

$k=3$

$x_q = (2B, 1R) \Rightarrow \text{Blue}$

$x_{q1} = (2B, 3R) \Rightarrow \text{Red}$

Classification
↓
voting

Regression
└ Average

Can I take even number of neighbours?

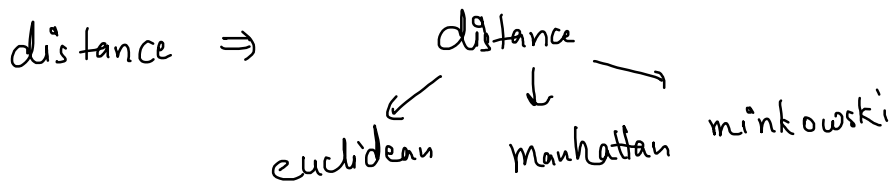
$k=4$ $x_{q2} = (2B, 2R) \Rightarrow$

$P(B) = \frac{1}{2}$

$P(R) = \frac{1}{2}$

Hyperparameters

$k \Rightarrow$ number of neighbours $\Rightarrow k \uparrow$ or $k \downarrow \Rightarrow$ helps you in improving the model.



effects of k

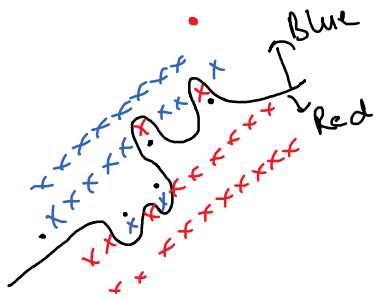
$\Rightarrow k=1$

blue

\rightarrow Non-smooth decision boundary

\rightarrow n/a mistakes

⇒ $K=1$



→ No mistakes

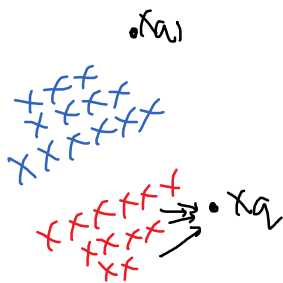
→ Perfect boundary

→ Training accuracy ↑

→ Testing accuracy ↓

↑ } overfitting
 ↓ } low bias & high variance
 ↓
 explore?

2) $K=N$



$x_q = (13B, 12R) \Rightarrow \text{Blue}$

$x_{q1} = (13B, 12R) = \text{Blue}$

$x_{q3} = (13B, 12R) = \text{Blue}$

Irrespective of x_q position, the answer will always be blue.

→ Training accuracy ↓

→ Testing accuracy ↓

↓ } underfitting

3) $K=5$



$x_{q1} = (4B, 1R) = \text{Blue}$

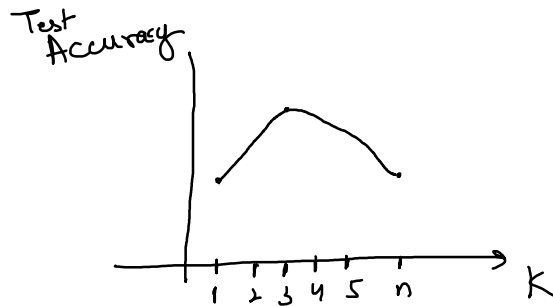
$x_{q2} = (4R, 1B) = \text{Red}$

Training accuracy is like (80% - 90%)

