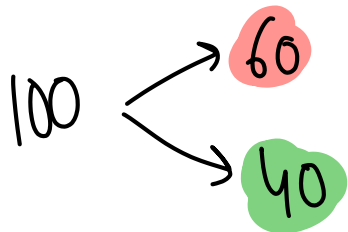


KNN-2

KNN

## Bias - Variance Tradeoff

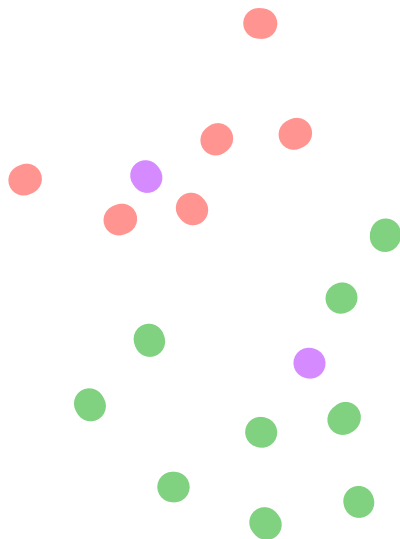


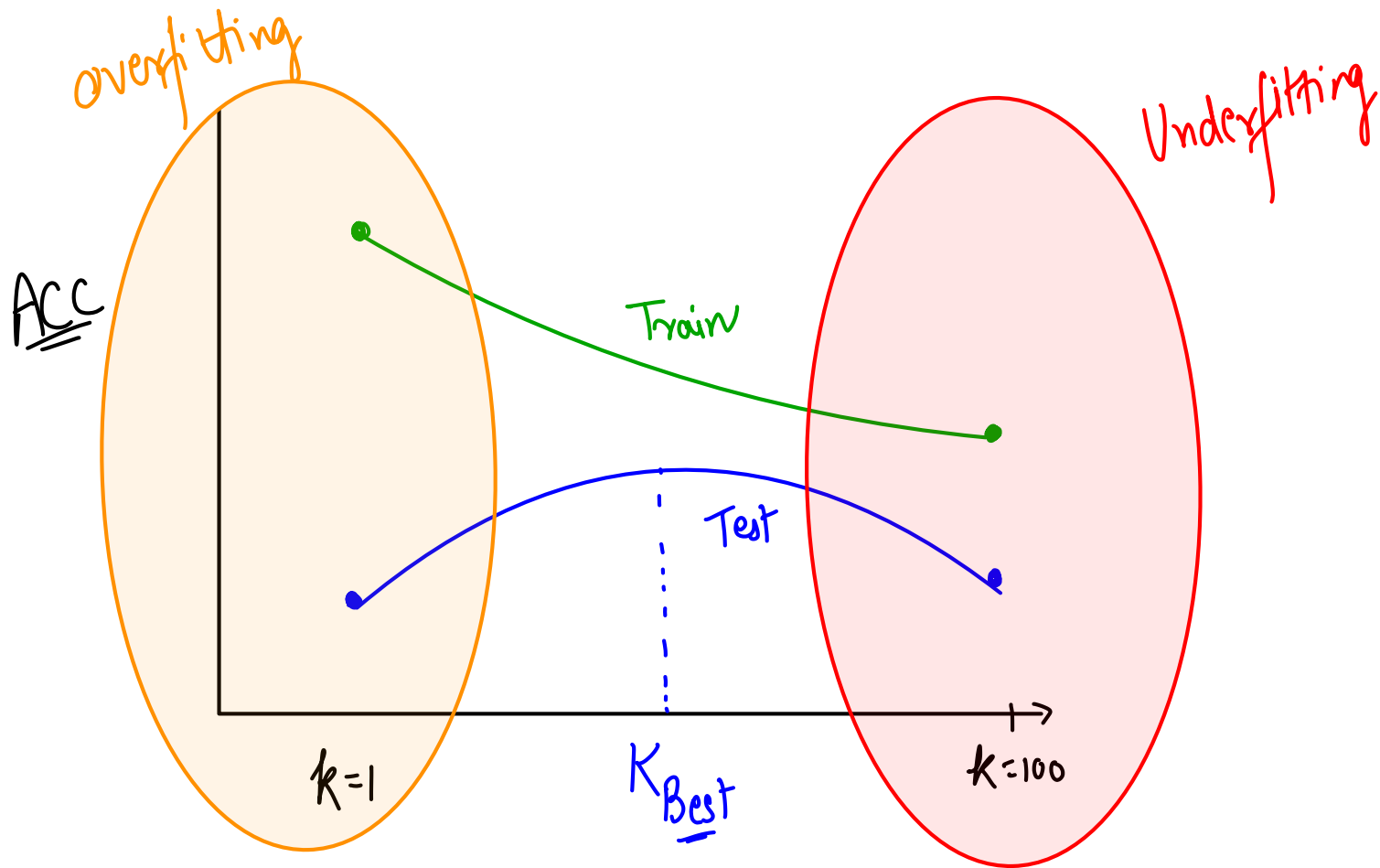
$k = 1$

Overfitting  
high variance

$k = 100$

Underfitting  
high bias





# Impact of Outlier

$k=1$   
low

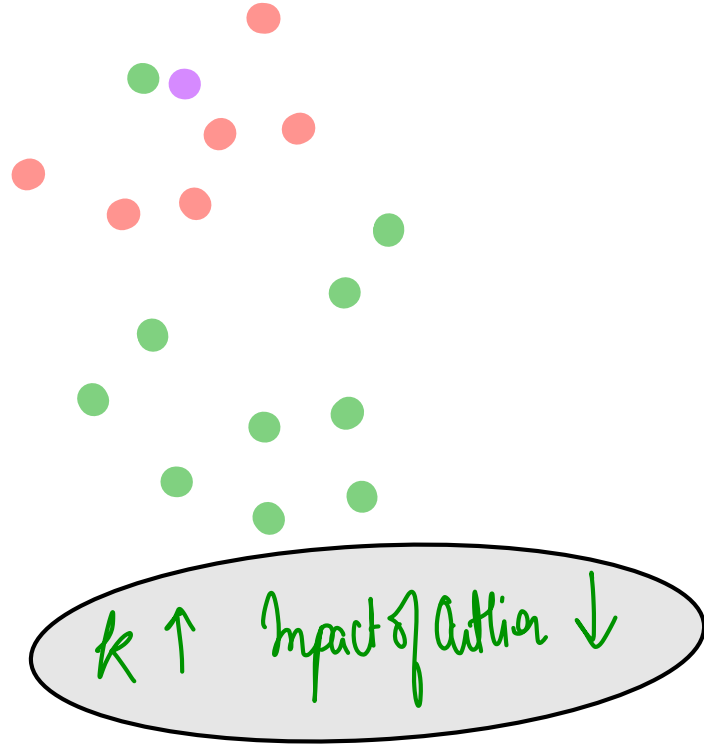
Impact of Outlier  $\uparrow\uparrow\uparrow$

$k=7$   
moderate


Impact of Outlier  $\uparrow$

$k=100$   
high

Impact of Outlier  $\downarrow\downarrow\downarrow$



## Quiz time!

 Quiz Ended!

**what to say if model has high varaince and low bias ?**

44 users have participated



**A**

**Model overfits**

80%

**B**

**Model underfits**

20%

# Quiz time!

🕒 Quiz Ended!

**k → hyperparameter, then what data to use for hyperparameter tuning ?**

34 users have participated



A

validation

82%

B

training

6%

C

test


3%

D

entire data

9%

## Quiz time!

 Quiz Ended!

