

KNN (k-Nearest Neighbors)

①

What is KNN (k-Nearest Neighbors) Algorithm?

The k-Nearest Neighbor (KNN) algorithm is a popular machine learning technique used for classification and regression tasks. It relies on the idea that similar data points tend to have similar labels or values.

During the training phase, the KNN algorithm stores the entire training dataset as a reference. When making predictions, it calculates the distance between the input data point and all training examples, using a chosen distance metric such as Euclidean distance.

Next, the algorithm identifies the k nearest neighbors to the input data point based on their distances. In the case of classification, the algorithm assigns the most common class label among the k neighbors as the predicted label for the input data point. For regression, it calculates the average or weighted average of the target value of the k neighbors to predict the value for the input datapoint.

The KNN algorithm is straightforward and easy to understand, making it a popular choice in various domains. However, its performance can be affected by the choice of k and the distance metric, so careful parameter tuning is necessary for optimal results.

When Do We Use the KNN Algorithm?

KNN can be used for both classification and regression predictive problems. However, it is more widely used in classification problem in the industry. To evaluate any technique, we generally look at 3 important aspects:

1. Ease of interpreting output
2. Calculation time
3. Predictive Power

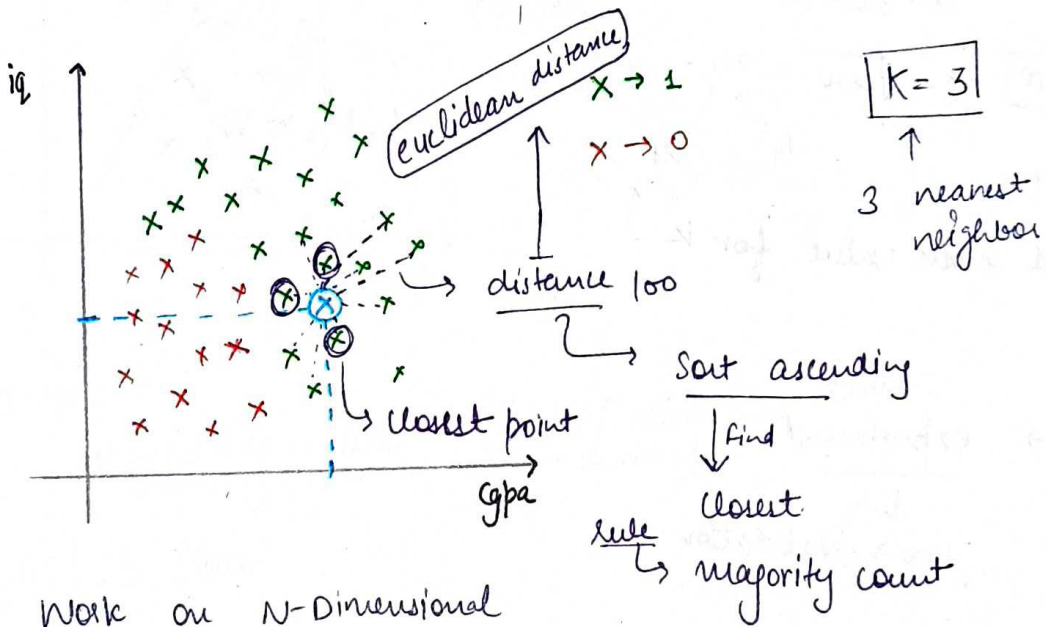
KNN classifier fair across all parameter of consideration, It is commonly used for its ease of interpretation and low calculation time.

