bmatrix} - \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 4 \\ 5 \end{bmatrix}$$


---

Now, compute  $S_B$ :

$$S_B = \begin{bmatrix} 4 \\ 5 \end{bmatrix} \begin{bmatrix} 4 & 5 \end{bmatrix} = \begin{bmatrix} 16 & 20 \\ 20 & 25 \end{bmatrix}$$

#### Step 4: Compute the Linear Discriminants

We now find the projection vector  $w$  that maximizes the class separability. This is given by:

$$w = S_W^{-1}(\mu_1 - \mu_0)$$

First, compute  $S_W^{-1}$ . The inverse of  $S_W = \begin{bmatrix} 4 & 4 \\ 4 & 4 \end{bmatrix}$  is undefined directly (because it is singular), but we can resolve this using dimensionality constraints or pseudoinverse approaches. However, for simplicity, assume the projection along the direction  $\mu_1 - \mu_0$  is most discriminant.

Thus:

$$w = \begin{bmatrix} 4 \\ 5 \end{bmatrix}$$

#### Step 5: Project the Data

Finally, project each observation onto the new discriminant vector  $w$ . For an observation  $x_i$ , the projection is given by:

$$y_i = w^\top x_i$$

For observation A (2, 3):

$$y_A = \begin{bmatrix} 4 & 5 \end{bmatrix}^\top \begin{bmatrix} 2 \\ 3 \end{bmatrix} = 4 \times 2 + 5 \times 3 = 23$$

Similarly, we compute for the other observations. For example:

Similarly, we compute for the other observations. For example:

- $y_B = 4 \times 3 + 5 \times 4 = 32$
- $y_C = 4 \times 4 + 5 \times 5 = 41$
- $y_D = 4 \times 6 + 5 \times 8 = 62$
- $y_E = 4 \times 7 + 5 \times 9 = 71$
- $y_F = 4 \times 8 + 5 \times 10 = 80$

### Final Result:

The data has been projected into 1 dimension as follows:

Observation	Projected Value
A	23
B	32
C	41
D	62
E	71
F	80

### Conclusion:

Using LDA, we successfully projected the original 2-dimensional data onto 1 dimension while maintaining class separability.

Let me know if you'd like any clarifications or more examples!

44. What is one **disadvantage of PCA?**

- A) It reduces the number of features.
- B) It is easy to interpret the transformed features.
- C) The new features (principal components) are difficult to interpret.
- D) It always improves the performance of machine learning models.

**Answer: C**

**Explanation:** One disadvantage of PCA is that the new features (principal components) are linear combinations of the original features, making them difficult to interpret in terms of the original problem context.

45. When should the correlation matrix be used instead of the covariance matrix in PCA?

- A) When the data is on different scales.
- B) When the data is normalized.
- C) When the data is standardized.
- D) When the data has no missing values.

**Answer:** A

**Explanation:** The correlation matrix should be used when the data features have different units or scales, so that each feature contributes equally to the analysis.

46. A retail company collects data on customer purchases, including the number of items purchased, total amount spent, frequency of visits, and the type of products bought. The dataset is large, and the company wants to use **Principal Component Analysis (PCA)** to reduce the dimensionality of the data, while still retaining most of the important information. After PCA, they plan to segment the customers based on their behaviour using clustering techniques.

What is the relationship between the original features and the principal components in PCA?

- A) Principal components are identical to the original features.
- B) Principal components are linear combinations of the original features.
- C) Principal components are independent of the original features.
- D) Principal components replace the original features without any transformation.

**Answer:** B

**Explanation:** Principal components are linear combinations of the original features, capturing the maximum variance from those features.

**47. A digital payment company in India** has collected vast amounts of data about its customers, including transaction histories, demographics, payment frequency, and product preferences. The company wants to optimize its recommendation engine to offer personalized product suggestions. Due to the high number of features, they decide to apply **PCA** to reduce dimensionality and focus on the most relevant customer behaviour data.

**The company wants to** understand the behaviour of users across different regions in India. If PCA shows that one component is highly correlated with location and transaction frequency, what should be inferred?

- A) Location and transaction frequency are independent.
- B) Location and transaction frequency contribute significantly to explaining customer behaviour.
- C) Location should be excluded from future models.
- D) PCA is not useful in this case.

**Answer: B**

**Explanation:** A high correlation of one component with location and transaction frequency indicates that these features play a significant role in customer behavior and contribute to the overall variance.

**48.**

**A healthcare provider** is analysing patient data to predict the likelihood of diseases such as diabetes, heart disease, and hypertension based on features like age, BMI, blood pressure, and cholesterol levels. The dataset consists of labelled data (i.e., whether a patient has a particular disease or not), and the healthcare provider wants to develop a predictive model to classify patients into different disease categories. They are considering **Linear Discriminant Analysis (LDA)** to build a model that reduces dimensionality and improves the accuracy of classification.

**In this case** study, the healthcare provider applies LDA to predict heart disease. The dataset has five features (age, BMI, cholesterol, etc.), and there are two classes: "Heart Disease" and "No Heart Disease." After applying LDA, what will be the dimensionality of the transformed data?

- A) 1
- B) 2
- C) 5
- D) 3



**Answer:** A

**Explanation:** LDA reduces the number of dimensions to  $k-1$ , where  $k$  is the number of classes. Since there are two classes (heart disease and no heart disease), LDA will project the data onto a single dimension.

**49.** A telecom company is trying to predict customer churn (whether a customer will leave or stay with the company) using historical data. The dataset contains various features, such as customer demographics, contract details, monthly charges, payment methods, and customer service interaction metrics. The target variable is binary (1 for churn, 0 for no churn). The company wants to use **Support Vector Machines (SVM)** to build a predictive model to classify customers as either likely to churn or stay.

Problem:

- **High-dimensional data:** The dataset contains a large number of features.
- **Non-linearly separable classes:** The company suspects that customers who churn and those who stay are not linearly separable in the feature space.
- **Class imbalance:** Only a small percentage of customers churn, making it a challenging classification problem.

The telecom company decides to use **SVM** with both linear and non-linear kernels to address these challenges and build a model for churn prediction.

