

Agenda:

↳ Recap on Statistics

↳ Pre-processing

Types of data:—

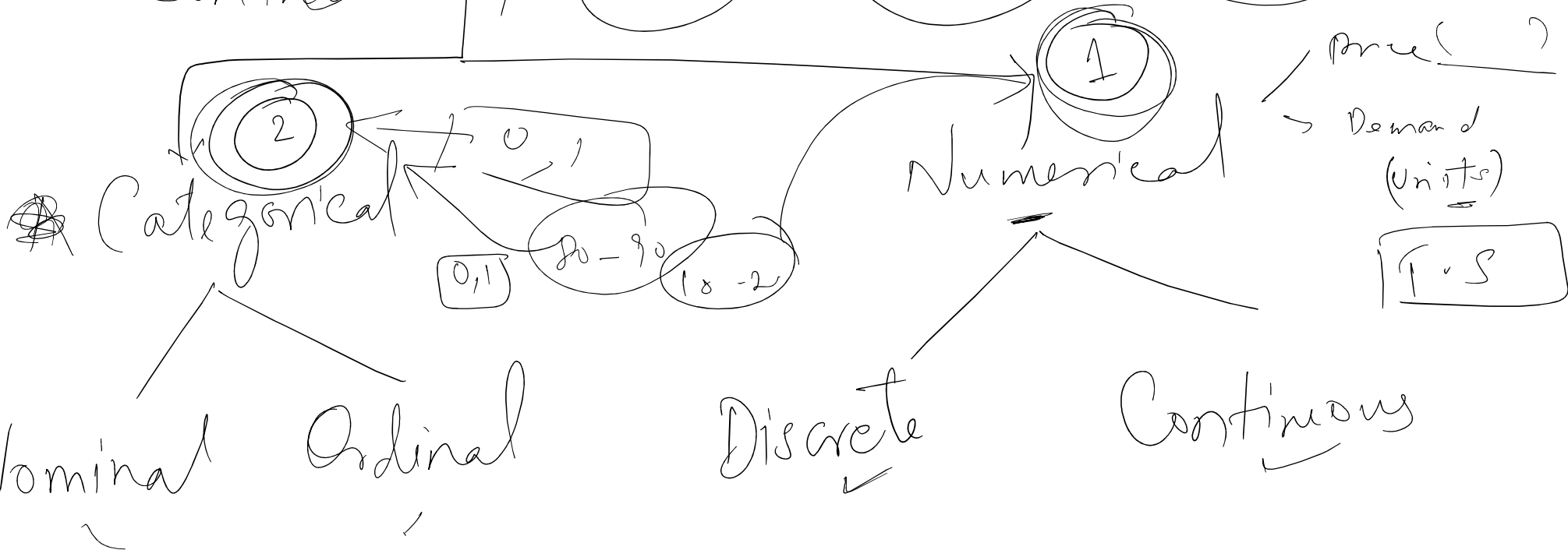
continuous

&

discrete