KNN Intuition

(2)

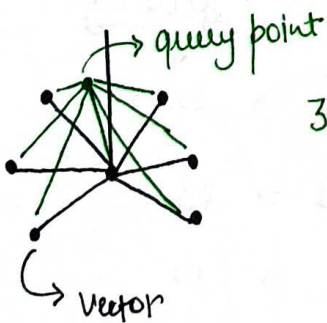
100 student

cgpa	iq	placement
8	80	1
7	70	0
⋮	⋮	⋮



* Work on N-Dimensional

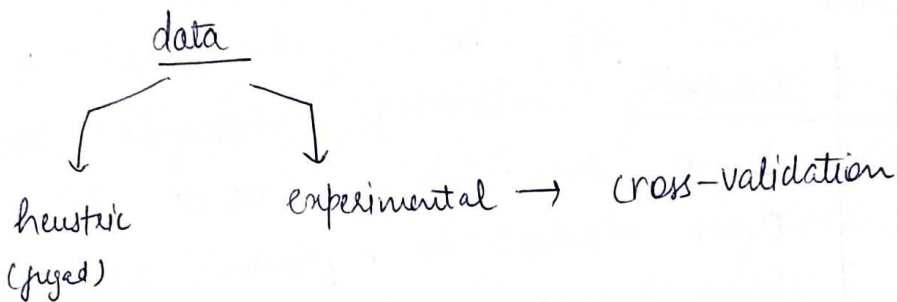
$x_1 \mid x_2 \mid x_3 \dots x_n$



3 vectors \rightarrow majority count \rightarrow class label of query point

1 1 0 \rightarrow ①
 closest point \uparrow majority class

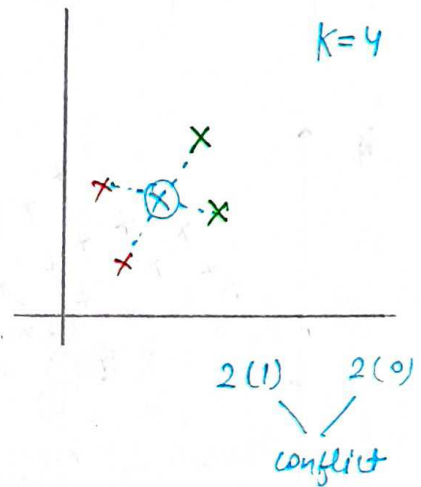
How to Select K?



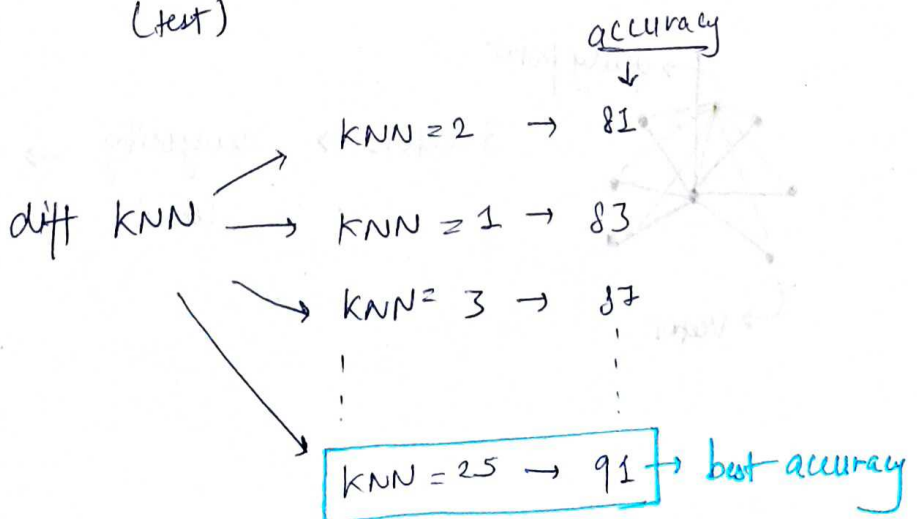
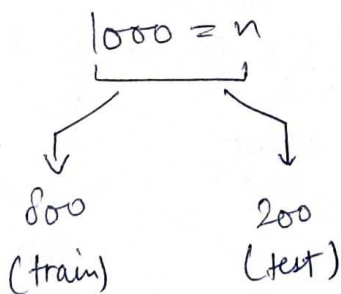
\sqrt{n} $n \rightarrow$ observation

