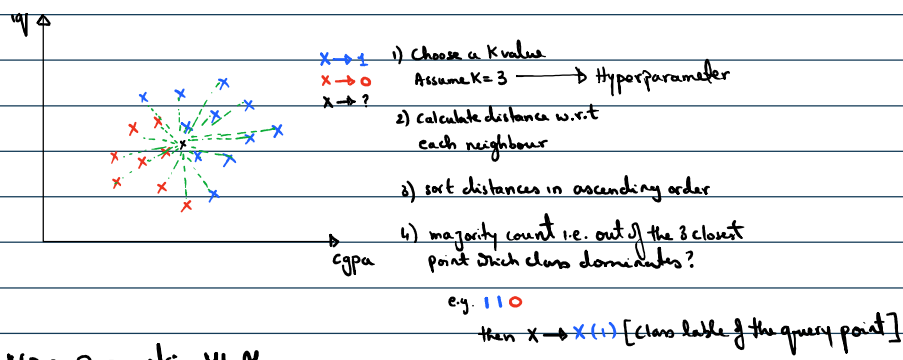


Algorithm Working

Sunday 14 April 2024 3:33 AM



\rightarrow Non-parametric ML Algo

\rightarrow This logic is applicable in n-D also

\rightarrow evaluation metric: accuracy score

* bring all your data on the same scale [KNN relies solely on distances]

```
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.metrics import accuracy_score
```

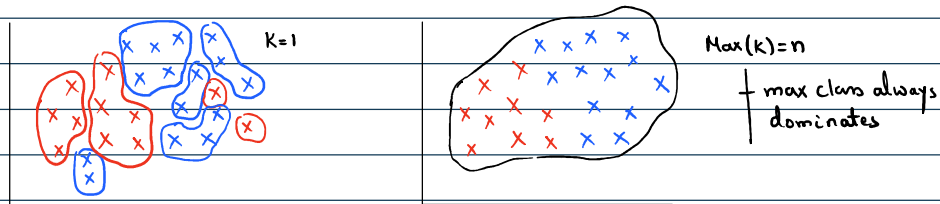
\rightarrow How to choose K?

- avoid even numbers
- experiment/hyperparameter tuning

\rightarrow Decision Surface/Boundary

- visually divides the coordinate space into 'n' regions, where n is the # of classes in the target variable.

\rightarrow Overfitting & Underfitting in KNN



\rightarrow Limitations of KNN

- Large datasets, as KNN is a lazy learner i.e. NO work during the training stage, all work is done when we have a query point.
need to calculate distances \rightarrow sort \rightarrow majority [slow/Latency]

- High dimensional data \rightarrow Curse of dimensionality

The concept of calculating distances in higher dimensions is not reliable, but KNN totally depends on distances. If the distances are distorted so are our results.

- Outliers
need to choose K wisely else we risk overfitting

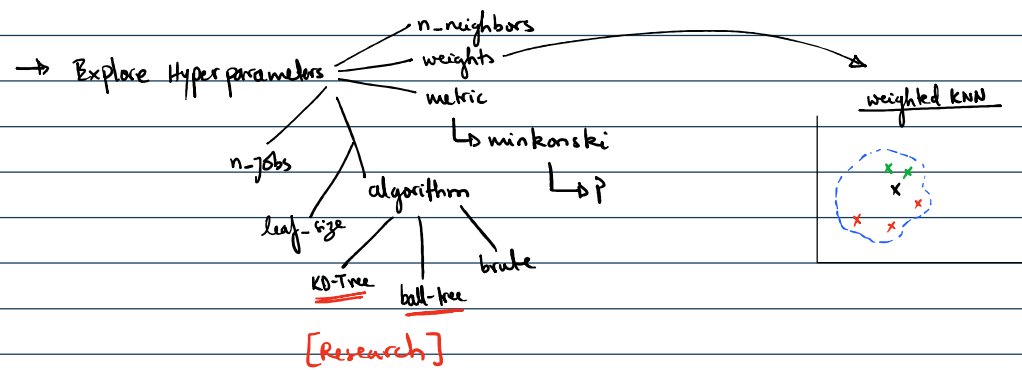
- Non homogeneous scales [features with higher values dominate & calculated distances are not reliable]