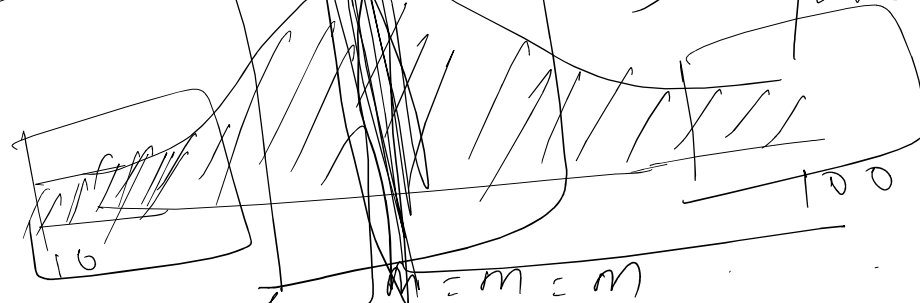
ordinal

nominal



Frequency :-

- # of occurrences over a certain period
- majority of data pts



=> Students

performance

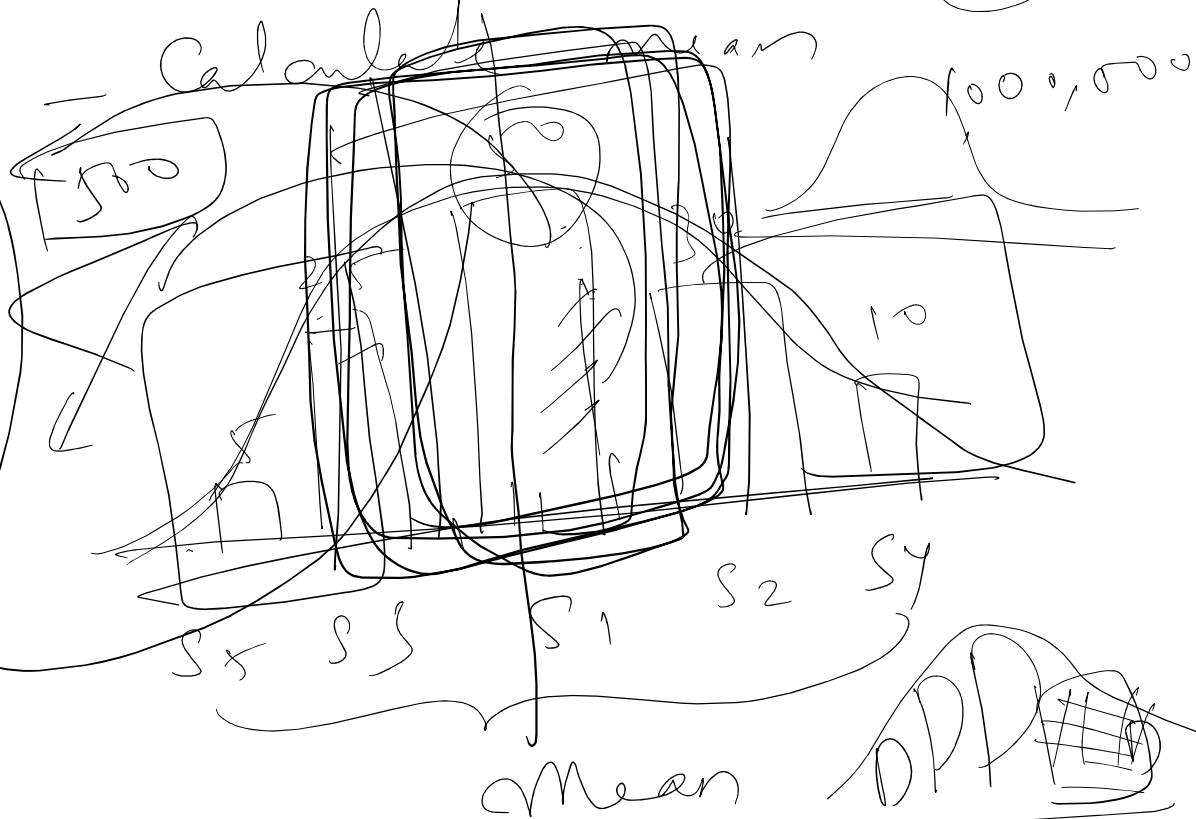
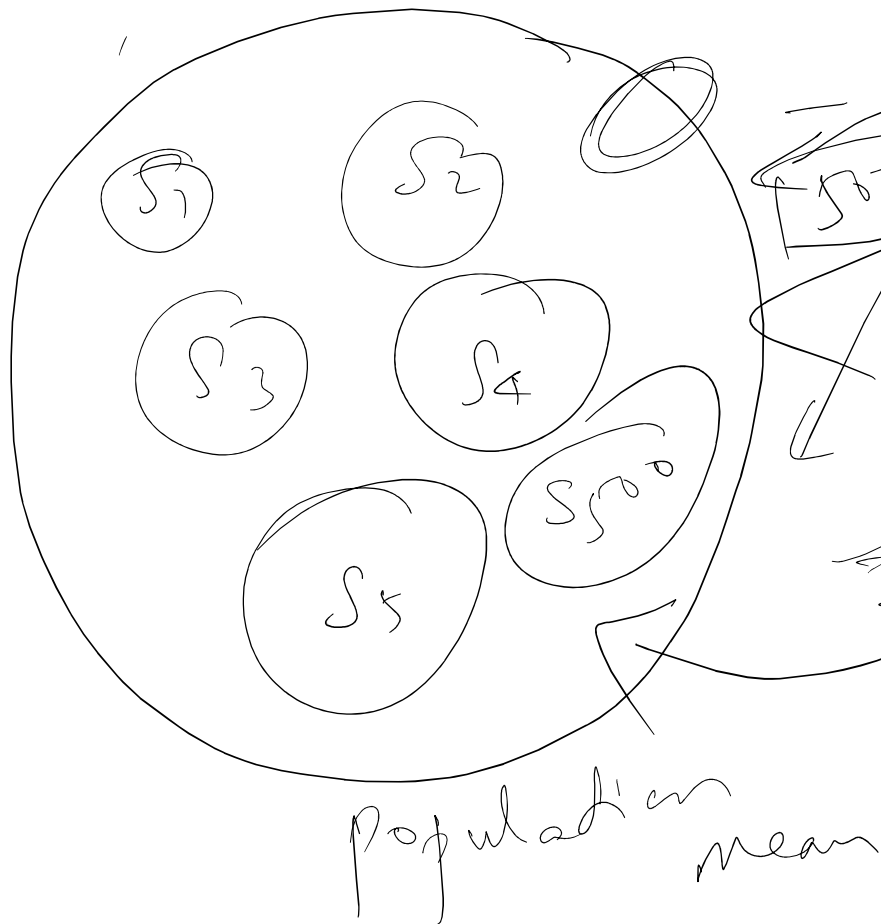
=> Normal