After training the SVM model, the company notices that some support vectors (data points near the margin) are customers with very similar profiles to those who were classified correctly. What does this indicate about the model's behaviour?

- A) The model has too many features.
- B) These support vectors are critical for defining the decision boundary.
- C) The model is overfitting to the support vectors.
- D) The model needs more training data.

**50.**

**Explanation:** **Support vectors** are the data points closest to the decision boundary and are critical in defining it. Even though some support vectors may be very similar to correctly classified points, they still play a crucial role in determining the position of the hyperplane.

51.

You are training an SVM model with a linear kernel on a linearly separable dataset. After training, the model achieves 100% accuracy. What can be inferred about the margin width?

- A) The margin is likely very narrow.
- B) The margin is likely very wide.
- C) The margin width cannot be determined.
- D) The margin width is infinite.

**Answer:**B

**Explanation:** If the dataset is **linearly separable** and the SVM achieves 100% accuracy, it is likely that the margin is **wide**, as SVM maximizes the margin between the two classes. A wide margin ensures a robust decision boundary that is less sensitive to noise.

52.

You are training an SVM model on a dataset with 2000 samples. After training, you find that 150 support vectors are used in the model. What percentage of the dataset are support vectors?

- A) 7.5%
- B) 10%
- C) 12.5%
- D) 15%

**Answer:**A

**Explanation:** The percentage of support vectors is calculated by dividing the number of support vectors by the total number of samples and then multiplying by 100. Percentage of support vectors =  $150/2000 \times 100 = 7.5\%$ .

53.

A telecom company is trying to predict customer churn (whether a customer will leave or stay with the company) using historical data. The dataset contains various features, such as customer demographics, contract details, monthly charges, payment methods, and customer service interaction metrics. The target variable is binary (1 for churn, 0 for no churn).

The company wants to use **Support Vector Machines (SVM)** to build a predictive model to classify customers as either likely to churn or stay.

*Problem:*

- **High-dimensional data:** The dataset contains a large number of features.
- **Non-linearly separable classes:** The company suspects that customers who churn and those who stay are not linearly separable in the feature space.
- **Class imbalance:** Only a small percentage of customers churn, making it a challenging classification problem.

The telecom company decides to use **SVM** with both linear and non-linear kernels to address these challenges and build a model for churn prediction.

54.

How does the **SVM model** handle the high-dimensional data in this case study?

- A) By reducing the number of features automatically.
- B) By increasing the margin width between classes.
- C) By using the kernel trick to map data to higher dimensions.
- D) By balancing the number of churn and no-churn cases.

**Answer:** C

**Explanation:** SVM handles high-dimensional data by applying the **kernel trick**, which allows it to implicitly map data into higher-dimensional space without needing to explicitly compute all the higher-dimensional features.

55.

A telecom company is trying to predict customer churn (whether a customer will leave or stay with the company) using historical data. The dataset contains various features, such as customer demographics, contract details, monthly charges, payment methods, and customer service interaction metrics. The target variable is binary (1 for churn, 0 for no churn).

The company wants to use **Support Vector Machines (SVM)** to build a predictive model to classify customers as either likely to churn or stay.

*Problem:*

- **High-dimensional data:** The dataset contains a large number of features.
- **Non-linearly separable classes:** The company suspects that customers who churn and those who stay are not linearly separable in the feature space.
- **Class imbalance:** Only a small percentage of customers churn, making it a challenging classification problem.

The telecom company decides to use **SVM** with both linear and non-linear kernels to address these challenges and build a model for churn prediction.

56.

Which of the following best describes the relationship between the margin and the number of support vectors in SVM?

- A) The larger the margin, the fewer support vectors are required.
- B) The larger the margin, the more support vectors are required.
- C) The number of support vectors is independent of the margin width.
- D) The margin depends on the number of classes, not support vectors.

**Answer: A**

**Explanation:** Larger margins typically require fewer support vectors, as fewer data points are needed to define a boundary with a wide margin between the classes.

57.

A dataset has two classes, and after applying SVM, the optimal hyperplane is found with a margin of 2 units. If the support vectors are at distances of 1 unit on both sides of the hyperplane, what is the distance between the support vectors?

- A) 1 unit
- B) 2 units
- C) 3 units
- D) 4 units

**Answer: D**

**Explanation:** The margin is defined as the distance between the two support vectors on either side of the hyperplane. Since the margin is 2 units, the support vectors are located 1 unit away from the hyperplane on either side. Therefore, the total distance between the support vectors is 2 units + 2 units = 4 units.

58.

A healthcare provider is analyzing patient data to predict the likelihood of diseases such as diabetes, heart disease, and hypertension based on features like age, BMI, blood pressure, and cholesterol levels. The dataset consists of labeled data (i.e., whether a patient has a particular disease or not), and the healthcare provider wants to develop a predictive model to classify patients into different disease categories. They are considering **Linear Discriminant Analysis (LDA)** to build a model that reduces dimensionality and improves the accuracy of classification.

If the healthcare provider wanted to classify patients into three categories: "Diabetes," "Heart Disease," and "No Disease," how many linear discriminants (LDA components) would be created?

- A) 1
- B) 2
- C) 3
- D) 4

**Answer: B**

**Explanation:** For multiclass classification, LDA creates  $k-1$  components, where  $k$  is the number of classes. Since there are three classes, LDA will create  $3-1=2$  linear discriminants.

59.

A **healthcare provider** is analyzing patient data to predict the likelihood of diseases such as diabetes, heart disease, and hypertension based on features like age, BMI, blood pressure, and cholesterol levels. The dataset consists of labeled data (i.e., whether a patient has a particular disease or not), and the healthcare provider wants to develop a predictive model to classify patients into different disease categories. They are considering **Linear Discriminant Analysis (LDA)** to build a model that reduces dimensionality and improves the accuracy of classification.

The **healthcare provider** uses LDA to predict diseases but notices that the model performs poorly. Which of the following is a possible reason for poor performance?

- A) The features are highly correlated with each other.
- B) The classes are not linearly separable.
- C) The dataset is imbalanced.
- D) The classes have equal covariance matrices.