**k-NN algorithm does more computation on test time rather than train time.**

37 users have participated



A

TRUE

65%

B

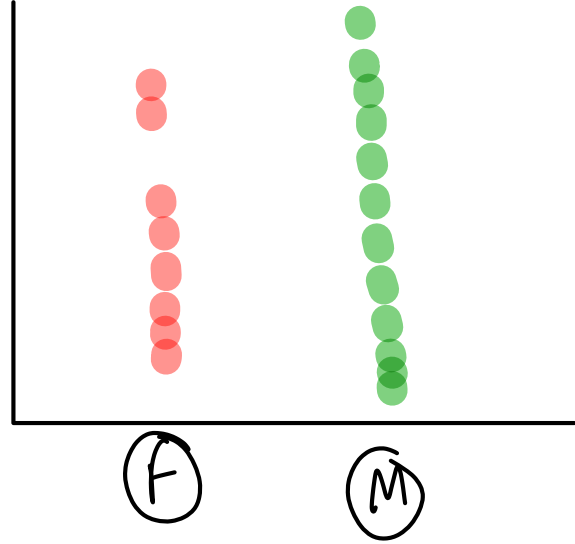
FALSE

35%

# KNN for Categorical features

gender	Blood group	Risk of high driving
		2

Age



Solution  $\Rightarrow$  Encoding

→	<u>One</u>	→	×
→	label	→	×
→	Target	→	✓



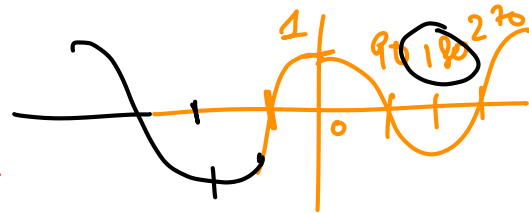
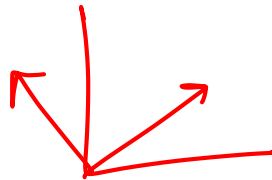
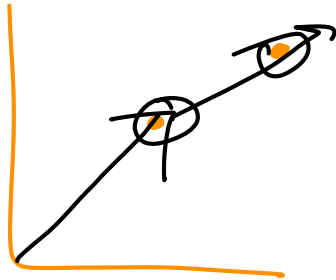
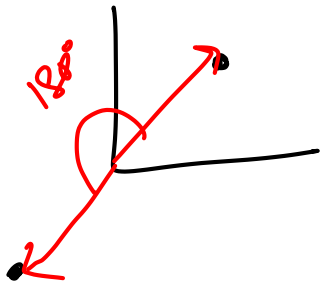
1000  $\rightarrow$  1000 Columns  $\Rightarrow$  1000 new features  
 $\rightarrow$  " " dimensions.

\* Euclidean will fail

Cosine Distance  $\Rightarrow$

$$\cos \theta = \frac{x_1 \cdot x_2}{\|x_1\| \|x_2\|}$$

$\begin{bmatrix} -1 & 1 \end{bmatrix}$

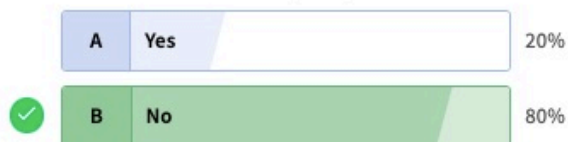


## Quiz time!

⌚ Time Left: 13s

**quiz (what do you think) if One Hot Encoding increases data dimension to ( $d=1000$ ), will Eculidean Distance work ?**

35 users have participated



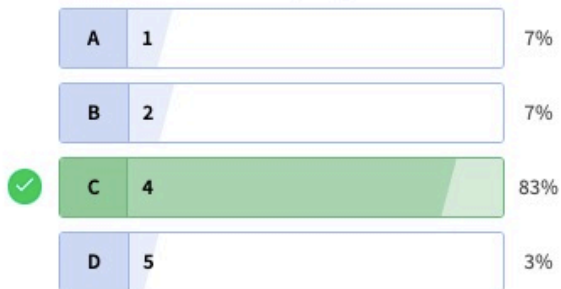
[End Quiz Now](#)

## Quiz time!

🕒 Quiz Ended!

Which of the following will be Manhattan Distance between the two data point A(1,1,3) and B(1,3,5)?

30 users have participated



Google KNN

$[010]$   
 $[010110]$  Qutub Minar ①

$[010]$   
 $[00001]$  Red fort ②

$[010]$   
 $[1101010]$  Rashtrapati Bhawan, ③

Marine Drive, Gateway of India.

$[010101]$   
 $[111]$  ④

$[101010]$   
 $[111]$  ⑤

key	value
delhi	1, 2, 3
Mumbai	4, 5 -

Hashing Table { }

LSH  $\Rightarrow$  Locality Sensitive Hashing.

