

20th April Assignment

April 29, 2023

1 Assignment 73

Q1. What is the KNN algorithm?

Ans.KNN is also known as k-nearest neighbour algorithm. it is a supervised machine learning algorithm. used for classification and regression. k refers to number of nearest neighbour. To make a prediction for a new data point, the algorithm calculates the distance between that data point and all other data points in the dataset. The K nearest data points are then used to make the prediction, based on the majority class (for classification) or average value (for regression) of the K neighbors. The KNN algorithm is non-parametric, meaning it does not make any assumptions about the underlying distribution of the data. It is also instance-based, meaning it does not learn a model from the training data but rather stores the entire dataset and uses it to make predictions for new data points.

Q2. How do you choose the value of K in KNN?

Ans.There is no fixed rule to choose the value of K in KNN, and it depends on the characteristics of the dataset and the problem at hand. Cross-validation and domain knowledge are two common approaches to select the optimal value of K.

Q3. What is the difference between KNN classifier and KNN regressor?

Ans.In KNN classification, the goal is to predict the class membership of a new data point based on its nearest neighbors. The output of the KNN classifier is a class label, which is a categorical variable. For example, in a binary classification problem, the output can be either 0 or 1, representing the two possible classes.

In KNN regression, the goal is to predict a continuous variable (i.e., a numeric value) for a new data point based on its nearest neighbors. The output of the KNN regressor is a numeric value, which can be any real number. For example, in a regression problem that predicts the price of a house based on its features, the output can be any positive real number.

Q4. How do you measure the performance of KNN?

Ans. Classification Metrics:

1. Accuracy: The proportion of correctly classified instances to the total number of instances.
2. Precision: The proportion of correctly classified positive instances to the total number of predicted positive instances.
3. Recall: The proportion of correctly classified positive instances to the total number of actual positive instances.
4. F1-score: The harmonic mean of precision and recall.

Regression Metrics:

1. Mean Absolute Error (MAE): The average absolute difference between the predicted and actual values.
2. Mean Squared Error (MSE): The average squared difference between the predicted and actual values.
3. Root Mean Squared Error (RMSE): The square root of the MSE.
4. R-squared (R2): The proportion of variance in the target variable that is explained by the model.

To evaluate the performance of the KNN algorithm, you can split the dataset into training and testing sets, fit the KNN model on the training set, and then make predictions on the testing set. You can then use the above-mentioned evaluation metrics to measure the performance of the KNN algorithm.

Q5. What is the curse of dimensionality in KNN?

Ans. The curse of dimensionality is a problem in KNN where the performance of the algorithm degrades as the number of dimensions in the dataset increases. This can lead to sparsity and increased computational cost, and various approaches such as dimensionality reduction and feature selection can be used to address this problem.

Q6. How do you handle missing values in KNN?

Ans. The ways to handle missing values in KNN are:

1. Deletion: One approach to handling missing values is to simply delete any data points with missing values. However, this approach can lead to a loss of data, especially if the percentage of missing values is high.
2. Imputation: Imputation is the process of filling in missing values with estimated values. There are several imputation techniques that can be used, including mean imputation, median imputation, mode imputation, and k-NN imputation.
3. Distance-based imputation: Distance-based imputation is a technique that is specific to the KNN algorithm. In this approach, the distance between the nearest neighbors is used to impute missing values. For example, if a feature value is missing for a data point, the missing value can be replaced with the weighted average of the feature values of the nearest neighbors, with weights based on their distance from the data point.

4. Algorithm-specific techniques: Some algorithms may have specific techniques for handling missing values. For example, the `KNNImputer` class in `scikit-learn` is an implementation of the KNN algorithm specifically designed for imputing missing values.

Q7. Compare and contrast the performance of the KNN classifier and regressor. Which one is better for which type of problem?

Ans. The KNN algorithm can be used for both classification and regression problems, and the performance of the KNN classifier and regressor can be evaluated using different metrics. The choice of which one to use depends on the specific problem and the characteristics of the data.

Q8. What are the strengths and weaknesses of the KNN algorithm for classification and regression tasks, and how can these be addressed?

Ans. Strengths of the KNN algorithm:

1. Non-parametric: The KNN algorithm is non-parametric, which means that it does not make any assumptions about the underlying distribution of the data. This makes it more flexible than some other algorithms that assume specific distributions.
2. Intuitive: The KNN algorithm is easy to understand and implement, making it a good choice for beginners or for problems where interpretability is important.
3. Good performance on small datasets: The KNN algorithm works well on small datasets, especially when the number of features is relatively small.

Weaknesses of the KNN algorithm:

1. Computationally expensive: The KNN algorithm can be computationally expensive, especially for large datasets, since it requires computing distances between all pairs of data points.
2. Sensitive to the choice of k: The KNN algorithm is sensitive to the choice of k, which can significantly impact the performance of the algorithm.
3. Prone to overfitting: The KNN algorithm can be prone to overfitting, especially when the number of features is large or when there are noisy or irrelevant features.
4. Curse of dimensionality: The KNN algorithm can suffer from the curse of dimensionality in high-dimensional spaces, which can cause the algorithm to perform poorly or require a large number of nearest neighbors to make accurate predictions.

To address the weaknesses of the KNN algorithm, several techniques can be used:

1. Use dimensionality reduction techniques, such as PCA or t-SNE, to reduce the number of dimensions in the dataset.
2. Use feature selection techniques to select the most informative features that can improve the performance of the algorithm.
3. Use cross-validation techniques to tune the value of k and prevent overfitting.
4. Use distance metrics that are more appropriate for the specific problem, such as the cosine similarity metric for text data.

Q9. What is the difference between Euclidean distance and Manhattan distance in KNN?

Ans.Euclidean distance and Manhattan distance are two commonly used distance metrics in KNN algorithm, where Euclidean distance is more appropriate for continuous variables and sensitive to the outliers, and Manhattan distance is more appropriate for both continuous and categorical variables and robust to outliers.

Q10. What is the role of feature scaling in KNN?

Ans.Feature scaling is an important preprocessing step in KNN algorithm that can significantly impact the performance of the algorithm. The role of feature scaling is to transform the features of the dataset to have the same scale and reduce the impact of the features with larger magnitude on the distance calculation.