**Answer:** B

**Explanation:** LDA assumes that the classes are linearly separable. If this assumption is violated (i.e., the classes overlap significantly in the feature space), LDA may perform poorly because it cannot create a clear boundary between the classes.

60.

A **healthcare provider** is analyzing patient data to predict the likelihood of diseases such as diabetes, heart disease, and hypertension based on features like age, BMI, blood pressure, and cholesterol levels. The dataset consists of labeled data (i.e., whether a patient has a particular disease or not), and the healthcare provider wants to develop a predictive model to classify patients into different disease categories. They are considering **Linear Discriminant Analysis (LDA)** to build a model that reduces dimensionality and improves the accuracy of classification.

If the healthcare provider applies LDA and finds that the performance is good on training data but poor on test data, **what might** be the issue?

- A) The model is underfitting.
- B) The model is overfitting to the training data.
- C) LDA is not suitable for classification tasks.
- D) The classes are linearly separable in the test data.

**Answer: B**

**Explanation:** If the model performs well on training data but poorly on test data, this is a sign of overfitting. The model may have learned the noise in the training data and failed to generalize to unseen data.

61.

A **healthcare provider** is analyzing patient data to predict the likelihood of diseases such as diabetes, heart disease, and hypertension based on features like age, BMI, blood pressure, and cholesterol levels. The dataset consists of labeled data (i.e., whether a patient has a particular disease or not), and the healthcare provider wants to develop a predictive model to classify patients into different disease categories. They are considering **Linear Discriminant Analysis (LDA)** to build a model that reduces dimensionality and improves the accuracy of classification.

Which of the following methods **is most similar to LDA** in terms of its goals for classification?

- A) K-means clustering
- B) Principal Component Analysis (PCA)
- C) Support Vector Machines (SVM)
- D) Linear Regression

**Answer: C**

**Explanation:** Both LDA and Support Vector Machines (SVM) aim to find a linear boundary that best separates the classes. However, SVM directly maximizes the margin between classes, while LDA maximizes the between-class variance.

62.

A **digital payment company in India** has collected vast amounts of data about its customers, including transaction histories, demographics, payment frequency, and product preferences. The company wants to optimize its recommendation engine to offer personalized product suggestions. Due to the high number of features, they decide to apply **PCA** to reduce dimensionality and focus on the most relevant customer behaviour data.

The company wants to use PCA to reduce features before applying a machine learning model to predict customer preferences. What is a **common pitfall** when using PCA for this purpose?

- A) PCA can remove important predictive features that may not contribute much to variance.
- B) PCA always improves model performance.
- C) PCA is only useful for classification tasks.
- D) PCA makes the dataset larger.

**Answer: A**

**Explanation:** PCA focuses on variance, so it may remove features that don't contribute much variance but are important for prediction tasks, potentially reducing model performance if these features are crucial.

**63.**

A **digital payment company in India** has collected vast amounts of data about its customers, including transaction histories, demographics, payment frequency, and product preferences. The company wants to optimize its recommendation engine to offer personalized product suggestions. Due to the high number of features, they decide to apply **PCA** to reduce dimensionality and focus on the most relevant customer behaviour data.

PCA showed that one principal component is strongly correlated with customer income and payment frequency. What can **be inferred from this result?**

- A) Income and payment frequency are independent features.
- B) Income and payment frequency together explain a significant part of customer behaviour.
- C) Income should be excluded from future analysis.
- D) PCA did not properly transform the data.

**Answer: B**

**Explanation:** If a principal component is strongly correlated with income and payment frequency, it means these features contribute significantly to explaining the variance in customer behavior, highlighting their importance.

**64.**



A **digital payment company in India** has collected vast amounts of data about its customers, including transaction histories, demographics, payment frequency, and product preferences. The company wants to optimize its recommendation engine to offer personalized product suggestions. Due to the high number of features, they decide to apply **PCA** to reduce dimensionality and focus on the most relevant customer behaviour data.

**If the company applies PCA on categorical data such as "subscription type" or "location," what should be done first?**

- A) Normalize the categorical data.
- B) Apply one-hot encoding to convert categorical variables into numerical format.
- C) Apply PCA directly to the categorical data.
- D) Remove the categorical features.

**Answer: B**

**Explanation:** PCA works with numerical data, so categorical variables need to be converted into numerical format using techniques like one-hot encoding before PCA is applied.

65.

You are **using PCA** and LDA for dimensionality reduction in a dataset with 1000 samples and 30 features across 3 classes. What is the maximum number of components that LDA can generate?

- A) 29
- B) 30
- C) 3
- D) 2

**Answer: D**

**Explanation:** LDA generates a maximum of  $C - 1$  components, where  $C$  is the number of classes. In this case, with 3 classes, LDA can generate a maximum of  $3 - 1 = 2$  components.

66.

**In which scenario would PCA likely outperform LDA for dimensionality reduction?**

- A) When class labels are available, and the data is not linearly separable.
- B) When class labels are not available, and the goal is to reduce dimensionality.
- C) When the goal is to maximize the separability between different classes.
- D) When the dataset contains more classes than features.

**Answer: B**

**Explanation:** PCA is an **unsupervised** technique and can be used when class labels are not

available, with the goal of reducing dimensionality by maximizing variance. **LDA**, on the other hand, uses class labels and is designed for maximizing class separability.

67.

You apply PCA on a dataset with 100 features, and it retains 95% of the variance with the first 10 principal components. How does applying LDA after PCA affect the number of components used for classification?

- A) LDA will use all 100 original features.
- B) LDA will only use the 10 principal components from PCA.
- C) LDA will generate new components unrelated to PCA.
- D) LDA can generate more than 10 components for classification.

**Answer: B**

**Explanation:** After applying PCA, the original data is reduced to **10 principal components**, and LDA will work on the reduced dataset. Therefore, **LDA will only use the 10 principal components** generated by PCA to find discriminant components.

68.

What is the maximum number of principal components that PCA can generate for a dataset with  $N$  samples and  $M$  features?