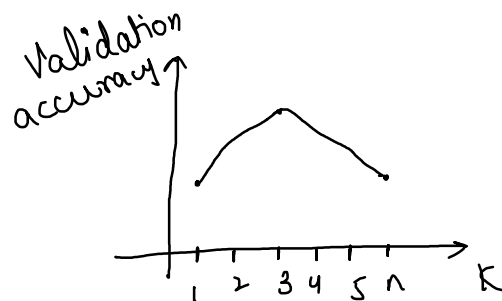
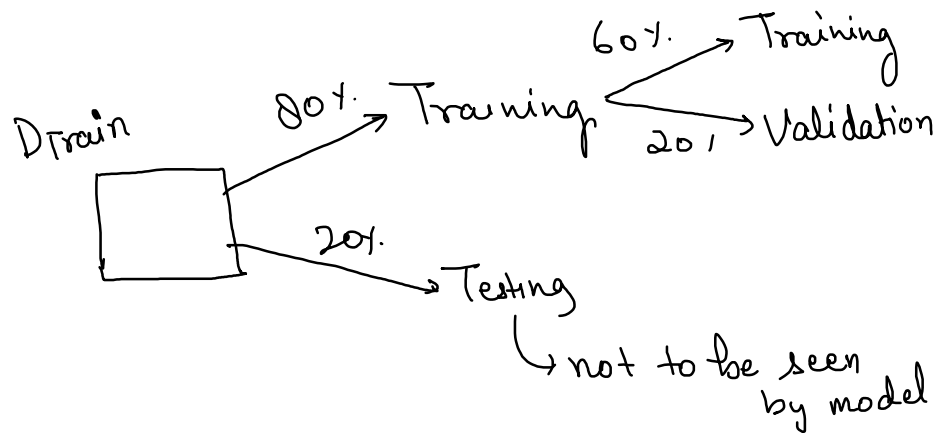


→ Training accuracy is fine (80% - 90%)
 → Test accuracy is fine (80%)
 ↓
 Right fit
 or
 best fit

curve of k with accuracy:



choosing the right value of $k \Rightarrow$ Cross-validation



CROSS-VALIDATION \Rightarrow K-fold \Rightarrow k'-fold Thumb rule \Rightarrow $K'=10$ (industry)

$K' \Rightarrow$ how many partitions will be done of your dataset

$K'=4$

D_{train}

D_1	D_2	D_3	D_4
-------	-------	-------	-------

(# neighbours) K	Training	Validation	Accuracy
1	$D_1 \ D_2 \ D_3$	D_4	a'_1
1	$D_2 \ D_3 \ D_4$	D_1	a'_2
1	$D_1 \ D_2 \ D_4$	D_3	a'_3
1	$D_1 \ D_3 \ D_4$	D_2	a'_4
<hr/>			
2	$D_1 \ D_2 \ D_3$	D_4	a^2_1
2	$D_2 \ D_3 \ D_4$	D_1	a^2_2
2	$D_1 \ D_2 \ D_4$	D_3	a^2_3
2	$D_1 \ D_3 \ D_4$	D_2	a^2_4

validation accuracy list = $[a'_{avg}, a^2_{avg}, a^3_{avg}, a^4_{avg}, a^5_{avg}, a^6_{avg}, a^7_{avg}]$

$\begin{matrix} K=1 \\ \uparrow \\ a'_{avg} \end{matrix}$
 $\begin{matrix} K=2 \\ \uparrow \\ a^2_{avg} \end{matrix}$
 $\begin{matrix} K=3 \\ \uparrow \\ a^3_{avg} \end{matrix}$
 $\begin{matrix} K=4 \\ \uparrow \\ a^4_{avg} \end{matrix}$
 $\begin{matrix} K=5 \\ \uparrow \\ a^5_{avg} \end{matrix}$
 $\begin{matrix} K=6 \\ \uparrow \\ a^6_{avg} \end{matrix}$
 $\begin{matrix} K=7 \\ \uparrow \\ a^7_{avg} \end{matrix}$

from the list, choose value k that has highest avg validation accuracy

Advantages:

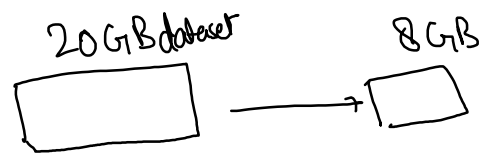
Disadvantages:

Advantages:

- It is very easy to understand
- No assumptions in the algorithm

Disadvantages:

- Lazy learner
↳ does all calculations at time of execution



- space issues
- very slow algorithm
- Imbalanced dataset

Application: Healthcare

Evaluation Metrics (CLASSIFICATION)

