1. write the Gaussian Distribution Emperical formula.
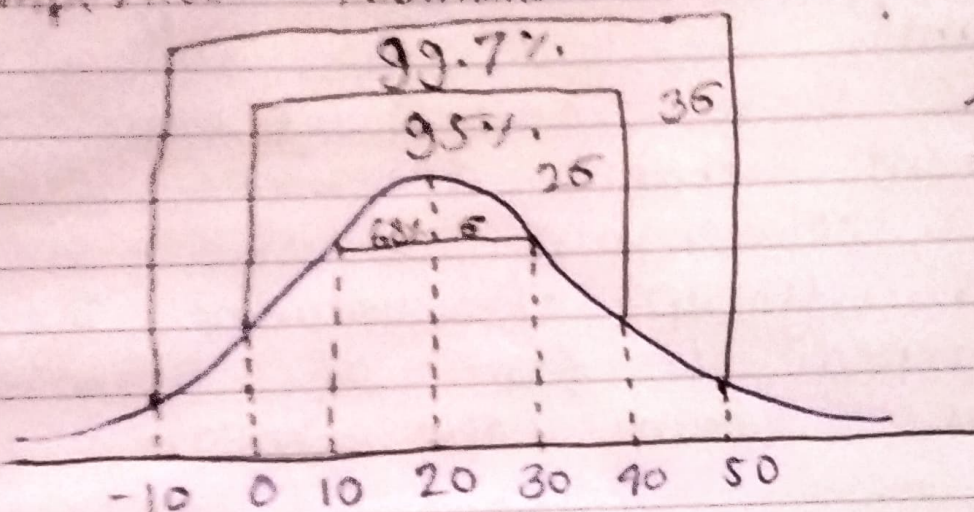


$u = 20$
$\sigma = 10$

If a data is in Gaussian Distribution the first standard deviation $(\sigma)$ will contain 68% of data.

For Second Standard deviation $(2\sigma)$ will contain 95% of all data available.

For Third Standard deviation $(3\sigma)$ will contain 99.7% of all data.

Emperical formula: 68-95-99.7

027-338 • WK-05

2022

27

JANUARY • THURSDAY

2022　　　　　　　　　　　　　JANUARY
M T W T F S S M T W T F S S
1 2 3 4 5 6 7 8 9
10 11 12 13 14 15 16 17 18 19 20 21 22 23
24 25 26 27 28 29 30 31

2. What is Z-Score, and why is it important?

→ A Z-Score is a numerical measurement that describes, a value's relationship to the mean of a group of values. Z-Score is measured in term of standard deviations from the mean.

If a Z-Score is 0, it indicates that the data point's score is identical to the mean.

A Z-Score of 1.0 would indicate a value indicate a value that is one standard deviation from the mean.

Z-Score may be positive or negative, with a positive value indicating the score is above the mean and a negative score indicating it is below the mean.

$$Z = \frac{X_i - \mu}{\sigma}$$

2022    FEBRUARY

M T W T F S S M T W T F S S
1 2 3 4 5 6 7 8 9 10 11 12 13
14 15 16 17 18 19 26 21 22 23 24 25 26 27
28

028-337 • WK-05

JANUARY • FRIDAY

28

2022

3. What is an outlier, Exactly?

-) An outlier is an observation that lies an abnormal distance from other values in a random sample from a population. In a sence, this definition leaves it up to the analyst to decide what will be considered abnormal. Before abnormal observations can be ~~sign~~ singled out, it is necemary to ~~chater~~ characterize normal observations.

We take use five ~~of~~ number Summary for finding the outliers.