let $\boxed{400}$ $\rightarrow \sqrt{400} = 20$
 \uparrow data
 $19 \quad 21$

* avoid even value for k

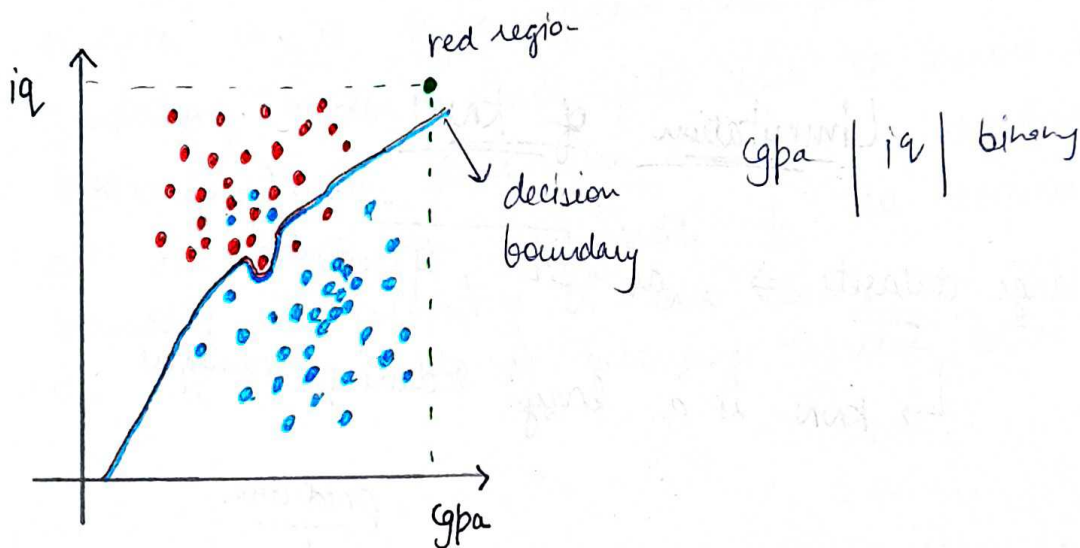
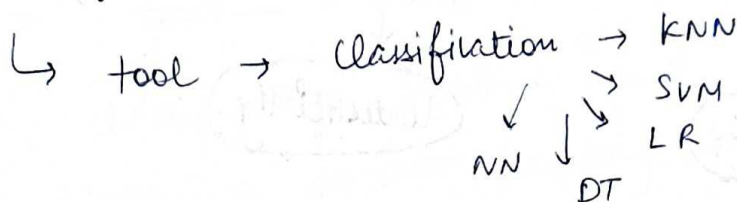


Best \rightarrow experiment
 \downarrow
cross-validation



Decision Surface

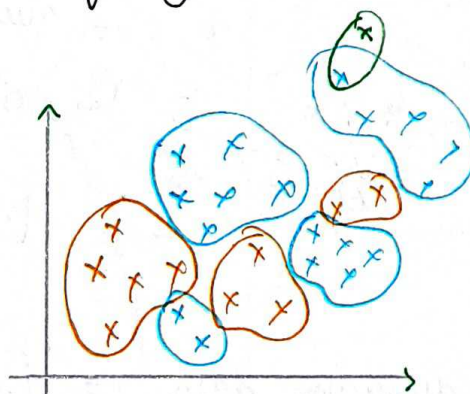
③



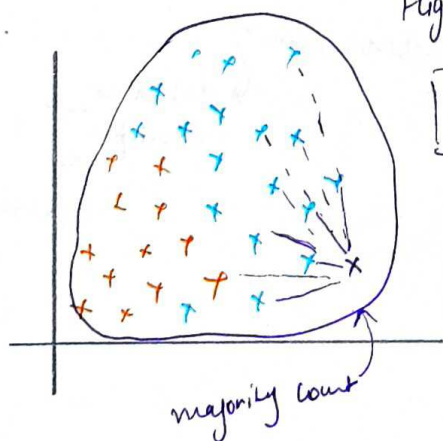
Overfitting and Underfitting in KNN

gpa | iq | placed

200 student



minor changes → overfitting
 ↓
 High variance



Highest value ↓

k=?

