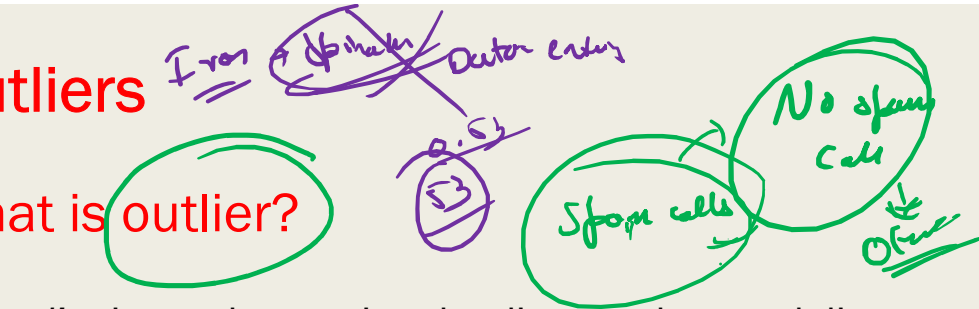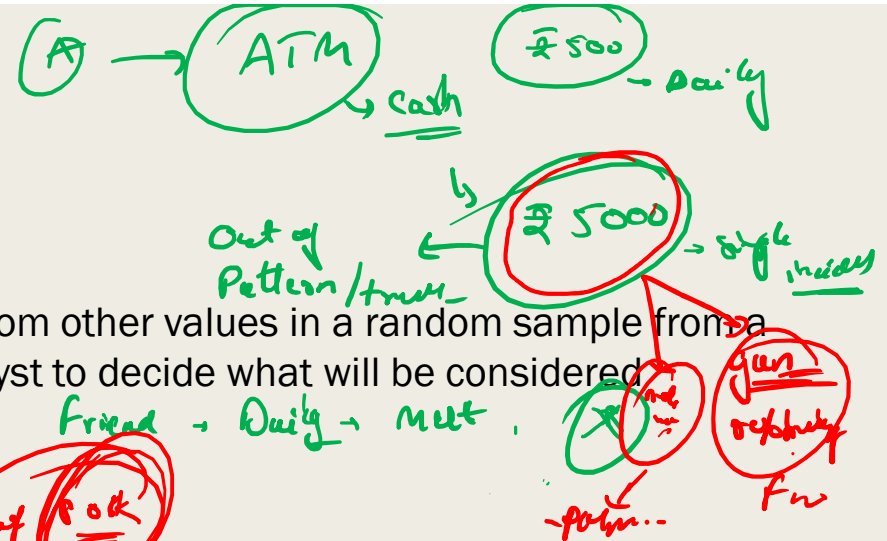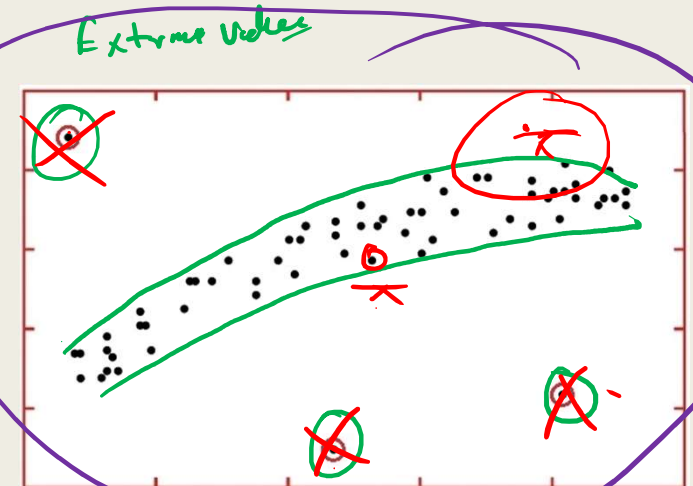# Outliers

## What is outlier?

An *outlier* is an observation that lies an abnormal distance from other values in a random sample from a population. In a sense, this definition leaves it up to the analyst to decide what will be considered abnormal.