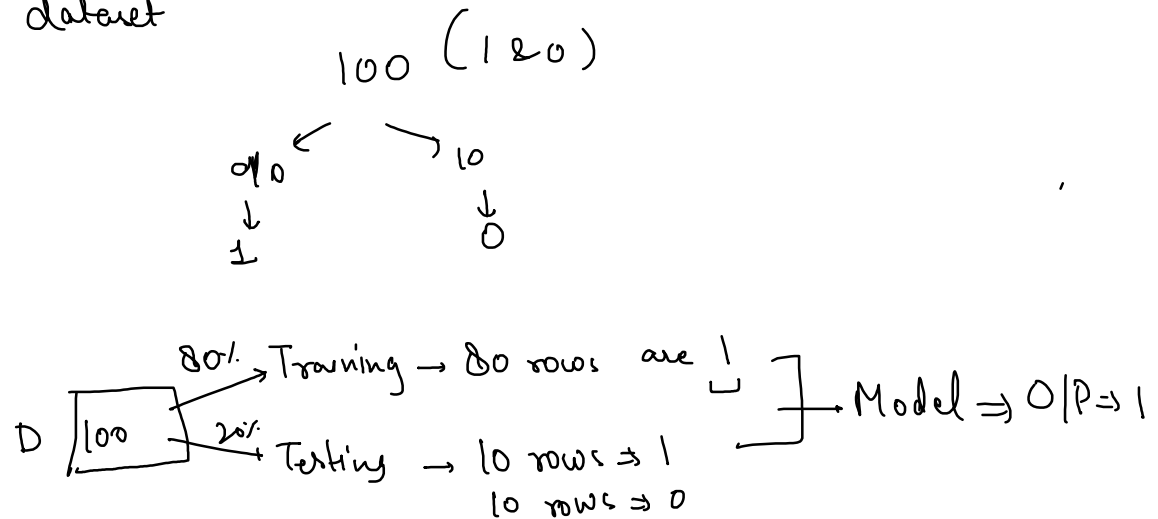
CONFUSION MATRIX

	Actual Values		
	0	1	
Predicted values	0	FN	⇒ Total predicted 0 (-ves)
	1	TP	⇒ Total predicted 1 (+ve)
	↓ Total actual -ve (0)	↓ Total actual positive (1)	

$$1) \text{ Accuracy} \Rightarrow \frac{TP + TN}{TP + FP + TN + FN}$$

Accuracy isn't reliable under two situations:

→ Imbalanced dataset



⇒ Prob ⇒ 0.5 (complete confusion)

$$\rightarrow \downarrow \text{Precision} \Rightarrow \frac{TP}{TP + FP \uparrow}$$

→ Recall ↑
(Sensitivity)
(TPR)

$$\frac{TP}{TP + FN \downarrow}$$

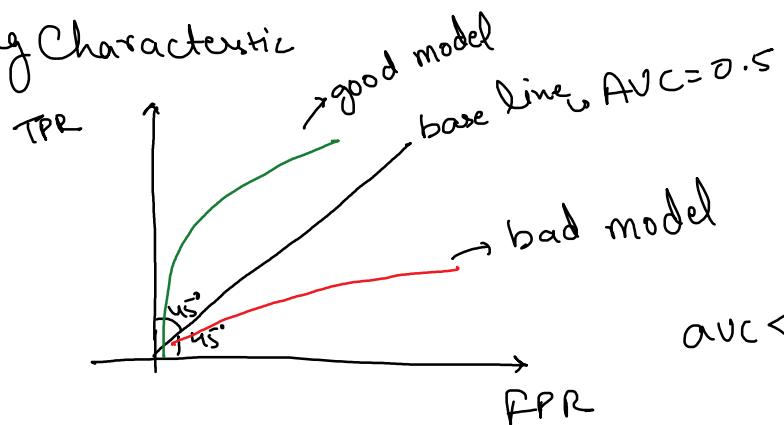
$$\rightarrow FPR = \frac{FP}{TN + FP}$$

$$\rightarrow \text{f1-score} \Rightarrow \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\rightarrow \text{specificity} = 1 - \text{FPR} = 1 - \frac{\text{FP}}{\text{TN} + \text{FP}} = \frac{\text{TN} + \text{TP} - \text{FP}}{\text{TN} + \text{FP}} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

ROC - AUC curve

Receiver Operating Characteristic



$\text{auc} < 0.5$
bad model

$\text{auc} > 0.7$
good model