k=200

blue > orange

n=200

simplistic

↓
underfitting

$k = \text{low}$
↓
overfitting

$k = \text{high}$
↓
Underfitting

Limitation of KNN

1) large datasets $\Rightarrow n = 5L, f = 100 \rightarrow (500000, 100)$

↳ KNN is a lazy learning technique



prediction

↓
Slow \rightarrow dataset

\rightarrow query point \rightarrow prediction

↳ training \rightarrow nothing

$[5L \text{ distance}] \rightarrow [\text{sort}] \rightarrow [\text{majority}]$

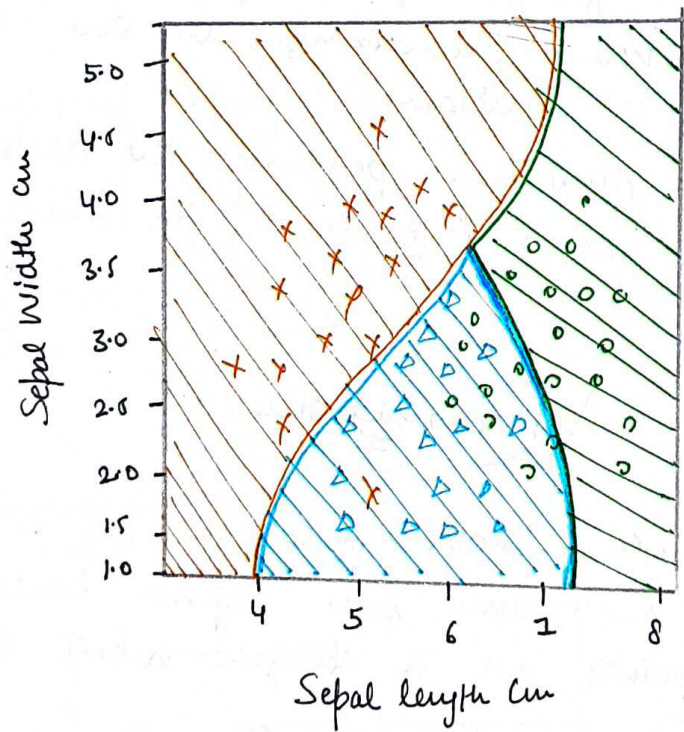
2) Higher dimension data $\rightarrow [f = 500]$

↓
curse of dimension \rightarrow distance concept \rightarrow unreliable

What is Decision Boundary?

⑤

In a classification problem with two or more classes, a decision boundary or decision surface is a hypersurface that partitions the underlying vector space into two or more sets, one for each class. The classifier will classify all the points on one side of the decision boundary as belonging to one class and all those on the other side as belonging to the other class.



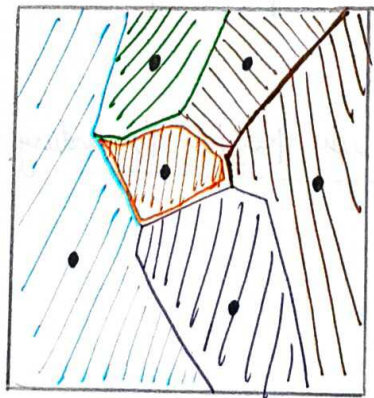
Decision Region Boundary: Sepal length and Sepal width