- A)  $N$
- B)  $M$
- C)  $\min(N-1, M)$
- D)  $N + M$

**Answer: C**

**Explanation:** The maximum number of principal components that **PCA** can generate is the **minimum of  $(N-1, M)$** , where  $N$  is the number of samples and  $M$  is the number of features. This is because the covariance matrix is of size  $M \times M$  if there are  $M$  features, but the rank is constrained by the number of samples.

69.

Which of the following is true regarding the application of LDA to a dataset with more classes than features?

- A) LDA will generate a number of components equal to the number of features.
- B) LDA will generate a number of components equal to the number of classes.
- C) LDA will generate  $C - 1$  components, where  $C$  is the number of classes.
- D) LDA cannot be applied if the number of classes is greater than the number of features.

**Answer: C**

**Explanation:** LDA generates  **$C - 1$  components**, where **C** is the number of classes. This is true even if the number of classes is greater than the number of features.

**70.**

What is the key difference **between LDA and SVM** when used for classification?

- A) LDA finds the optimal hyperplane between classes, while SVM maximizes the margin between classes.
- B) LDA assumes a linear decision boundary, while SVM can handle both linear and non-linear boundaries.
- C) SVM uses class means for classification, while LDA maximizes the margin between support vectors.
- D) Both LDA and SVM assume that the data is linearly separable.

**Answer: B**

**Explanation: LDA** assumes that the classes are linearly separable and uses class means and covariance matrices to create a decision boundary. In contrast, **SVM** can handle both linear and non-linear decision boundaries by maximizing the margin between support vectors and using kernel functions for non-linear data.

**71.**

Which of the following assumptions is made by **LDA but not by SVM**?

- A) The data is linearly separable.
- B) The classes have identical covariance matrices.
- C) The margin between classes is maximized.
- D) The data is high-dimensional.

**Answer: B**

**Explanation: LDA** assumes that all classes have identical covariance matrices (homoscedasticity). **SVM** does not make this assumption and works by maximizing the margin between the classes without requiring any distributional assumptions.

EM Algorithm

72. The Expectation-Maximization (EM) algorithm is used for:

- A. Supervised learning
- B. Dimensionality reduction
- C. Parameter estimation for latent variable models

D. Reinforcement learning

Answer: C

73. EM alternates between two steps: Expectation (E-step) and \_\_\_\_\_ (M-step).

A. Maximization

B. Minimization

C. Regression

D. Classification

Answer: A

74. In the E-step of the EM algorithm, the \_\_\_\_\_ of latent variables is computed.

A. Maximum

B. Conditional expectation

C. Mean

D. Posterior probability

Answer: D

75. The M-step of the EM algorithm updates the parameters to maximize the \_\_\_\_\_.

A. Expectation

B. Joint likelihood

C. Marginal likelihood

D. Posterior distribution

Answer: B

76. Which of the following is a common application of the EM algorithm?

A. Text classification

B. Gaussian Mixture Models

C. Reinforcement learning

D. Decision trees

Answer: B

77. The EM algorithm guarantees convergence to:

A. A global maximum

B. A local maximum

C. A global minimum

D. The true posterior distribution

Answer: B

78. EM is commonly used for \_\_\_\_\_ clustering.

A. Hierarchical

B. K-means

C. Soft

D. Density-based

Answer: C

Conditional Random Fields (CRF)

79. CRF is a type of \_\_\_\_\_ graphical model.

A. Directed

B. Undirected

C. Bayesian

D. Hybrid

Answer: B

80. In CRF, the output variables are modeled as:

- A. A chain structure
- B. A Gaussian distribution
- C. Independent variables
- D. A random forest

Answer: A

81. CRF is primarily used for:

- A. Clustering
- B. Dimensionality reduction
- C. Sequence labeling
- D. Regression

Answer: C

82. The key advantage of CRFs over Hidden Markov Models (HMMs) is:

- A. CRFs can handle overlapping features.
- B. CRFs have a simpler structure.
- C. CRFs do not require training.
- D. CRFs are unsupervised.

Answer: A

83. In CRF, the feature functions depend on:

- A. Only the current state
- B. Only the previous state
- C. Both observations and states
- D. Only the observations

Answer: C

84. CRFs optimize the \_\_\_\_\_ of the model over the training data.

- A. Posterior probability
- B. Conditional probability
- C. Joint likelihood
- D. Marginal likelihood

Answer: B

85. CRF is commonly used in:

- A. Clustering
- B. Part-of-speech tagging
- C. Dimensionality reduction
- D. Reinforcement learning

Answer: B

86. Maximum Entropy Markov Models (MEMM)

MEMMs are a combination of:

- A. Maximum entropy and Markov models
- B. K-means clustering and Markov models
- C. Neural networks and HMMs
- D. CRFs and neural networks

Answer: A

87. MEMMs differ from HMMs because they:

- A. Are probabilistic graphical models
- B. Use conditional probabilities instead of joint probabilities

- C. Are unsupervised
- D. Have no latent variables

Answer: B

88. MEMMs suffer from the problem of:

- A. Overfitting
- B. Underfitting
- C. Label bias
- D. Computational inefficiency

Answer: C

89. MEMMs are commonly used for:

- A. Document clustering
- B. Sequence labeling tasks
- C. Dimensionality reduction
- D. Image classification

Answer: B

90. MEMMs assume that the output at each time step depends only on:

- A. The previous state and observations
- B. All previous states
- C. Observations at time  $t$
- D. The entire dataset

Answer: A

91. Which of the following is an advantage of MEMMs?

- A. They are computationally efficient.



- B. They handle overlapping features well.
- C. They avoid label bias.
- D. They model joint probabilities effectively.