distribution = Symmetrical

Central Limit Theorem

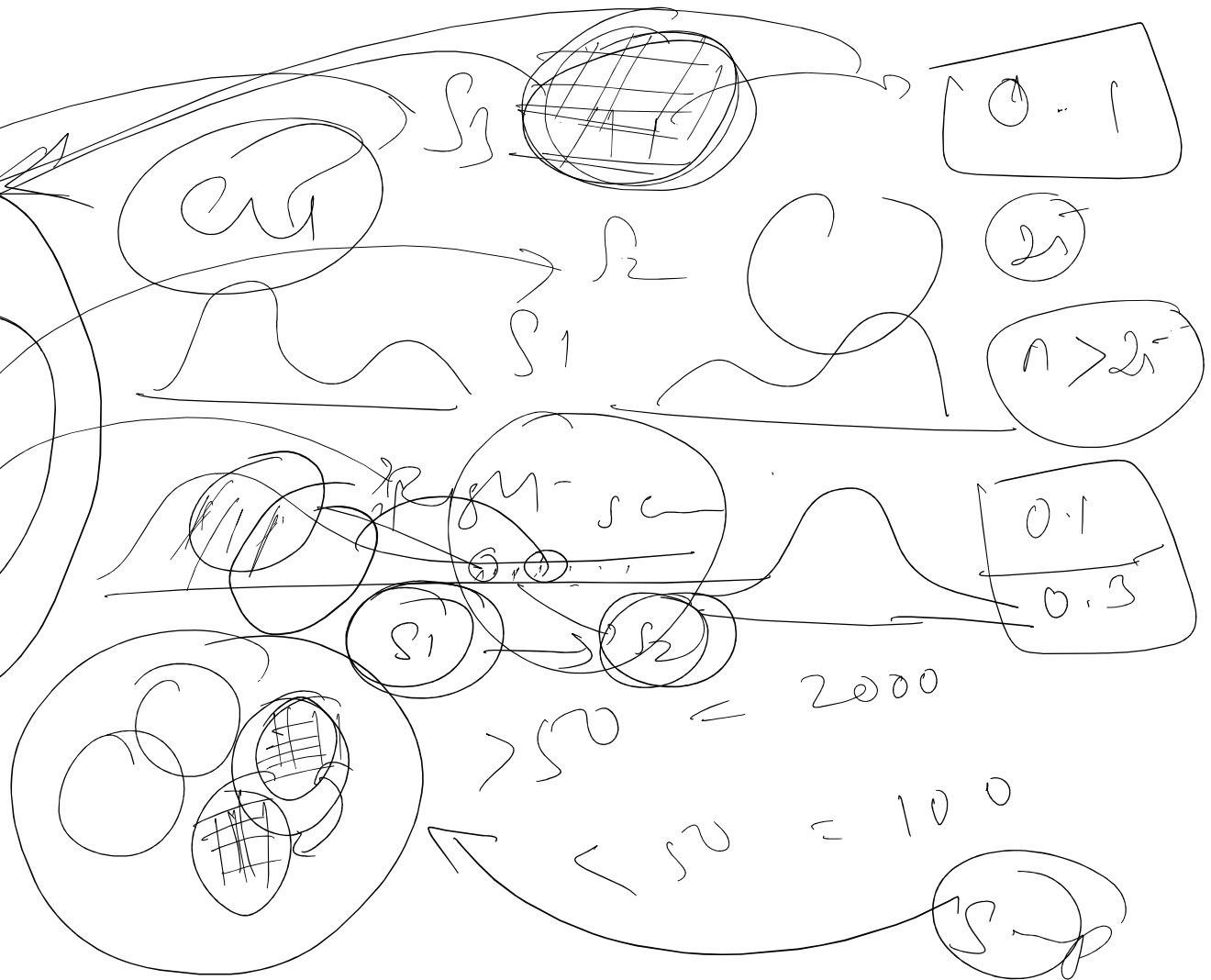
$$n \geq 25$$

$$n > 25$$



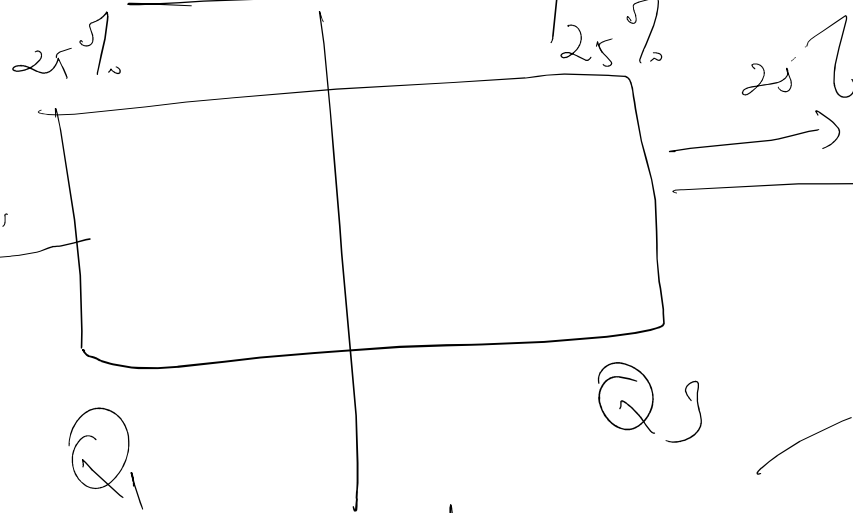
Sampling:

Population



Outlier:

Box - Whisker plot.