Important Points

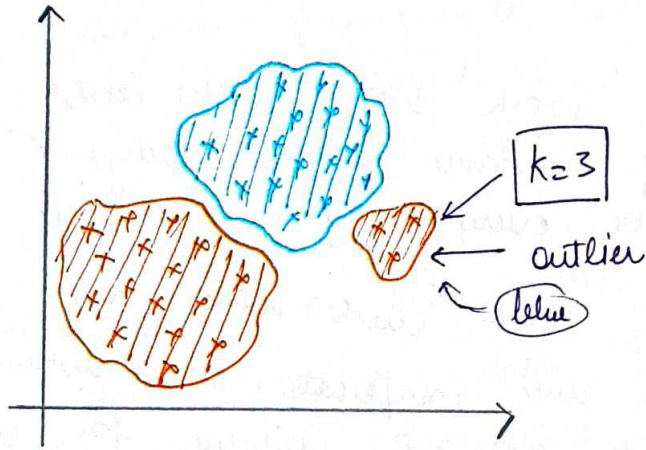
1. We can draw decision boundary for all the classification algorithms including Neural Networks
2. Decision boundary can be both linear (as in the case of SVM) or non-linear (as in the case of Decision tree classifier or knn)
3. Decision boundaries are not always clear cut. That is, the transition from one class in the feature space to another is not discontinuous, but gradual. This effect is common in fuzzy logic based classification algorithms, where membership in one class or another is ambiguous.
4. For higher dimension problems the decision boundary acts as a hyperplane (for linear ones)

Voronoi Diagram

In mathematics, a Voronoi diagram is a partitioning of a plane into regions based on distance to points in a specific subset of the plane.



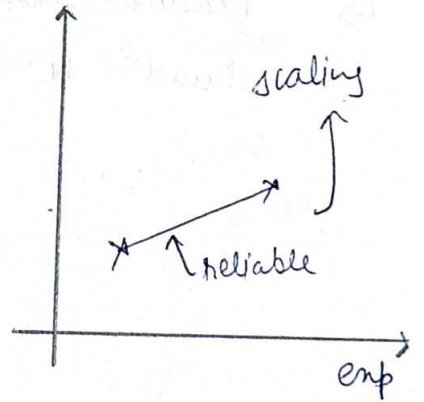
3) Outliers



4) Non-homogeneous scales

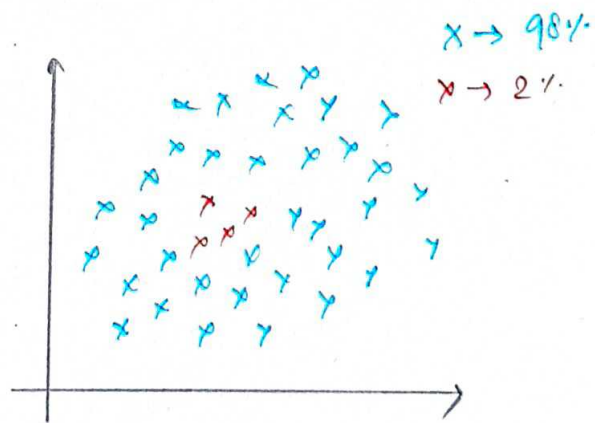
experience	salary	fine
0-25	20k-16	

dominant because of scale



5) Imbalance dataset

Yes → 98%
 No → 2%
 → biased



6) Inference (fail KNN or not good KNN for Inference)

↳ black box model

- 7) Does not work well with large dataset as calculating distance between each data instance would be very costly.
- 8) Does not work well with high dimensionality as this will complicate the distance calculating process to calculate distance for each dimension.
- 9) Sensitive to noisy and missing data
- 10) Feature Scaling - Data in all the dimension should be scaled properly.

3. Steps to plot Decision Boundary for KNN

⑥

(Assume 2 input cols)

1. Train the classifier on the training set.
2. Create a uniform grid (with the help of Numpy Meshgrid) of points that densely cover the region of input space containing the training set.
3. Classify each point on the grid. Store the result in an array A , where A_{ij} contains the predicted class for the point at row i , column j on the grid.
4. Plot the array as an image, where each pixel corresponds to a grid point and its color represents the predicted class. The decision boundary can be seen as contours where the image changes color.
5. Finally plot the training data with their respective color on the same contours.