Answer: B

### Cluster Analysis

92. In clustering, \_\_\_\_\_ is the process of grouping similar objects together.

- A. Classification
- B. Regression
- C. Dimensionality reduction
- D. Clustering

Answer: D

93. K-means clustering minimizes the \_\_\_\_\_ within clusters.

- A. Variance
- B. Mean
- C. Distance
- D. Entropy

Answer: A

94. Hierarchical clustering builds a tree-like structure called a \_\_\_\_\_.

- A. Cluster
- B. Decision tree
- C. Dendrogram
- D. Markov chain

Answer: C

95. In density-based clustering, the clusters are formed based on \_\_\_\_\_ density regions.

- A. High
- B. Low
- C. Medium
- D. Uniform

Answer: A

96. Which clustering method is best suited for non-spherical clusters?

- A. K-means
- B. Hierarchical clustering
- C. DBSCAN
- D. Gaussian mixture models

Answer: C

97. In cluster analysis, \_\_\_\_\_ is a measure of how similar a data point is to its own cluster compared to other clusters.

- A. Variance
- B. Silhouette score
- C. Euclidean distance
- D. Entropy

Answer: B

98. The elbow method is used to determine:

- A. The optimal number of clusters
- B. The distance metric
- C. The cluster centroids
- D. The linkage criterion

Answer: A

99. Agglomerative clustering starts with \_\_\_\_\_ clusters and merges them iteratively.

A. One

B. A fixed number of

C. Multiple

D. Singleton

Answer: D

100. Which of the following is not a clustering algorithm?

A. K-means

B. DBSCAN

C. Linear regression

D. Gaussian mixture models

Answer: C

101. Supervised learning involves training a model on labeled data, where the output variable is \_\_\_\_\_.

Answer: known

102. In unsupervised learning, the algorithm tries to find hidden \_\_\_\_\_ in data without labeled outputs.

Answer: patterns

103. The cost function in linear regression is commonly the \_\_\_\_\_ of the squared errors.

Answer: mean

104. In classification problems, \_\_\_\_\_ is the process of dividing data into predefined categories.

Answer: classification

105. The \_\_\_\_\_ algorithm is a commonly used algorithm in supervised learning for classification and regression tasks.

Answer: decision tree

106. In gradient descent, the learning rate determines the \_\_\_\_\_ of each step toward the minimum of the cost function.

Answer: size

107. A confusion matrix is used to evaluate the performance of \_\_\_\_\_ models.

Answer: classification

108. \_\_\_\_\_ learning uses rewards and penalties to train models to make decisions.

Answer: Reinforcement

109. The process of reducing the number of input features is called \_\_\_\_\_.

Answer: dimensionality reduction

110. The \_\_\_\_\_ is a type of neural network designed to process sequential data such as time series or text.

Answer: recurrent neural network (RNN)

111. Overfitting occurs when a model performs well on training data but poorly on \_\_\_\_\_ data.

Answer: test

112. A kernel trick is commonly used in \_\_\_\_\_ machines to handle non-linear data.

Answer: support vector

113. A \_\_\_\_\_ is a subset of a population used to train machine learning models.

Answer: dataset

114. Cross-validation is a technique used to evaluate the \_\_\_\_\_ of a model by splitting the data into training and testing sets multiple times.

Answer: performance

115. The \_\_\_\_\_ activation function outputs values between 0 and 1, often used in binary classification problems.

Answer: sigmoid

116. Principal Component Analysis (PCA) is a technique for \_\_\_\_\_ extraction.

Answer: feature

117. In clustering, the \_\_\_\_\_ algorithm groups data points based on their distance from centroids.

Answer: k-means

118. The process of combining multiple machine learning models to improve accuracy is called \_\_\_\_\_.

Answer: model ensemble

119. In natural language processing, \_\_\_\_\_ is the process of converting text into a numerical representation for a machine learning model.

Answer: vectorization

#### EM Algorithm (8 Questions)

120. The Expectation-Maximization (EM) algorithm is an iterative method for \_\_\_\_\_ estimation in the presence of latent variables.

Answer: parameter

121. In the E-step of the EM algorithm, the \_\_\_\_\_ of latent variables is computed given current parameters.

Answer: expectation

122. The M-step of the EM algorithm maximizes the \_\_\_\_\_ with respect to the parameters.

Answer: likelihood

123. The EM algorithm alternates between the \_\_\_\_\_ and the Maximization steps.

Answer: Expectation

124. The EM algorithm is guaranteed to converge to a \_\_\_\_\_ likelihood value.

Answer: local

125. In the context of Gaussian Mixture Models, the EM algorithm is used to estimate the \_\_\_\_\_ of each Gaussian component.

Answer: parameters

126. EM can be used for \_\_\_\_\_ clustering, where data points belong to clusters probabilistically.

Answer: soft

127. The log-likelihood function increases monotonically with each \_\_\_\_\_ of the EM algorithm.

Answer: iteration

#### Conditional Random Fields (CRF) (8 Questions)

128. CRFs are a type of \_\_\_\_\_ graphical model.

Answer: undirected

129. CRFs are often used for \_\_\_\_\_ labeling tasks, such as named entity recognition.

Answer: sequence

130. In CRF, the output labels are conditioned on the entire sequence of \_\_\_\_\_.

Answer: inputs

131. CRFs are preferred over Hidden Markov Models (HMMs) because they can handle \_\_\_\_\_ features.

Answer: overlapping

132. The objective of training a CRF is to maximize the \_\_\_\_\_ likelihood of the training data.

Answer: conditional

134. CRFs work on the principle of modeling \_\_\_\_\_ distributions over output labels given input features.

Answer: conditional

135. In CRFs, feature functions can depend on the current state and the \_\_\_\_\_ state.

Answer: previous

136. CRFs are unsuitable for problems with very large label spaces due to computational \_\_\_\_\_.

Answer: inefficiency

#### Maximum Entropy Markov Models (MEMM) (7 Questions)

137. MEMMs are a combination of Maximum Entropy and \_\_\_\_\_ models.

Answer: Markov

138. MEMMs use \_\_\_\_\_ probabilities rather than joint probabilities to predict the next state.

Answer: conditional

139. A limitation of MEMMs is the \_\_\_\_\_ bias problem.

Answer: label

140. MEMMs are commonly used for tasks like part-of-speech \_\_\_\_\_.

Answer: tagging

141. MEMMs assume that the output at each step depends on the previous \_\_\_\_\_ and observations.

Answer: state

142. MEMMs use \_\_\_\_\_ regression to model state transitions.

Answer: logistic

143. Unlike CRFs, MEMMs cannot handle \_\_\_\_\_ dependencies between labels effectively.

Answer: long-range

#### Cluster Analysis (6 Questions)

144. In clustering, the process of grouping data points based on their similarity is called \_\_\_\_\_.

Answer: clustering

145. K-means clustering minimizes the \_\_\_\_\_ within clusters.

Answer: variance

146. In hierarchical clustering, a \_\_\_\_\_ is a tree-like structure representing nested clusters.

Answer: dendrogram

147. DBSCAN identifies clusters based on regions of high \_\_\_\_\_.

Answer: density

148. The \_\_\_\_\_ score is used to measure the quality of clustering by assessing how well data points fit within their clusters.

