

Assignment -2

1. In logistic regression, what is the logistic function (sigmoid function) and how is it used to compute probabilities?

Ans:-

The sigmoid function is a mathematical function used to map the predicted values to probabilities. It maps any real value into another value within a range of 0 and 1. The value of the logistic regression must be between 0 and 1, which cannot go beyond this limit, so it forms a curve like the "S" form.

For $y=1$ The predicted probabilities will be:

$$P(X;b,w) = p(x)$$

For $y = 0$ the predicted probabilities will be: $1-p(X;b,w)=1-p(x)$

2. When constructing a decision tree, what criterion is commonly used to split nodes, and how is it calculated?

Ans:-

The commonly used criterion to split nodes in decision trees is to maximize the information gain or minimize the impurity of the nodes.

One popular impurity measure is the Gini impurity, which measures the probability of incorrectly classifying a randomly chosen element if it were randomly labeled according to the distribution of labels in the node.

Another common criterion is entropy, which measures the randomness or uncertainty in the distribution of class labels in a node.

The decision tree algorithm evaluates different splits based on these criteria and chooses the split that maximizes information gain or minimizes impurity. average of the child it have yes or No in the observations in the dataset.

$$\text{Information Gain} = \text{Entropy}(S) - [(\text{Weighted Avg}) * \text{Entropy}(\text{each feature})]$$

3. Explain the concept of entropy and information gain in the context of decision tree Construction?

Ans:-

Entropy is a measure of randomness or uncertainty in a set of data. In the context of decision tree construction, entropy is calculated for each node to determine the homogeneity of the class labels

Information gain is the measure of the effectiveness of a particular attribute in classifying the data. It quantifies the reduction in entropy (or impurity) achieved by splitting the data on a particular attribute. Randomness and disorder.

Follows as $H(x) = -\sum_{k \in K} p(k) \log(p(k))$

When constructing a decision tree, the algorithm selects the attribute that maximizes the information gain at each split, aiming to create nodes with the highest possible homogeneity.

4. How does the random forest algorithm utilize bagging and feature randomization to improve classification accuracy?

Ans:-

Random forest is an ensemble learning technique that combines multiple decision trees to improve classification accuracy.

Bagging (Bootstrap Aggregating) is used in random forests to train each decision tree on a random subset of the training data with replacement.

Feature randomization is employed by randomly selecting a subset of features at each split in each decision tree, which helps to decorrelate the trees and reduce overfitting.

By averaging the predictions of multiple trees trained on different subsets of data and features, random forests can reduce variance and improve generalization performance

In Bagging, we apply a similar concept. However, if we train multiple models without making any adjustments, these models may not exhibit the necessary diversity. To address this, we use a technique called bootstrap to train different models, ensuring that each model has a distinct capability. This approach helps us obtain a set of diverse models.

5. What distance metric is typically used in k-nearest neighbours (KNN) classification, and how does it impact the algorithm's performance?

Ans:-

The Euclidean distance metric is typically used in KNN classification, although other distance metrics such as Manhattan distance or Minkowski distance can also be used.

Euclidean distance measures the straight-line distance between two points in a multidimensional space.

The choice of distance metric can impact the performance of the KNN algorithm, as it determines how the "closeness" of neighbours is calculated and influences the decision boundaries.

6. Describe the Naïve-Bayes assumption of feature independence and its implications for classification.

Ans:-

The Naïve-Bayes algorithm assumes that features are conditionally independent given the class label.

This means that the presence or absence of a particular feature is independent of the presence or absence of other features, given the class label.

Despite its simplifying assumption, Naïve-Bayes classifiers can perform well in practice, especially on text classification tasks.

7. In SVMs, what is the role of the kernel function, and what are some commonly used kernel functions?

Ans:-

In SVMs (Support Vector Machines), the kernel function is used to transform the input features into a higher-dimensional space where a linear decision boundary can be found.

Kernel functions measure the similarity between pairs of data points in the input space.

Some commonly used kernel functions include linear kernel, polynomial kernel, radial basis function (RBF) kernel (Gaussian kernel), and sigmoid kernel.

The choice of kernel function can significantly impact the performance and flexibility of the SVM model.

8. Discuss the bias-variance, trade-off in the context of model complexity and overfitting?

Ans: The bias-variance ,trade-off refers to the balance between the bias (error due to the model's simplifying assumptions) and variance (sensitivity to small fluctuations in the training data) of a machine learning model. Increasing the model complexity typically reduces bias but increases variance, and vice versa.

Overfitting occurs when a model has low bias but high variance, capturing noise in the training data rather than the underlying patterns.

Underfitting occurs when a model has high bias and low variance, failing to capture the complexity of the underlying data.

Finding the right balance between bias and variance is essential to avoid overfitting or underfitting and achieve good generalization performance.

9. How does TensorFlow facilitate the creation and training of neural networks?

Ans:-

TensorFlow is an open-source machine learning library developed by Google for building and training neural networks and other machine learning models.

TensorFlow provides a flexible and efficient framework for constructing computational graphs and executing operations on large datasets, making it suitable for training deep learning models.

It offers high-level APIs like Keras for building neural networks with ease and low-level APIs for more flexibility and control over the model architecture and training process.

TensorFlow supports distributed computing and can run on CPUs, GPUs, and TPUs (Tensor Processing Units) to accelerate training and inference.

