

Q1) a) ① Lasso regression stands for Least Absolute Shrinkage and Selection Operator

② It is a type of linear regression that includes a regularization term.

③ The regularization term is L1 penalty.

④ The goal of lasso regression is to minimize the sum of RSS i.e. Residual sum of squares with the absolute values of coefficients by tuning parameter λ

$$\text{minimize } \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j| \right)$$

where,

y_i = observed response

\hat{y}_i = predicted response,

β_j = model coefficients

λ = regularization parameter

⑤ Feature Selection in Lasso:-

Lasso performs automatic feature selection by forcing coefficients of some features to be exactly zero.

The process is as:-

a) L1 penalty added to cost function encourages sparsity.

b) feature with coefficients shrinking to zero are removed model

c) This allows lasso to select only the most important features, making the model more efficient and reducing the multi colinearity.

(Q1b) ① Linear Regression:-
Models the relationship between a dependent variable and independent variable using a straight line.

② Ridge Regression:-
A variant of linear regression with L_2 regularization to prevent overfitting

③ Lasso Regression:-
Similar to Ridge but uses L_1 regularization, which shrink some zero and performing feature selection.

④ Elastic Net Regression:
Combines L_1 and L_2 by balancing feature selection.

⑤ Polynomial Regression:
Extends linear regression by modelling non-linear relationships using polynomial terms.

⑥ Logistic Regression:
Used for binary classification, models and outcomes of 0 or 1.

⑦ KNN:-

K-Nearest Neighbor Regression. It predicts the value based on average of K nearest neighbors.

Q1(c) ① Bias Variance Trade off

① Bias :-

The error is by ~~an~~ comes by assuming a model that don't capture the underlying models of data.

high bias leads to overfitting.

② Variance :-

The error is due to model's sensitivity to the data. High variance leads to overfitting, where the model becomes too complex.

* Relationship to Underfitting & Overfitting

- ① Underfitting occurs when the model has high bias and poor performance
- ② Overfitting occurs when the model has high variance and gives good performance

Q3) a) ① Kernel methods in SVM are used to transform non-linearly data into a higher-dimensional space where it becomes linearly separable.

- ② This is achieved without computing the coordinates of data in this higher-dimensional space.

* Common Kernel functions

- ① Linear kernel:-

No transformation is applied, and the data is assumed to be linearly separable.

- ② Polynomial Kernel:-

Maps the data into a higher-dimensional polynomial feature space.

- (3) RBF
Maps the data infinite dimensional space
- (3) The Kernel method allow SVM to learn non-linear decision boundaries by working in higher-dimensional space, without the need to compute the transformed data.
- (4) This helps SVM to classify the complex data.

Q3) Advantages of KNN:-

- (1) Simple and Easy:
It is easy to understand.
- (2) No Training Phase:
KNN is lazy learner it doesnot require a training phase.
- (3) Non-Parametric:
It makes no assumptions about data distribution, useful for complex and non-linear data.
- (4) flexible:
Can be used for both classification

disadvantages of KNN

- ① Computationally Expensive :-
The prediction phase is slow, it requires calculating distances to all training samples.
- ② Sensitive to Irrelevant features :-
Performance can degrade with high dimensional or noisy data.
- ③ Memory Intensive :-
Requires storing the dataset for prediction.
- ④ Choosing K :-
The choice of the right value K can affect the performance.

3.c) Metrics of KNN.

① Euclidean distance

The most widely used metric, measures the straight-line distance between two points.

② Manhattan distance

Measure the sum of absolute coordinates useful in grid-like differences of data.

③ Minkowski Distance:-
Generalizes Euclidean and Manhattan distances by introducing a parameter p

④ Cosine Similarity:-
measures the cosine similarity of data.

⑤ Chebyshev Distance
measures the maximum absolute difference along any dimension.

Q6.) a) Isolation forest

① Isolation forest model is an anomaly selection algorithm that works by isolating the selected splits

② How it works:-

1) Isolation:-
The algo builds multiple random trees by selecting random features and splitting the data at random values.

2) Isolation of anomalies:-
Anomalies are isolated in some steps because they are different

③ Scoring:

The number of steps required to isolate a data point is used as a measure of its anomaly score.

④ The advantages of Isolation forest are efficient and No assumptions.

Q. 6 (b) Hierarchical Clustering is a method of cluster analysis that seeks to build a hierarchy of clusters:-

Two types:-

① Bottom - Up:-

Starts with each data point and merges the clusters.

② Top - Down:-

Starts with all points in a single cluster and splits into smaller clusters.

* The process of Clustering is:

① Start

② Merge

③ Repeat.

Eg:-
Consider a dataset of 5 points A, B, C, D, E

Step 1:-
Start with individual
 $\{A\}, \{B\}, \{C\}, \{D\}, \{E\}$.

Step 2:-
merge the closest pair
 $\{A, B\}$.

Step 3:-
calculate distance between
 $\{A, B\}, \{C\}, \{D\}, \{E\}$ and merge

Step 4:-
Repeat the process until one cluster
remains.

Q.10) ① Micro - Average Precision and Recall.

① Micro - Avg Precision :-
Calculated by considering the total true
positives, false positives and false
negatives then calculate precision.

