

## Missing Values

How to fill the missing values

Statistics used :-

⇒ median

mode [categorical]

⇒ mean

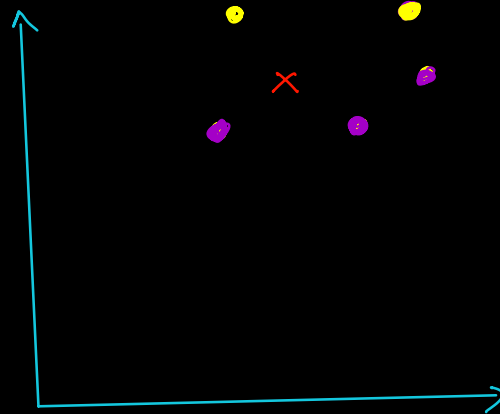
Distance Based :-

⇒ K-nn values

## K-nn Based Imputation

Euclidean or manhattan

	$a_1$	$a_2$	$a_3$	$a_4$
	$(a_1, u_1)$			
$u_1$	( )		( )	( )
$u_2$	( )	( )	( )	( )
$u_3$				
$u_4$				



$$\text{Euclidean} = \sqrt{(y_2 - y_1)^2 + (x_2 - x_1)^2}$$

$$\text{Manhattan} = |x_2 - x_1| + |y_2 - y_1|$$

⇒ mean for numerical  
mode for categorical

⇒ K should be odd for categorical

## Framework

⇒ Create a small subset of data with complete observation.

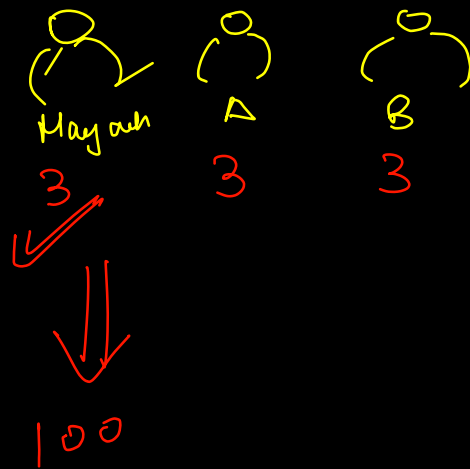
⇒ Delete some values manually

⇒ Use multiple method to fill.

⇒ See where are they failing / error calculation

⇒ Use the best method

	Person	Gender	Company	Prime Customer
	1) 350	—	—	○
	2) 200	—	✓—	○
⇒	3) NaN	—	—	○
	4) 200	—	✓	○
⇒	5) NaN	—	—	○

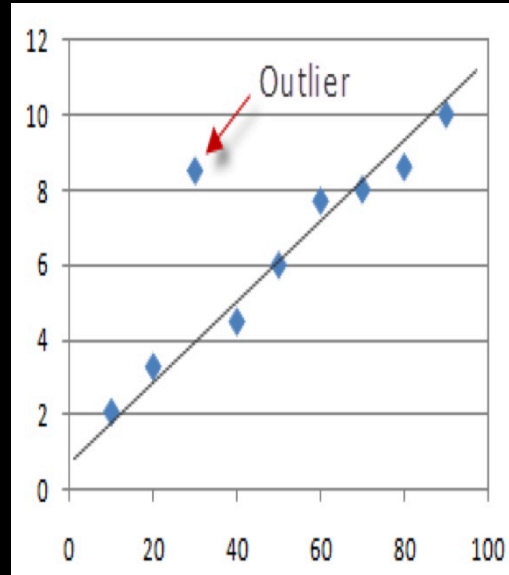


100 lectures

$$(\text{original value} - \text{calculated value})^2$$

# Outlier Analysis

Observation which is inconsistent with rest of dataset.