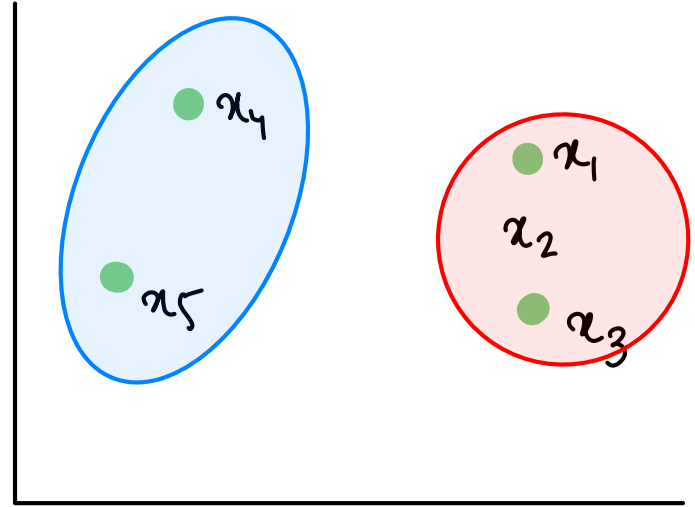
# Missing data

## KNN BASED IMPUTATION.

median, modu, mean, fixed value.


Imputation: Finding the best possible guess.

	$f_1$	$f_2$	$f_3$	$f_4$	$f_5$
$x_1$					
$x_2$					
$x_3$					
$x_4$					
$x_5$					



$$\text{Purple scribble} = \frac{+ \text{Blue dot}}{2}$$

# Quiz time!

 Quiz Ended!

Select the true statements


s1- kNN is less time intensive when LSH is used

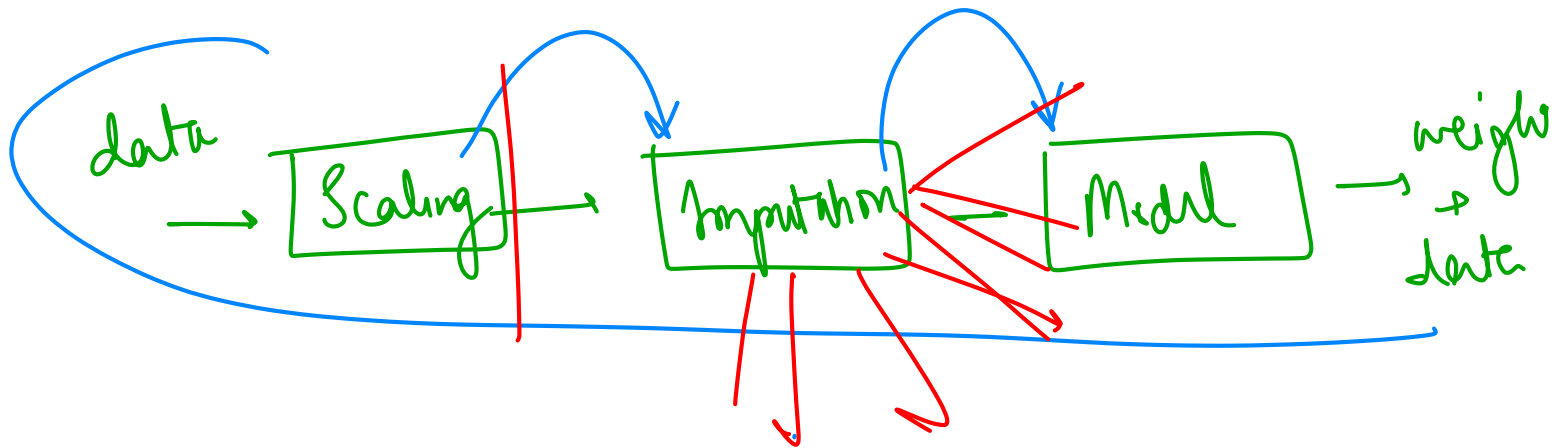
s2- k must be odd

s3- kNN used for imputing

s4- For high dimension, euclidean not used

35 users have participated

A	s1	9%
B	s2	3%
C	s3	6%
	D all of the above	83%



$\text{Scaler} = \text{StandardScaler}()$   
 $X_{\text{sm}} = \text{Scaler.fit\_transform}(X)$   
 $X_{\text{final}} = \text{KNNImputer}(X_{\text{sm}}, k=5)$   
 $\text{KNN neighbor}(\text{un}(X_{\text{final}}))$