

Variance

Variance measures the dispersion of a set of data points around their mean.

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

$$\text{Sample Variance } s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

~~Sample Variance~~ ~~Population Variance~~

• The reason behind $n-1$ in Sample Variance is in population Variance we will use total N which is 100% of data.

• In Sample Variance we are taking $n-1$, because, we don't have 100% of values.

• In a population of 1000, if u take sample as 100, then to calculate ~~std~~ Sample Variance the denominator is $100-1$ that is 99.

• we have more uncertainty due to incomplete data.

Reason for squaring the values:-

* The reason for squaring the values because, the main aim of Variance is to find the dispersion of data, which is distance, so the distance cannot be negative.

* It amplifies the effect of large differences

* Non-negative values don't cancel out

Standard deviation (Preferred more than variance)

* It means how much ~~to~~ or how far the data has been spreaded with respect to the mean.

Population

* Stand. dev = $\sqrt{\sigma^2}$

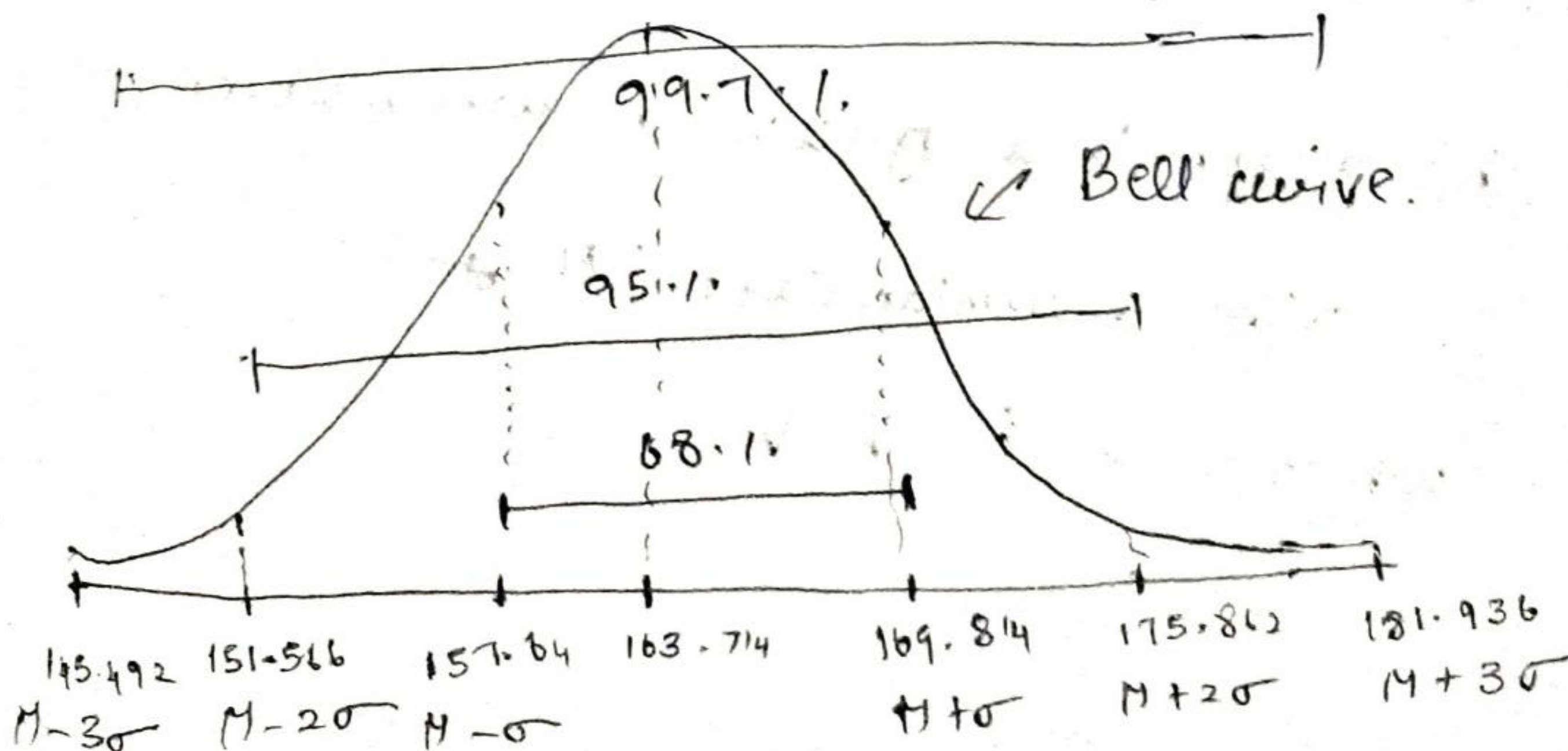
Sample

std. deviation = $\sqrt{s^2}$

$X = \{160, 165, 168, 170, 172, 175, 178\}$

mean = 163.714

$\sigma = 6.074$



• In that graph the data is normally distributed means it forms a bell curve

• The feature which is normally distributed gives better accuracy.