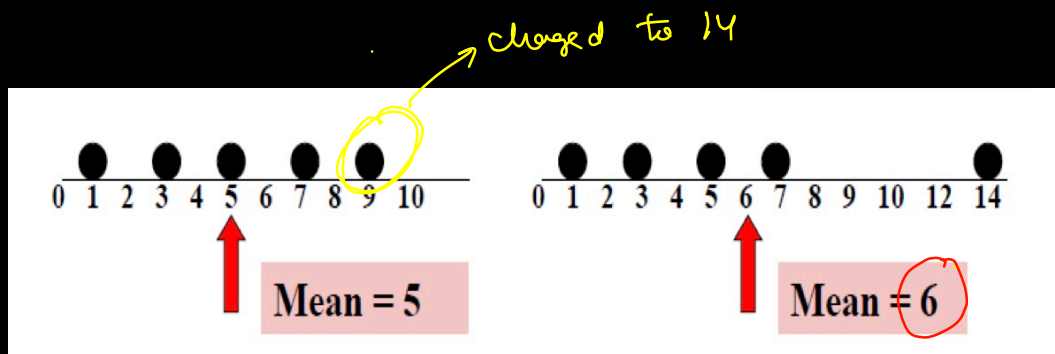


Cause of Outliers

⇒ Poor data quality / Contamination

⇒ Low quality / malfunctioning equipment

⇒ Correct but exceptional data



IQR, Box plot

Upper fence

Lower fence

## Applications

- ⇒ Fraud Detection
- ⇒ Medical Treatment
- ⇒ Income / buying habits

## Variable Importance / Feature Selection

⇒ subset of relevant features to be used in model training / construction

Real life

Cricket

Rohit

Kohli

Sachin

Bowling avg

8

7

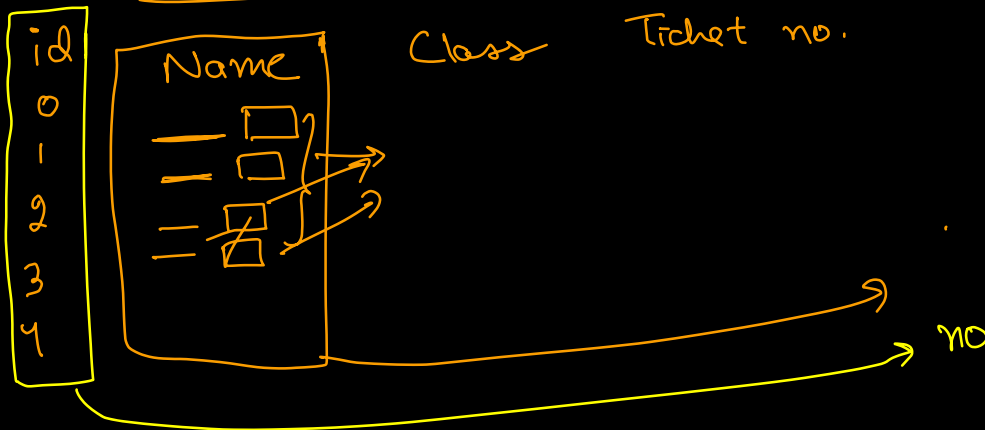
9



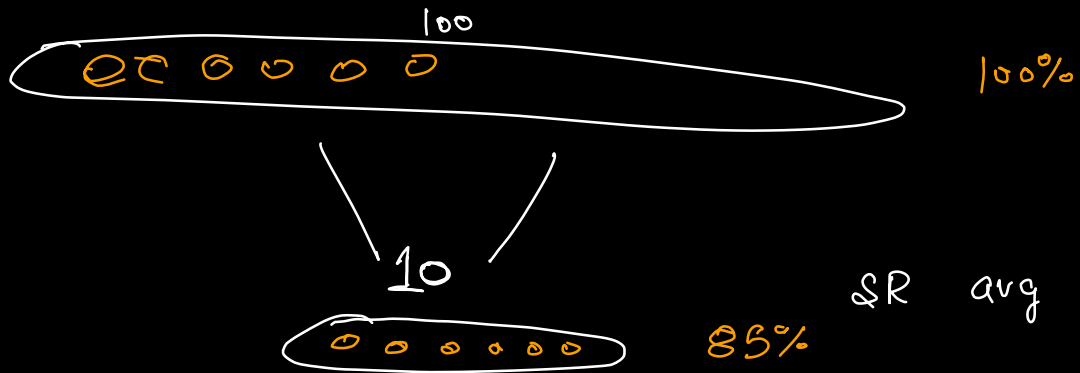
minutes played

seconds played  
≤ 60

## Tikanic Dataset



# Dimensionality Reduction



non-linear model  
speed

## Methods

Filter  
Wrapper  
Embedded

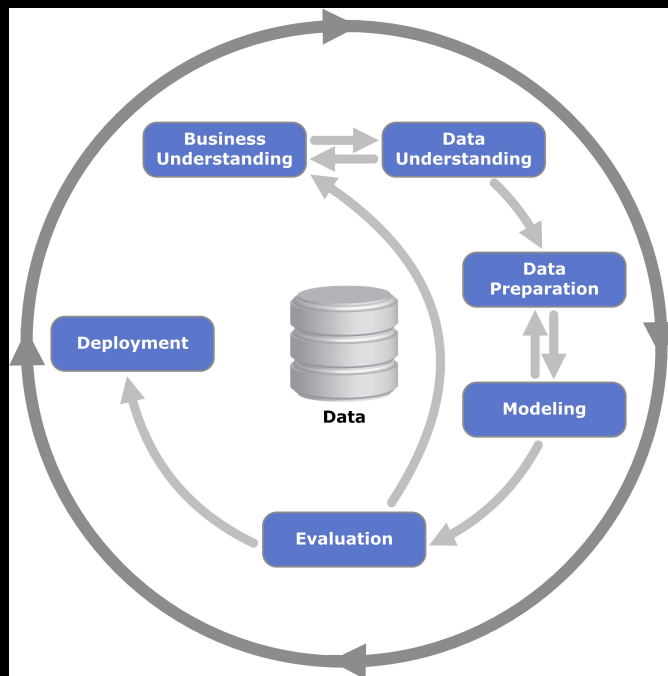
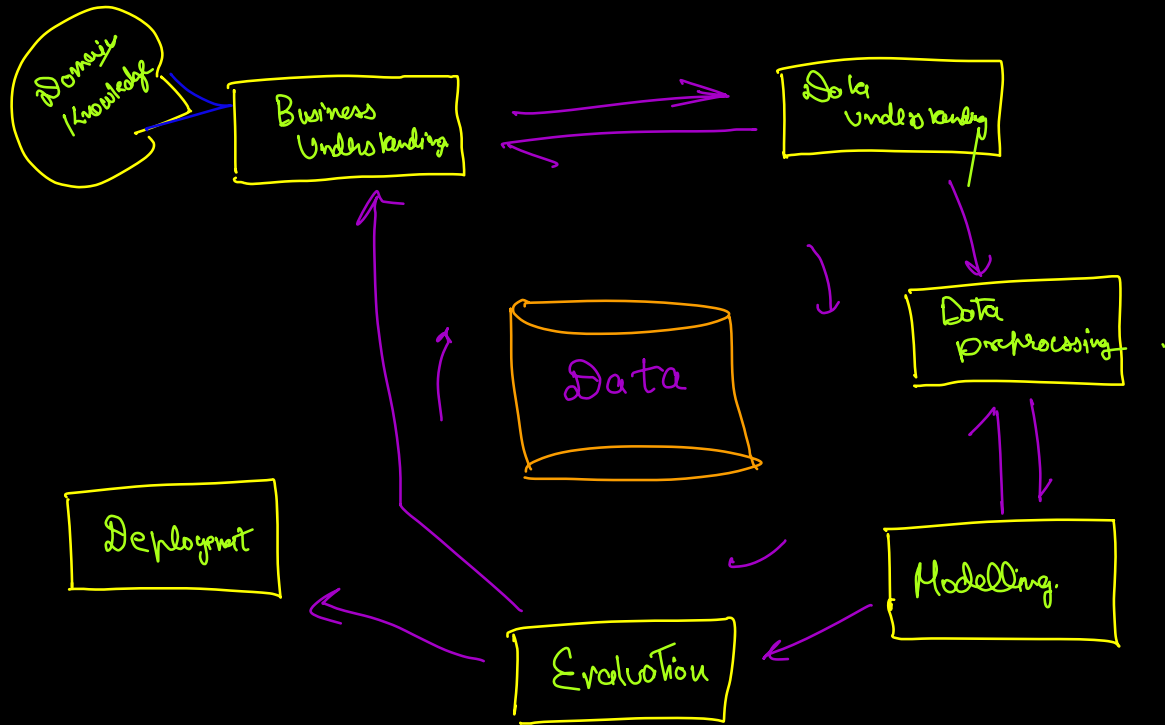
## Dimensionality Reduction

PCA  
RF  
Co-relation  
SVD

## Normalization / Feature Scaling



# CRISP - DM



# Data Science Libraries

Numpy  $\Rightarrow$  numerical python