10. Explain the concept of Cross-validation is a technique used to assess the performance and generalization ability of a machine learning model.

ANS:-

It involves splitting the dataset into multiple subsets (folds), training the model on a subset of the data, and evaluating it on the remaining fold(s).

Common types of cross-validation include k-fold cross-validation, where the dataset is divided into k equal-sized folds, and leave-one-out cross-validation, where each data point is used as a separate test set.

Cross-validation helps to estimate the model's performance on unseen data and detect potential issues such as overfitting or data leakage.

11. What techniques can be employed to handle overfitting in machine learning models?

Ans:-

Overfitting occurs when a machine learning model learns the training data too well, capturing noise or random fluctuations in the data instead of the underlying patterns. To address overfitting, several techniques can be employed:

Cross-validation:

Cross-validation is a technique used to assess the performance and generalization ability of a model by splitting the dataset into multiple subsets (folds), training

the model on a subset of the data, and evaluating it on the remaining fold(s). Crossvalidation helps to estimate the model's performance on unseen data and detect potential issues such as overfitting or data leakage.

Regularization:

Regularization techniques add a penalty term to the loss function, discouraging overly complex models with large weights. Two common types of regularization are.

L1 regularization (Lasso): Adds the sum of the absolute values of the model weights to the loss function, encouraging sparsity and feature selection.

L2 regularization (Ridge): Adds the sum of the squares of the model weights to the loss function, penalizing large weights and preventing overfitting.

Ensemble Methods: Ensemble methods combine multiple models to improve predictive performance and reduce overfitting. Techniques such as bagging, boosting, and random

forests train multiple models on different subsets of the data or using different algorithms and aggregate their predictions to make the final prediction

12. What is the purpose of regularization in machine learning, and how does it work?

Ans:-

Regularization is a technique used to prevent overfitting by adding a penalty term to the loss function, discouraging overly complex models.

L1 regularization (Lasso) adds the sum of the absolute values of the model weights to the loss function, encouraging sparsity and feature selection.

L2 regularization (Ridge) adds the sum of the squares of the model weights to the

13. Describe the role of hyper-parameters in machine learning models and how they are tuned for optimal performance.?

Ans:-

Hyperparameters are parameters that are not learned directly from the data during the training process but are set prior to training. They control the behavior and performance of

the machine learning model, influencing aspects such as model complexity, regularization, and optimization.

14. What are precision and recall, and how do they differ from accuracy in classification evaluation?

Ans:- Precision:

Precision measures the proportion of true positive predictions (correct positive predictions) among all instances predicted as positive by the model.

It represents the ability of the classifier to avoid false positives, i.e., instances that were predicted as positive but are actually negative.

Precision is calculated as the ratio of true positives (TP) to the sum of true positives and

false positives (FP):

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

A high precision indicates that the model has a low rate of false positives and is confident when it predicts a positive class.

Recall: Recall, also known as sensitivity or true positive rate, measures the proportion of true positive predictions among all actual positive instances in the dataset. It represents the ability of the classifier to capture all positive instances and avoid false

negatives, i.e., instances that were predicted as negative but are actually positive.

Recall is calculated as the ratio of true positives (TP) to the sum of true positives and false

negatives (FN):

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

A high recall indicates that the model can identify most of the positive instances in the dataset, even at the cost of higher false positive predictions.

Accuracy:

Accuracy measures the overall correctness of the model's predictions, irrespective of class labels. It represents the proportion of correct predictions (true positives and true negatives) among all instances in the dataset. Accuracy is calculated as the ratio of the sum of true positives and true negatives (TP + TN)

to the total number of instances:

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

While accuracy is a useful metric for balanced datasets, it may not accurately reflect the performance of a model on imbalanced datasets where the classes are disproportionately represented.

15. Explain the ROC curve and how it is used to visualize the performance of binary classifiers.

Ans:-

True Positive Rate (TPR):

TPR, also known as sensitivity or recall, measures the proportion of positive instances that are correctly identified by the classifier. It is calculated as:

$$\text{TPR} = \text{TP} / (\text{TP} + \text{FN})$$

False Positive Rate (FPR):

FPR measures the proportion of negative instances that are incorrectly classified as positive by the classifier. It is calculated as:

$$\text{FPR} = \text{FP} / (\text{FP} + \text{TN})$$

ROC Curve:

The ROC curve is created by plotting the TPR (sensitivity) on the y-axis against the FPR (1 -specificity) on the x-axis for different threshold values.

Each point on the ROC curve represents the trade-off between correctly identifying positive instances (high TPR) and incorrectly classifying negative instances as positive (high FPR) at a particular threshold setting.

A classifier with perfect performance would have an ROC curve that passes through the upper-left corner (TPR = 1, FPR = 0), indicating high sensitivity and low false positive rate across all threshold values.

Area Under the ROC Curve (AUC-ROC):

The Area Under the ROC Curve (AUC-ROC) quantifies the overall performance of the classifier across all possible threshold settings. AUC-ROC ranges between 0 and 1, where a value closer to 1 indicates better classifier performance. A classifier with an AUC-ROC of 0.5 represents random guessing, while an AUC-ROC of 1 represents perfect classification. AUC-ROC provides a single scalar value to compare the performance of different classifiers, making it a useful metric for model selection and evaluation.

$$\text{ROC} = (\text{current value} / \text{previous value} - 1) * 100$$