# 21th April Assignment

April 29, 2023

## 1 Assignment 74

**Q1. What is the main difference between the Euclidean distance metric and the Manhattan distance metric in KNN? How might this difference affect the performance of a KNN classifier or regressor?**

Ans.The main difference between the Euclidean distance metric and the Manhattan distance metric in KNN (K-Nearest Neighbors) is the way they measure the distance between data points.

Euclidean distance measures the straight-line distance between two points in a two-dimensional or multi-dimensional space. It calculates the square root of the sum of the squared differences between corresponding coordinates. Mathematically, it can be represented as:

d(x, y) = sqrt((x1 - y1)^2 + (x2 - y2)^2 + ... + (xn - yn)^2)

On the other hand, the Manhattan distance measures the distance between two points by calculating the sum of the absolute differences between corresponding coordinates. It calculates the distance as the sum of the differences in the horizontal and vertical directions (hence, the name Manhattan distance). Mathematically, it can be represented as:

d(x, y) = |x1 - y1| + |x2 - y2| + ... + |xn - yn|

The difference in the way the two distance metrics measure the distance between data points can affect the performance of a KNN classifier or regressor. In some cases, the choice of distance metric can lead to different results. For instance, if the data points are clustered together, the Euclidean distance metric may be more effective as it takes into account the overall distance between the points. However, if the data points are dispersed or have outliers, the Manhattan distance metric may be more robust as it only considers the absolute differences between coordinates.

Therefore, the choice of distance metric in KNN should be based on the nature of the data and the problem being solved. In general, it is advisable to try both distance metrics and select the one that performs better for the given dataset and problem.

**Q2.** **How do you choose the optimal value of k for a KNN classifier or regressor? What techniques can be used to determine the optimal k value?**

**Ans. The technique which we use to determine optimal K values are:**

1. Cross-validation: Cross-validation is a common technique used to evaluate the performance of a KNN model with different k values. In k-fold cross-validation, the data is split into k equal parts, and the model is trained on k-1 parts and tested on the remaining part. This process is repeated k times, and the performance of the model is averaged over the k iterations. The optimal k value is the one that produces the highest accuracy or lowest error rate on the test data.

2. Grid search: Grid search is a brute-force technique used to search for the optimal k value by evaluating the model's performance for a range of k values. The range of k values is specified beforehand, and the model's performance is evaluated using cross-validation. The optimal k value is the one that produces the highest accuracy or lowest error rate.

3. Elbow method: The elbow method is a graphical technique used to determine the optimal k value by plotting the error rate or accuracy against the k values. The plot looks like an elbow, and the optimal k value is the point where the error rate or accuracy starts to plateau or does not improve significantly.

4. Distance plot: The distance plot is a technique used to determine the optimal k value by plotting the distance between each data point and its kth nearest neighbor. The plot looks like a step function, and the optimal k value is the one that produces the steepest drop in distance.

**Q3.** **How does the choice of distance metric affect the performance of a KNN classifier or regressor? In what situations might you choose one distance metric over the other?**

**Ans.Minkowski distance, cosine similarity, and Mahalanobis distance, can also be used in KNN. The choice of distance metric should be based on the characteristics of the data and the problem being solved. For instance:**

**If the data has high-dimensional features, cosine similarity can be a better metric as it measures the angle between the feature vectors rather than the distance.**

**If the data has features with different units or scales, normalization can be applied, and Manhattan distance can be used as it is insensitive to scale.**

**If the data has outliers or high variance in some features, Mahalanobis distance can be used as it takes into account the covariance matrix of the data.**

**Q4. What are some common hyperparameters in KNN classifiers and regressors, and how do they affect the performance of the model? How might you go about tuning these hyperparameters to improve model performance?**

**Ans. The common hyperameter in KNN classifier and regressor are:**

1. k: The number of nearest neighbors to consider. A larger k value can lead to smoother decision boundaries, while a smaller k value can lead to overfitting.

2. Distance metric: The distance metric used to calculate the distance between data points. Different distance metrics can be more or less suitable depending on the nature of the data.

3. Weight function: The weight function used to weight the contribution of each neighbor to the prediction. A uniform weight function treats all neighbors equally, while a distance-weighted function gives more weight to closer neighbors.

4. Algorithm: The algorithm used to find the nearest neighbors. The brute-force algorithm is straightforward but computationally expensive, while tree-based algorithms, such as kd-tree or ball tree, can be faster but may require more memory.

**Q5. How does the size of the training set affect the performance of a KNN classifier or regressor? What techniques can be used to optimize the size of the training set?**

**Ans.The size of the training set can have a significant impact on the performance of a KNN classifier or regressor. In general, increasing the size of the training set can improve the model's accuracy and reduce overfitting, but it can also increase the computational complexity and training time.**

**Q6. What are some potential drawbacks of using KNN as a classifier or regressor? How might you overcome these drawbacks to improve the performance of the model?**

**Ans. The potential drawbacks of using KNN as a classifer or regressor are:**

1. Computational complexity: KNN requires computing the distance between each data point and all the training samples, which can be computationally expensive for large datasets.

2. Curse of dimensionality: As the number of dimensions in the data increases, the number of neighbors needed to accurately estimate the target variable increases exponentially, making KNN less effective in high-dimensional spaces.

3. Imbalanced datasets: KNN can be biased towards the majority class in imbalanced datasets, leading to poor performance on the minority class.

4. Sensitivity to outliers: KNN is sensitive to outliers and noise in the data, which can affect the prediction accuracy.