- Imbalanced dataset

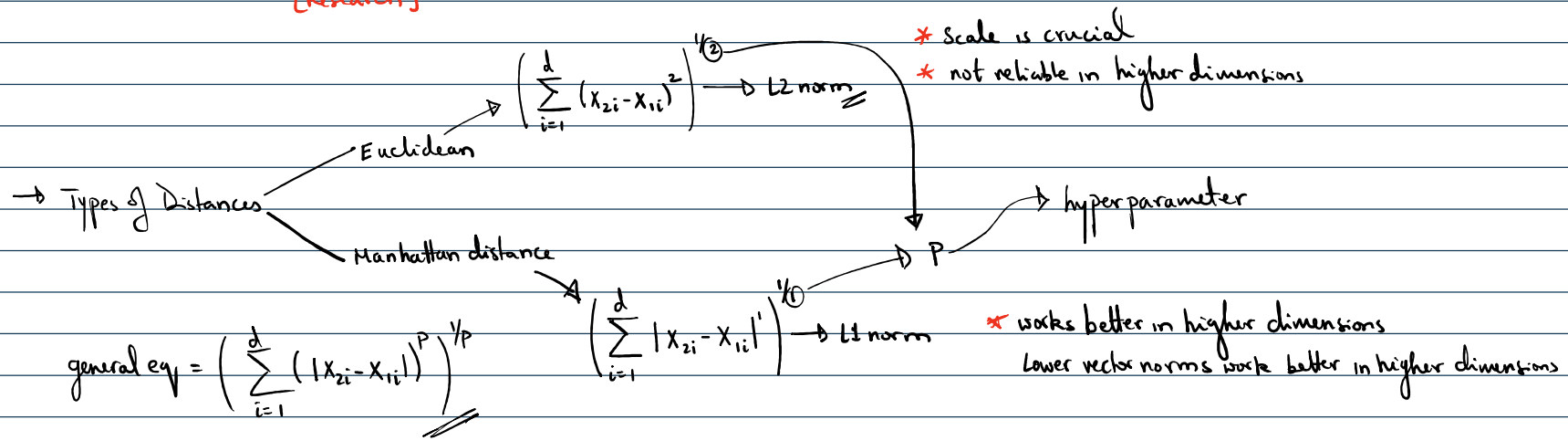
- good for prediction but not for inference [acts as a black box model]

i.e. How much does feature 1 contribute to the prediction compared to feature 2, unable to deduce feature importance

- KNN can be used for regression
- decide a K value
 - find K nearest neighbors
 - Average/Median of K nearest neighbors



uniform (majority count) $x \rightarrow x$
 distance \rightarrow closer the point \uparrow weight = $\frac{1}{\text{distance}}$ [add weights, pick class with highest weight]



→ Time/space complexity $O(nd)$

After some digging, I have some good answers. First let me tell you that as mentioned by some users like @anasvaf, you should only use odd number for binary classification. This is absolutely untrue. Firstly, when we use majority voting for binary classification, on some tie, it is entirely dependent on the actual library to choose the action. For example, in scikit-learn, it takes the mode of the variable. This means if in the training dataset, number of datapoints for class 1 is more, then 1 would be used on the tie. But there is a better solution.

What we can use weighted-KNN instead of normal KNN. In the weighed KNN, if there is a tie, we can see total distance for 1 labeled datapoints from k points and 0 labelled points. If total distance for 1 is more, then class would be 0 and vice-versa.

There are other good techniques too to handle tie in KNN but to be very honest, KNN is not a good learning algorithm specially due to its space of time complexity on large datasets.