Formula :- Micro precision = $\frac{\sum \text{True positives}}{\sum (\text{True positives} + \text{false positives})}$

② Micro-Average Recall: -
Similar to precision, it aggregates true positives and false negatives

Formula: -
$$\text{Micro Recall} = \frac{\sum \text{True Positives}}{\sum (\text{True Positives} + \text{False Negatives})}$$

③ Micro-Average F-score.

① The harmonic mean of micro-average precision and recall, offering a single score that balance both metrics

Formula :-
$$\text{Micro-F1} = 2 \times \frac{\text{Micro-Precision} \times \text{Micro-Recall}}{\text{Micro-Precision} + \text{Micro-Recall}}$$

Q7) a) ① RNN stands for Recurrent Neural Networks.

② They are a type of neural network used in data

③ They allow them to maintain a memory of previous inputs.

④ This makes them particularly useful for tasks like time series prediction.

Eg:-
Predicting the Next Word in Sentence.

- ① Input sequence:-
" I love programming "
- ② The RNN takes "I" as input and processes it.
- ③ Then takes "love" as input, updates the state again and so on.
- ④ At every step, the RNN uses the memory from previous words to predict the next word.

Q7.) b) The activation functions.

① Sigmoid:-

Output value between 0 and 1, it is a binary classification.

$$\text{formula : } \sigma(x) = \frac{1}{1 + e^{-x}}$$

② Tan h:-

Output value between -1 and 1, a scaled version of sigmoid

formule:- $\tanh(x) = \frac{2}{1 + e^{-2x}} - 1$

③ ReLU (Rectified Linear Unit)

Outputs the input if positive, otherwise zero, its mostly used in activation function.

formule:- $\text{ReLU}(x) = \max(0, x)$

④ Leaky ReLU:-

It allows a small or non-zero gradient for negative inputs.

formule:- $\text{leaky ReLU}(x) = \max(\alpha x, x)$

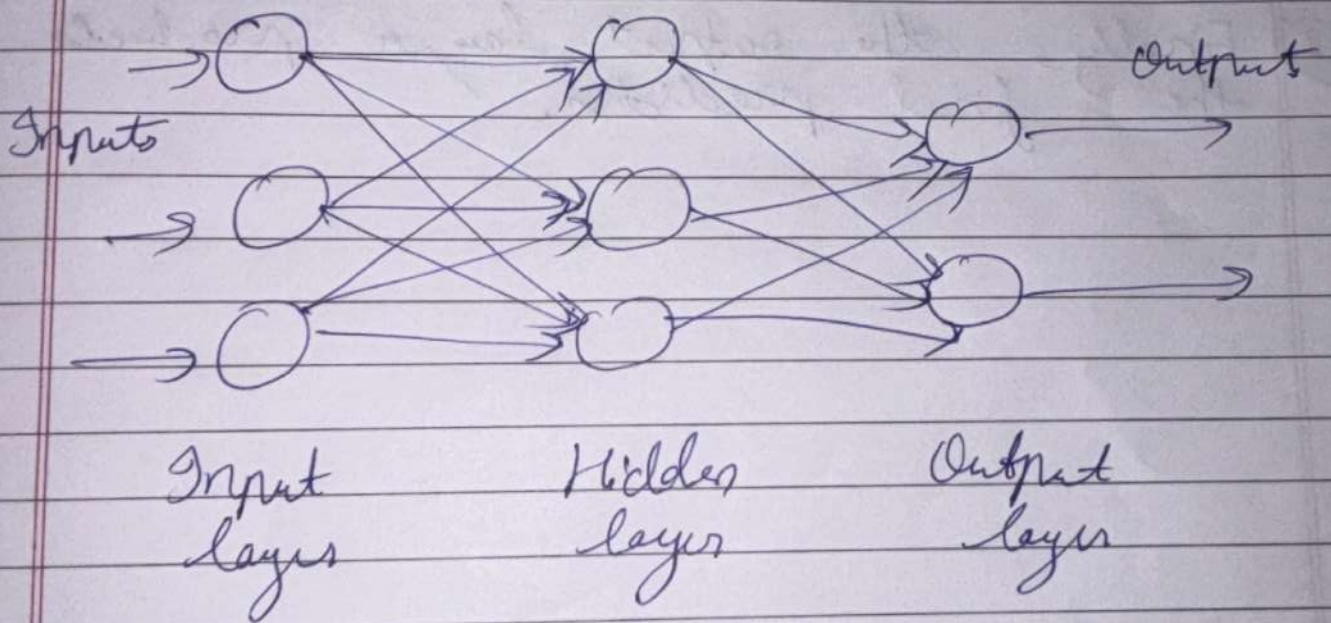
⑤ Softmax:-

Used for multi-class classification, it converts logits into probability distribution over multiple classes

formule:- $\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$

Q7.9) ① MLP (Multilayer Perceptron)

- ② It is a type of feedforward neural network consisting of multiple layers of neurons.
- ③ It has input layers, hidden layers, output layers.



1.7 Structure of MLP:-

① Input layer:-

The first layer that gets the input data.

② Hidden layer:-

One or more layers where the activation functions are performed.

③ The output layer is
The produced output typically passes
activation function like softmax.

④ How it works :-

① Data is passed from input layer to
the first hidden layer then process
to the output layer.

② Finally, the output layer produces
the final prediction.