Dataset :-

Dataset :- 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27

ii) Calculate the median

Next, we need to calculate the median. The median is the center of the data. If the data set has an odd number of data points, then the mean is centermost number. On the other hand, if the data set has an even number of values, then we will need to take the arithmetic average of the two centermost values. We will calculate the average by adding the two numbers together and then dividing that number by two.

$$Q_2 \text{ (median)} = 10^{th} \text{ position value}$$
$$= 5$$

iii) Calculate upper & lower limit :-

$$Q_1 (25\%) = \frac{Percentile}{100} \times (n+1)$$

$$= \frac{25}{100} \times (19+1)$$

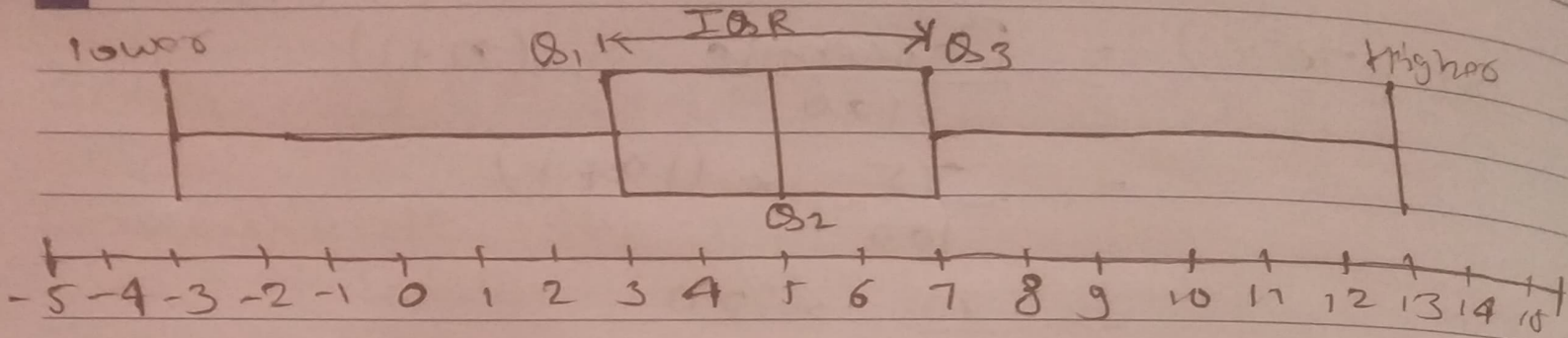$$= \frac{25}{100} \times 20 = 5^{th} \text{ index}$$
$$= 3 \text{ (value)}$$

$$Q_3 \,(70\%) = \frac{\text{Percentile} \times (n+1)}{100}$$

$$= \frac{75}{100} \times (19+1)$$

$$= \frac{75}{100} \times 20$$

$$= 15\text{th} \quad \text{index}$$

$$\Rightarrow 7$$

iv) Calculate the difference:-

$$IQR = Q_3 - Q_1$$

$$= 7 - 3$$
$$= 4$$

Lower fence:- $Q_1 - 1.5 \times IQR$
$$= 3 - 1.5 \times 4$$
$$= 3 - 6$$
$$= -3$$

Higher fence $= Q_3 + 1.5 \times IQR$
$$= 7 + 1.5 \times 4$$
$$= 7 + 6 \qquad = 13$$

lower    Q₁ ←—— IQR ——→ Q₃    higher

-5 -4 -3 -2 -1 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Q₂

If any value goes beyond lower fence and higher fence. it is considered as outliers.
in the given dataset 27 is considered as outliers.

4. What are our options for dealing with outliers in our dataset?

→ One of the most important steps as part of data preprocessing is detecting and treating the outliers as they can negatively affect the statical analysis and the training process of a machine learning algorithm resulting in lower accuracy.

An outlier may occur due to the experimental error / human error. They may indicate an experimental error or heavy skewness in data.

IF the dataset in small, we can just look at it and find out the outliers, but if it is big then we need to have the mathematical and visualization technique. These techniques can be as follows:

5. write the Sample and population variances equation and explain Bessel correction.

→

Sample variance:-

$$\frac{\sum_{i=1}^{n} (x_i - \mu)^2}{n-1}$$

Population variance:-

$$\frac{\sum_{i=1}^{N} (x_i - \mu)^2}{N}$$

Bessel's Correction explains that "the use of $n-1$ instead of $n$ in the formula for Sample variance and Standard variance, where $n$ is the number of data in Sample. This corrects the bias of the estimation of the population variance.