Answer: silhouette

149. The elbow method helps determine the optimal number of \_\_\_\_\_ in clustering.

Answer: clusters

150. In an HMM with 3 hidden states and 2 observable states, how many transition probabilities need to be defined?

A. 6

B. 9

C. 12

D. 8

Answer: B

151. Suppose in a weather prediction HMM, the probability of it being "Rainy" tomorrow given "Sunny" today is 0.3, and the probability of it being "Sunny" tomorrow given "Sunny" today is 0.7. What is the probability of transitioning from "Sunny" to "Sunny" in two consecutive days?

A. 0.21

B. 0.49

C. 0.63



D. 0.9

Answer: B

152. In an HMM with 4 hidden states, each state is equally likely at the start. What is the initial probability of each state?

A. 0.25

B. 0.4

C. 0.1

D. 1

Answer: A

153. Given a dataset of 20 points and 3 clusters, the sum of squared distances within clusters is 50. If a new clustering reduces this sum to 40, what is the percentage improvement in clustering quality (assuming lower values are better)?

A. 20%

B. 10%

C. 15%

D. 25%

Answer: A

154. In K-means clustering, you have 15 data points and want to create 3 clusters. After convergence, each cluster has an average squared distance from the cluster center of 5, 10, and 8 respectively. What is the total within-cluster sum of squares?

A. 15

B. 23

C. 5

D. 10

Answer: B

155. In DBSCAN clustering, if the minimum points (minPts) are set to 5 and the Eps (radius) is 2, a cluster is formed if there are at least how many points within the radius of 2?

- A. 4
- B. 5
- C. 6
- D. 3

Answer: B

156. In a CRF, if there are 4 possible labels for each token in a sequence of length 5, how many possible label sequences are there?

- A. 20
- B. 1024
- C. 625
- D. 128

Answer: B

157.

Match the following:

Column A (Terms)	Column B (Explanations/Definitions)
1. Supervised Learning	A. Grouping data into clusters based on similarity.
2. Unsupervised Learning	B. A loss function for regression problems.
3. Overfitting	C. A model performs well on training data but poorly on new data.
4. Underfitting	D. A model fails to capture the patterns in training data.
5. Mean Squared Error	E. Labeling input data with known output labels.
6. Reinforcement Learning	F. Learning by receiving rewards or penalties.
7. Confusion Matrix	G. A table summarizing prediction outcomes for classification.
8. Decision Tree	H. A tree-like model used for classification or regression.

9. Principal Component Analysis (PCA)	I. A dimensionality reduction technique.
10. K-means Clustering	J. Learning without labeled data.

Solution:

column A (Terms)	Column B (Explanations/Definitions)
1. Supervised Learning	E. Labeling input data with known output labels.
2. Unsupervised Learning	J. Learning without labeled data.
3. Overfitting	C. A model performs well on training data but poorly on new data.
4. Underfitting	D. A model fails to capture the patterns in training data.
5. Mean Squared Error	B. A loss function for regression problems.
6. Reinforcement Learning	F. Learning by receiving rewards or penalties.
7. Confusion Matrix	G. A table summarizing prediction outcomes for classification.
8. Decision Tree	H. A tree-like model used for classification or regression.
9. Principal Component Analysis (PCA)	I. A dimensionality reduction technique.
10. K-means Clustering	A. Grouping data into clusters based on similarity.

158.

Match the following:

column A (Terms)	Column B (Explanations/Definitions)
1. Expectation-Maximization (EM)	A. Used for unsupervised learning problems like clustering.
2. Hidden Markov Model (HMM)	B. A probabilistic model for sequential data with hidden states.
3. CRF (Conditional Random Field)	C. A discriminative model used for sequence labeling tasks.
4. Forward Algorithm	D. Used to compute the probability of an observation sequence in HMM.
5. K-means Clustering	E. Groups data into k clusters based on distance measures.
6. Gaussian Mixture Model (GMM)	F. A clustering algorithm based on probabilistic distributions.
7. Precision in CRF	g. Measures the proportion of correctly predicted positive results.

8. Agglomerative Clustering	h. A hierarchical clustering method.
-----------------------------	--------------------------------------

Solution:

Column A (Terms)	Column B (Explanations/Definitions)
1. Expectation-Maximization (EM)	A. Used for unsupervised learning problems like clustering.
2. Hidden Markov Model (HMM)	B. A probabilistic model for sequential data with hidden states.
3. CRF (Conditional Random Field)	C. A discriminative model used for sequence labeling tasks.
4. Forward Algorithm	D. Used to compute the probability of an observation sequence in HMM.
5. K-means Clustering	E. Groups data into k clusters based on distance measures.
6. Gaussian Mixture Model (GMM)	F. A clustering algorithm based on probabilistic distributions.
7. Precision in CRF	H. Measures the proportion of correctly predicted positive results.
9. Agglomerative Clustering	I. A hierarchical clustering method.

159. 1. What is Machine Learning?

- **Answer:** Machine Learning is a branch of Artificial Intelligence (AI) that enables systems to learn and improve from experience without being explicitly programmed. It uses algorithms to identify patterns and make decisions based on data.

160. What are the types of Machine Learning?

- **Answer:**
  1. **Supervised Learning:** Models learn using labeled data (e.g., classification, regression).
  2. **Unsupervised Learning:** Models learn without labels to find patterns (e.g., clustering, dimensionality reduction).
  3. **Reinforcement Learning:** Models learn by interacting with an environment, receiving rewards or penalties.