Empirical formulas:-

$$P(\mu - \sigma \leq x \leq \mu + \sigma) = 68.1\%$$

means 68.1% of data distributed between $\mu - \sigma$ and $\mu + \sigma$

$$P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) = 95.4\%$$

means 95.4% of data distributed between $\mu - 2\sigma$ and $\mu + 2\sigma$

$$P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) = 99.7\%$$

means 99.7% of data distributed between $\mu - 3\sigma$ and $\mu + 3\sigma$

• The graph is a gaussian distribution

$$X \sim \text{Gaussian dist}(\mu, \sigma)$$

where x is R.V

Coefficient of Variation:-

- It is also called as relative standard deviation
- $CV = \text{Standard deviation} / \text{mean}$
- Coefficient of Variation (CV) plays a important role in comparing 2 datasets
- comparing standard deviation is meaningless but the comparison of coefficient of variation is meaning ful.

Population formula $CV = \sigma / \mu$

Sample formula $\hat{CV} = s / \bar{x}$

Example:-

Lets take a example of ~~zudata~~ pizza price in 2 different as 2 datasets.

Dataset 1
America

currency	Price
\$	1.00
\$	2.00
\$	3.00
\$	3.00
\$	5.00
\$	6.00
\$	7.00
\$	8.00
	9.00, \$11.00

Dataset 2

~~Mexico~~ Mexico

currency	Price
MXN	18.81
"	37.62
"	56.43
"	56.43
"	94.05
"	112.86
"	131.67
"	150.48
	169.29
	206.91

To calculate

1. Find whether it is sample or population.
2. Find mean
3. Find Sample Variance
4. Find Sample Std. dev.

This dataset contains just ~~10~~ price info of 10 hotels in america and mexico, so it is a sample population.

	Dataset 1	Dataset 2
Mean	5.50	103.46
Sample Variance	10.72	3793.69
Sample std deviation	3.27	61.59

• We cannot compare 3.27 and 61.59 but we can compare its coefficients.

for sample Coefficient of variation	0.60	0.60
--	------	------

• so we obtained same results
• so both the datasets have same
Variability.

• coefficient of Variation does not have a
unit of measurement

• It is universal across datasets.

Percentiles and quantiles:-

heights = {168, 170, 150, 160, 182, 140, 175, 180, 170, 190}

First we need to sort

Sort = {140, 150, 160, 168, 170, 170, 175, 180, 182, 190}

1 2 3 4 5 6 7 8 9 10

5th percentile ~~value~~ = basically means 5th Value.

★ 5th percentile = 170

— 5th percentile means 50% of data points are less than 170 and 50% of data are more than 170.

★ 8th percentile = 180

— It means 80% of data are less than 180 and 20% of data are more than 180.

~~quantiles means 0th 25th 50th 100th percentiles~~

↓

~~4 divisions~~

quantile means 25.1th 50th 75th 100th
1st 2nd 3rd 4th quantiles

Percentile.

Lets consider Amazon delivery time report of the products

$X = \{4d, 5d, 4d, 4.5d, 5d, 5.2d, 5.3d, \dots\}$

Here in this type of cases 95th percentile and 99th percentile plays major role.

* If 95th percentile is 5.6 days, it means 95% of products are delivered within 5.6 days from the time of placing order.

* If 99th percentile is 7 days, it means 99% of products are delivered within 7 days from time of placing order.

95th % = 5.6 days

99th % = 7 days

~~So the supply unit should work on to reduce the time btw 95th and 99th %.~~

~~In 95~~

Considering 95th %, 5% product are not delivered in 5.6 days, it taking extra time, So Supply unit of amazon will work on this to reduce time.

Interquartile range:-

IQR also called as midspread or middle 50%, technically H-spread is the difference between the third quartile and first quartile

$$IQR = Q_3 - Q_1$$

75th - 25th

IQR has breakdown point of 25% due to which it is often preferred over total range.

IQR is used to build box plots

IQR identifies outliers

IQR gives central tendency of data.

∴ dataset of high IQR has more

Variability

∴ dataset of low IQR is preferable.

ex: $X = \{5, 2, 1, 7, 3, 4\}$

1 Step:- Sort it

$$X = \{1, 2, 3, 4, 5, 7\} \quad N=6$$

$$\text{median} = \frac{3+4}{2} = 3.5 = Q_2$$

Q_1 is median of n smallest values

that is Q_1 is median of ~~the~~ first 3 values

$$Q_1 = 2$$

Q_3 is median of n largest values

$$Q_3 = 5$$

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

Median Absolute deviation

~~It uses the~~

It calculates the spread of the data with respect to median.

$$\text{deviation} = |x_i - \text{median}|$$

Since the modulus is used the

value is absolute.