Fast computations on n-dimensional arrays.

- Based on only one data structure  
ndarray

np. Command

array like  $\Rightarrow$  super fast

list in python but they are very slow to process.

50x faster than  
your traditional  
python list

Why should we use numpy

$\left. \begin{array}{l} \text{list} \\ \text{tuple} \\ \text{dict} \\ \text{set} \end{array} \right\} \Rightarrow \text{heterogeneous}$

● np array are homogeneous

Create a numpy array

np. array (seq)

$\Rightarrow$  creates an array

Scaler  $\Rightarrow 10, 20$

Vector  $[10, 20, 30 \dots]$

Matrix  $\left[ \begin{array}{c} \\ \end{array} \right]$

Why numpy array faster than list

1. Locality of reference  $\rightarrow$  things are stored  
in numpy array in continuous blocks  
of memory



## Creation of NumPy Arrays

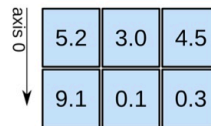
- `Np.array (Sequence)`

1D array



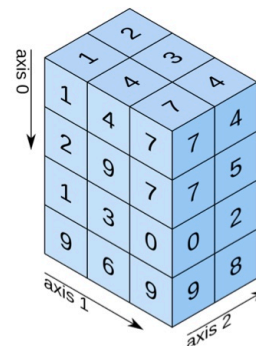
shape: (4,)

2D array



shape: (2, 3)

3D array



shape: (4, 3, 2)