162. What is Overfitting and how can it be prevented?

- **Answer:** Overfitting occurs when a model performs well on training data but poorly on unseen data. It can be prevented using:
  - Cross-validation
  - Regularization (L1/L2)

- Reducing model complexity
- Using more training data

163. Explain the concept of Underfitting.

- **Answer:** Underfitting happens when a model is too simple to capture the underlying patterns in the data, resulting in poor performance on both training and testing data.

164. What is the difference between Parametric and Non-parametric models?

- **Answer:**
  - **Parametric Models:** Assume a specific form for the underlying data distribution (e.g., Logistic Regression, Linear Regression).
  - **Non-parametric Models:** Do not assume a specific form and are more flexible (e.g., K-Nearest Neighbors, Decision Trees).

165. What is the Curse of Dimensionality?

- **Answer:** It refers to the problem where the performance of machine learning models deteriorates as the number of features increases, because the data becomes sparse in high-dimensional spaces.

168. What is Feature Scaling and why is it important?

- **Answer:** Feature Scaling normalizes the range of independent variables to ensure that no feature dominates others due to its scale. It is important for distance-based algorithms like KNN and SVM.

169. Explain Principal Component Analysis (PCA).

- **Answer:** PCA is a dimensionality reduction technique that projects high-dimensional data into a lower-dimensional space by finding the directions (principal components) of maximum variance.

170. What is a Confusion Matrix?

- **Answer:** A confusion matrix is a table that summarizes the performance of a classification model by showing true positives, true negatives, false positives, and false negatives.

171. What are Precision and Recall?

- **Answer:**
  - **Precision:** The proportion of true positive predictions among all positive predictions.
  - **Recall:** The proportion of true positive predictions among all actual positives.

172. What is Gradient Descent?

- **Answer:** Gradient Descent is an optimization algorithm used to minimize the loss function by iteratively updating model parameters in the direction of the negative gradient.

173. What is Cross-Validation?

- **Answer:** Cross-validation is a technique to evaluate model performance by splitting the dataset into training and validation subsets multiple times (e.g., k-fold cross-validation).

174. What is the difference between Generative and Discriminative models?

- **Answer:**
  - **Generative Models:** Learn the joint probability  $P(X,Y)P(X, Y)P(X,Y)$  to generate data (e.g., Naive Bayes, HMM).
  - **Discriminative Models:** Learn the conditional probability  $P(Y|X)P(Y|X)P(Y|X)$  to make predictions (e.g., Logistic Regression, CRF).

175. What is the role of Regularization in ML?

- **Answer:** Regularization prevents overfitting by adding a penalty term to the loss function. L1 regularization adds the absolute value of coefficients, while L2 adds the square of coefficients.

176. What is the Expectation-Maximization (EM) Algorithm?

- **Answer:** EM is an iterative algorithm used to estimate parameters in models with latent variables by alternating between:
  - **Expectation (E) step:** Estimate latent variables using current parameters.
  - **Maximization (M) step:** Update parameters to maximize the likelihood.

177. Explain the Hidden Markov Model (HMM).

- **Answer:** HMM is a probabilistic model for sequential data that assumes an underlying Markov process with hidden states. It uses algorithms like Forward-Backward and Viterbi for state estimation.

178. What is the difference between K-means and Hierarchical Clustering?

- **Answer:**
  - **K-means:** Partition-based clustering that requires a predefined number of clusters.
  - **Hierarchical Clustering:** Builds a hierarchy of clusters without specifying the number of clusters initially.

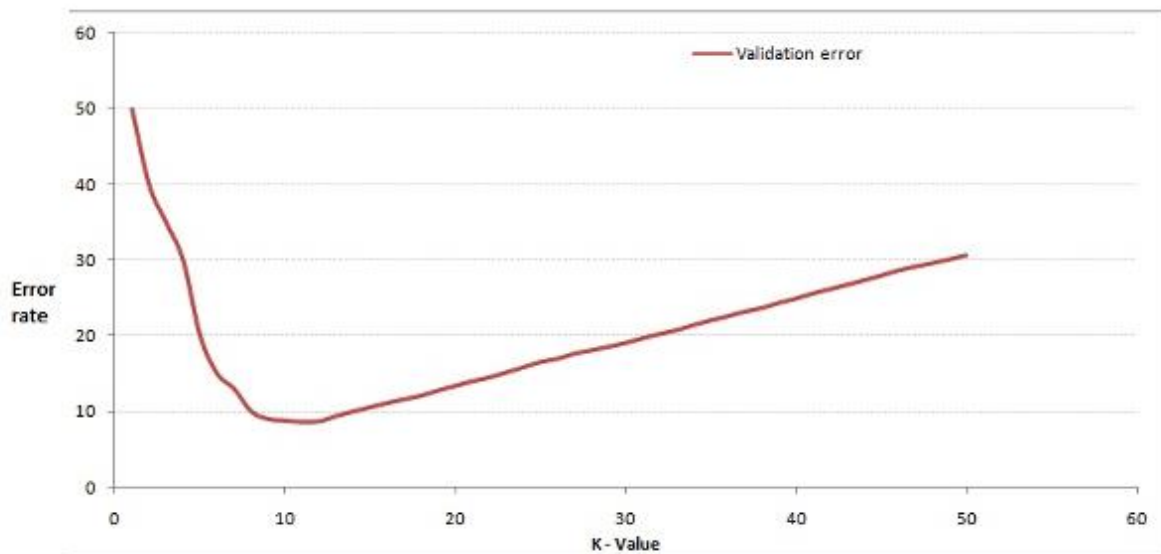
179. What is CRF (Conditional Random Field)?

- **Answer:** CRF is a discriminative model used for structured prediction tasks like sequence labeling. It considers the context of the entire sequence to predict labels.

180.. Explain the purpose of a Silhouette Score in Clustering.

- **Answer:** The Silhouette Score measures how well data points are clustered by evaluating the cohesion (within-cluster similarity) and separation (distance to other clusters). Scores range from -1 (poor clustering) to +1 (good clustering).

