

How does K affect model performance on evaluation sets?

As the of k increases the algorithm has more data to learn from. But at certain point the precision stops increasing. We can see from the graph obtained that for both models and we can select optimal k as 5 for SVM and 10 for NB as increment of K after that point has showed very or no improvement in the performance.

Therefore, by taking small values of K, there may not provide enough training data for the model to train to provide a very accurate result on the test set.

Similarly, by taking a very large value of K, may cause overfitting on the training data and also reduce the performance of the test set.

By varying K , you also change sample size for training. How can the sample size affect model performance?

The more the amount of data the algorithm has the more the algorithm can learn and improve accuracy. Hence very large datasets are preferred.

By increasing the value of K, we increase the training dataset gets split into K chunks having a different test and training set for each of the K iterations. By increasing K by a lot, we may cause the model to overfit and give bad results on a new test set. Similarly, small value of K may not provide enough data do train the model properly.