## Common Causes of Outliers

1. Data entry errors (human errors)
2. Measurement errors (instrument errors)
3. Experimental errors (data extraction or experiment planning/executing errors)
4. Intentional (dummy outliers made to test detection methods)
5. Data processing errors (data manipulation or data set unintended mutations)
6. Sampling errors (extracting or mixing data from wrong or various sources)
7. Natural (not an error, novelties in data)

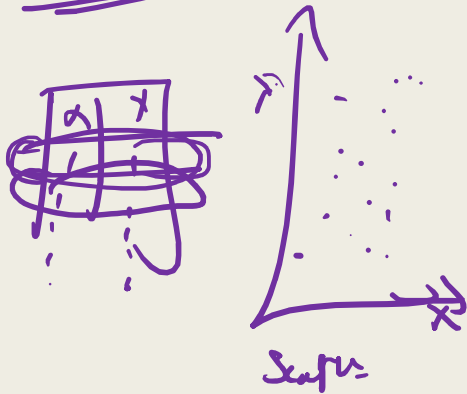# Common methods of determining an Outlier

500, 500, 500, 500, ( $5000 )

1. Sort the data and see for the extreme values
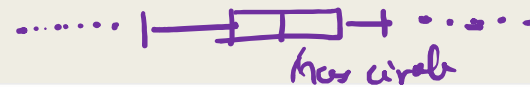2. Plotting - Boxplot, Scatterplot
3. IQR Method
4. - Z - score Method

N.D.

Has circle

No circle

Scatter



With Outliers

Outliers removed
A much better fit!

# Why do we need to treat outliers?

Outliers can impact the results of our analysis and statistical modelling in a drastic way.

# IQR Method

$< LL \rightarrow$ outlier

$> UL \rightarrow$ outlier

## A DATA VALUE IS CONSIDERED TO BE AN OUTLIER IF..

$LL =$

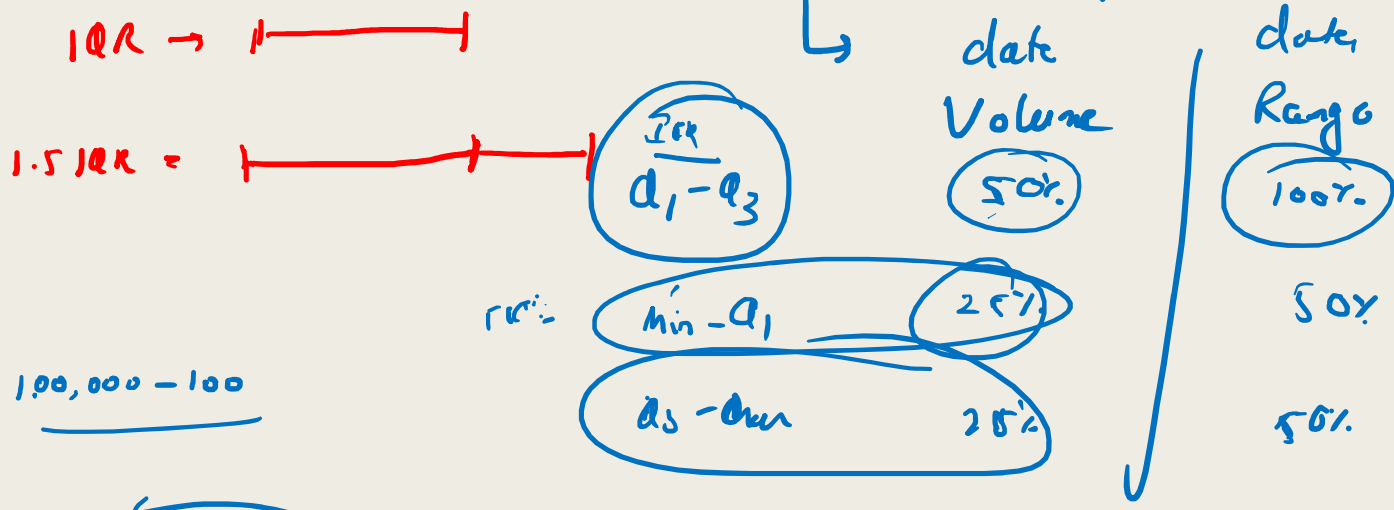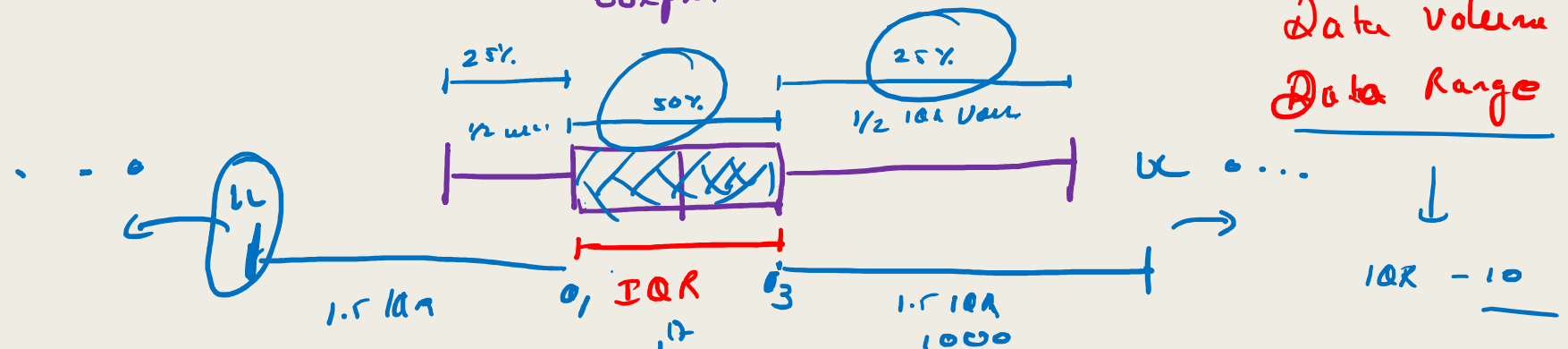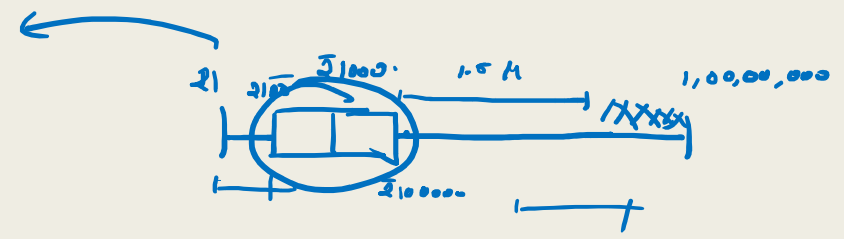DATA VALUE **<** $Q1 - 1.5(IQR)$

