

# Lecture : 09

" You are the average of the five people you spend most time with "

K-Nearest Neighbor (KNN)

→ It's supervised learning - classification Algorithm / Regression

Bias: the inability of model to truly capture the relationship in the training data.

(diff b/w our actual and predicted values)

Variance: It measures how much the prediction of a model vary for different training datasets.

(model's sensitivity to fluctuations in the data)

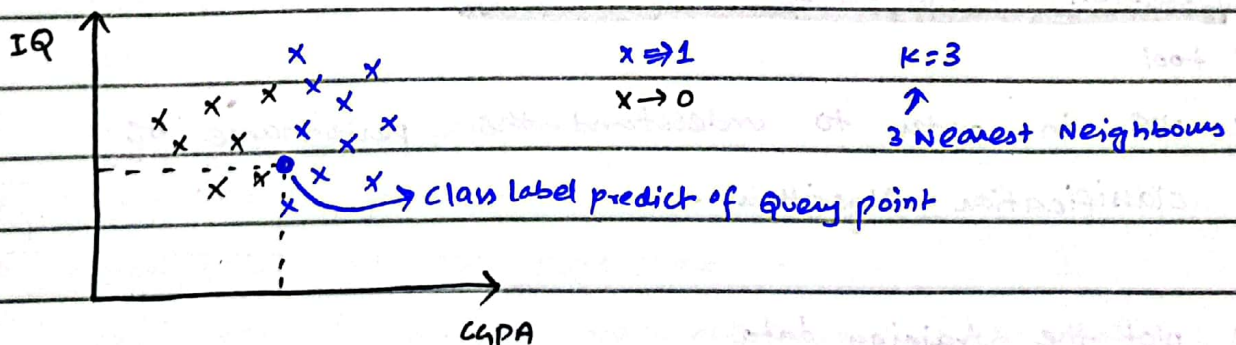
Overfitting: training error low and test error high

Underfitting: training error high

★ (Target would be low bias, low variance)

## Example

→ Let suppose we have CGPA and IQ of students and we have to predict about their placement selection.



→ calculate euclidean distances  $d(x,y) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

→ sort distances in ascending order

→ majority count



## How to select K?

### ① heuristic approach

$\sqrt{n} \rightarrow$  no of observation

$\hookrightarrow$  use odd number

(e.g)  $n = 400$

$$\sqrt{400} = 20$$

either 19, 21  
for K

### ② Experimentation approach

(e.g)  $n = 1000$

800  
training

200  
testing

knn  $\rightarrow 1$

knn  $\rightarrow 2$

knn  $\rightarrow 3$

$\vdots$

knn  $\rightarrow 25$

build different knn models for  $n=1, 2, 3, \dots$   
and select the best one.

## Decision Surface / Decision Boundary:

$\rightarrow$  tool

$\rightarrow$  use in order to understand the performance of  
classification Algorithm.

① plot the training data

② plot range of training data on x and y-axis

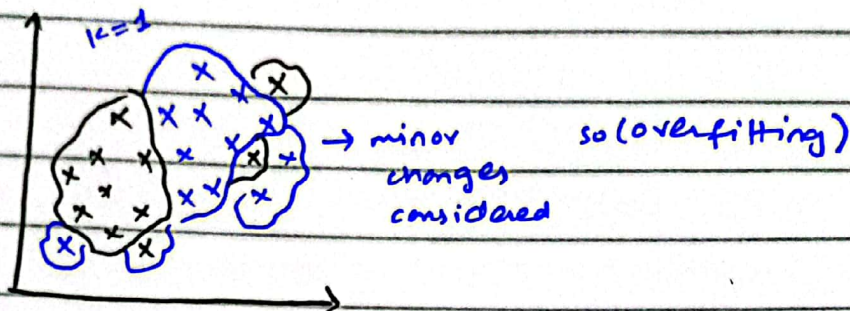
③ for the range generate a numpy meshgrid



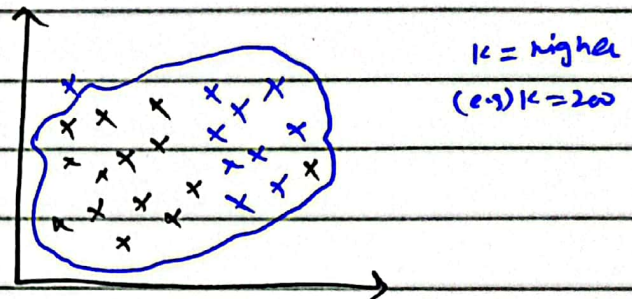
## overfitting and underfitting in kNN:

overfitting  $\rightarrow$  high variance ( $k \rightarrow$  low values)

underfitting  $\rightarrow k \rightarrow$  high values



(majority count)  
 $\rightarrow$  for a new query point  
majority count is considered  
so (underfitting)



## Limitations of kNN:

$\rightarrow$  Large datasets

$\hookrightarrow$  kNN is a lazy learning technique.  $\rightarrow$  Nothing happens in training phase  
 $\rightarrow$  prediction phase is slow

$\rightarrow$  high dimensional data

• if dataset has high dimensions (distance concept is not reliable)

$\rightarrow$  outliers

$\rightarrow$  Non-homogenous features scale

$\rightarrow$  Imbalanced dataset

$\rightarrow$  Inference <sup>and</sup> not for prediction

It does not provide inference (means which feature has high/low impact)

## kNN for Regression:

$\rightarrow$  Average of neighbours values