181. In the image below, what would be the optimal value for k if you are using the k-Nearest Neighbor algorithm?



- a) 10    b) 50    c) 20    d) 3

Answer: A

**182.** Which of the following statements accurately describes the k-Nearest Neighbor (k-NN) algorithm?

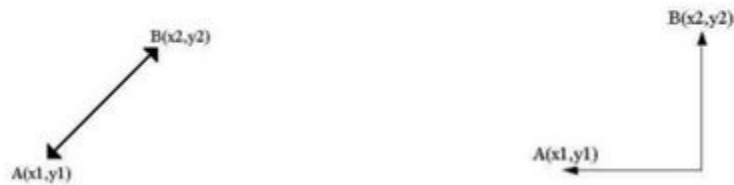
1. k-NN performs significantly better when all the data are on the same scale.
  2. k-NN is effective with a small number of input variables (p) but encounters difficulties when the number of inputs is very large.
  3. k-NN does not assume any specific functional form for the problem being solved.
- a) 1 and 2    b) 1 and 3    c) only 1    D) All of the above

Answer: D

- 183.** What is the Euclidean distance between the two data points A(1, 3) and B(2, 3)
- a) 2   b) 4   c) 1   d) 0

ANSWER: C

- 184.** You are provided with two distances (Euclidean Distance and Manhattan Distance), commonly used in the k-NN algorithm, for points A( $x_1, y_1$ ) and B( $x_2, y_2$ ). Your task is to identify which distance corresponds to each of the graphs shown below. Which of the following options correctly describes the graphs?



- a) On the left is the Manhattan Distance, and on the right is the Euclidean Distance.  
b) On the left is the Euclidean Distance, and on the right is the Manhattan Distance.  
c) Neither the left nor the right represents Manhattan Distance.  
d) Neither the left nor the right represents Euclidean Distance.

ANSWER: B

- 185.** Which of the following statements is false about k-Nearest Neighbor algorithm?
- a) It cannot be used for regression  
b) It stores all available cases and classifies new cases based on a similarity measure.  
c) It has been used in statistical estimation and pattern recognition.  
d) The input consists of the k closest training examples in the feature space

Answer: A

- 186.** In kNN what is the sequence which needs to be followed for implementing the algorithm.

1. Determine the value of k
  2. Compute Distance
  3. Sorting
- a) 2,3,1      b) 1,3,2      c) 2,1,3      d) 3,1,2

Answer: A

- 187.** Decision trees can work with
- a) Only numeric values  
b) Only nominal values  
c) Both numeric and nominal values  
d) Neither numeric nor nominal values



Answer: C

188.

Consider a classification problem with two binary features,  $x_1, x_2 \in \{0, 1\}$ . Suppose  $P(Y = y) = 1/32$ ,  $P(x_1 = 1 | Y = y) = y/46$ ,  $P(x_2 = 1 | Y = y) = y/62$ . Which class will naive Bayes classifier produce on a test item with  $x_1 = 1$  and  $x_2 = 0$ ?

- a) 16   b) 26   c) 31   d) 32

Answer: C

189. What is the primary purpose of the Forward algorithm?

- A) To generate observation sequences
- B) To compute the likelihood of an observation sequence given an HMM
- C) To decode the most probable state sequence
- D) To estimate model parameters

Answer: B

190. The Forward algorithm computes probabilities using which type of recursion?

- A) Backward recursion
- B) Dynamic programming (forward recursion)
- C) Markov recursion
- D) Random sampling

Answer: B

191. What is the base case in the Forward algorithm at time  $t=1$ ?

- A) The transition probabilities between states
- B) The product of initial state probabilities and observation likelihoods
- C) The sum of all transition probabilities
- D) The maximum likelihood estimation

Answer: B

192. The Backward algorithm is used to:

- A) Predict the next observation in a sequence
- B) Compute the likelihood of future observations
- C) Compute the probability of the partial observation sequence from  $t+1$  to  $T$
- D) Find the optimal state sequence

Answer: C

193. What is the main difference between the Forward and Backward algorithms?

- A) Forward algorithm works with states, while Backward works with observations
- B) Forward algorithm calculates probabilities from start to end, while Backward calculates from end to start
- C) Forward is for supervised learning, while Backward is for unsupervised learning
- D) Forward uses initial probabilities, while Backward uses transition probabilities

Answer: B

194. In the Forward algorithm, the recursion formula involves which of the following components?

- A) Transition probabilities, emission probabilities, and previous forward values
- B) Transition probabilities and final state probabilities
- C) Only emission probabilities
- D) Only transition probabilities

Answer: A

195. Which of the following is a primary use of the Forward-Backward algorithm?

- A) Classification of sequential data
- B) Parameter estimation for HMMs
- C) Dimensionality reduction
- D) Clustering data points

Answer: B

196. In the Backward algorithm, what is the base case at time  $t=T$ ?

- A) All probabilities are set to 0
- B) All probabilities are set to 1
- C) Probabilities are initialized randomly
- D) Probabilities are based on observation likelihoods

Answer: B

197. The Forward-Backward algorithm is commonly used in which machine learning task?

- A) Clustering
- B) Regression
- C) Sequence labeling
- D) Dimensionality reduction

Answer: C

198. Which of the following is a primary goal of clustering in data analysis?

- A) To reduce the dimensions of data
- B) To group similar data points together
- C) To perform classification based on predefined labels
- D) To increase the variance within the data

Answer: B) To group similar data points together

199. Which of the following clustering algorithms is based on the concept of distance between data points?

- A) K-Means
- B) DBSCAN
- C) Agglomerative Hierarchical Clustering
- D) All of the above

Answer: D) All of the above

200. Which of the following is true about DBSCAN (Density-Based Spatial Clustering of Applications with Noise)?

- A) It requires the number of clusters to be specified in advance
- B) It can identify clusters of arbitrary shape
- C) It only works with numerical data
- D) It assigns all data points to clusters, without considering noise

Answer: B) It can identify clusters of arbitrary shape