OR

$UL =$

DATA VALUE **>** $Q3 + 1.5(IQR)$

Boxplot

Data Volume
Data Range

25% · 50% · 25%
½ IQR Value

IQR — 10

$Q_1$   IQR   $Q_3$

1.5 IQR
1.5 IQR
1000

IQR →

IQR
$Q_1 - Q_3$

data
Volume
50%

data
Range
100%

1.5 IQR =

min $- Q_1$    25%

$Q_3 -$ max    25%

50%

1,00,000 — 100

50%

1.5 times of IQR

1 IQR     ½ IQR

$$1, \textcircled{2}, 3, 4, \textcircled{5}, 100$$

min



LL   $Q_1$   M   $Q_3$   UL

0 → Man

Min.   $Q_1$   Media   $Q_3$   UL

6

.5    6    3.5

−.2.5

$2^n$

Modified Boxplot

Max

●

$IQR = Q_3 - Q_1$
$= 5 - 2 = 3$

1   2   3   4   5   9.5   100

$LL = Q_1 - 1.5 \, IQR$
$= 2 - 1.5 * 3$
$= 2 - 4.5 = \boxed{-2.5}$

$UL = Q_3 + 1.5 \, IQR$
$= 5 + 1.5(3) = 5 + 4.5$
$= \boxed{9.5}$

00

LL

max

Min = 1
$Q_1 = 2$
med = 3.5
$Q_3 = 5$
max = 100

Q. Can you identify the outliers from the below dataset, using the IQR method?

26.0 ℃ , 15.0 ℃ , 20.5 ℃ , 31 ℃ , 350.0 ℃ , 31.0 ℃ , 30.5 ℃

-350 , 15 , 20.5 , 26 , 30.5 , 31 , 31

7

Outlier

$IQR = 31 - 15$
$= 16$

$1.5 * 16 = 24$

$UL = Q_3 + 1.5(IQR)$

$LL = Q_1 - 1.5(IQR)$

$IQR = Q_3 - Q_1$

$UL = 31 + 24 = 55$

$LL = 15 - 24 = -9$

MODIFIED BOXPLOT