Outlier 25%

V1

Median

IQR

Mean	5.10	50%
Median	5.8	50%
Min	4.5	30
Max	12	1,000,000,000

min	200	500	700	max
10				

Pandas - profiling: ()

Measure of Dispersion:

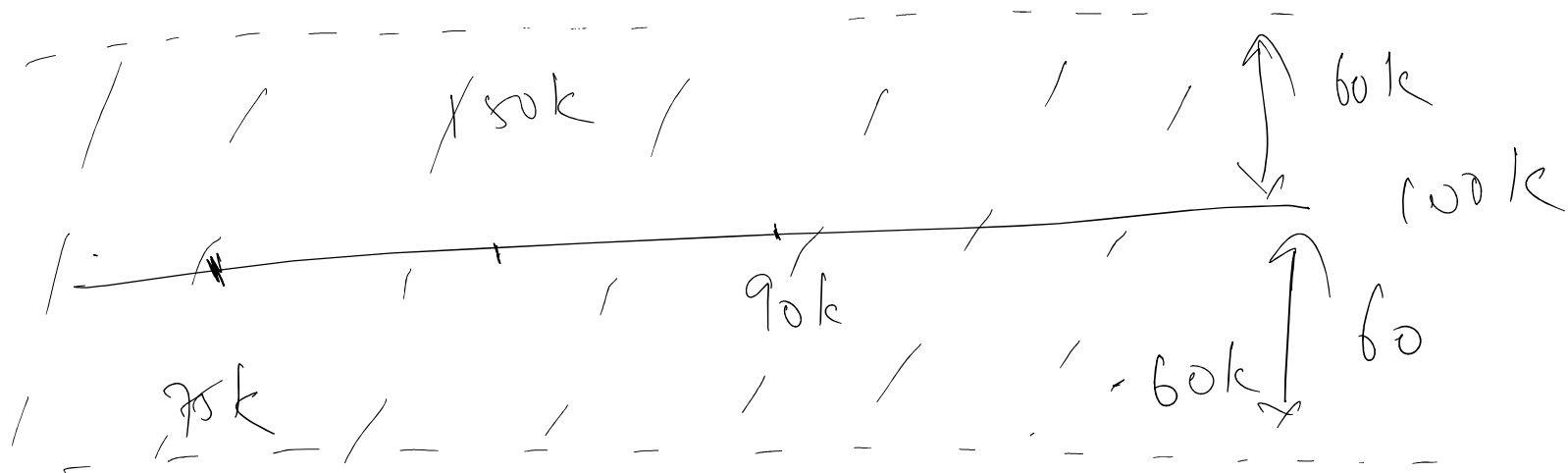
- Standard deviation

- Range (max - min)

(3, 7, 9, 10, 11, 15)
 $15 - 3 = 12$

Standard Deviation: $\pm 2\sigma$

68%





sample \leftarrow $\textcircled{M} \pm 5 \text{ ft}$

~~10~~
 $\sigma = 0.5$



$68\% \Rightarrow 1\sigma$

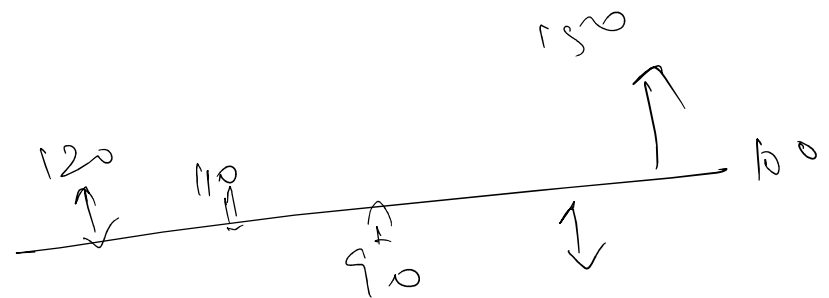
$95\% \Rightarrow 2\sigma$

$99.7\% \Rightarrow 3\sigma$

Coeff of Variation: (CV)

$$CV \Rightarrow \frac{\text{Std}}{\text{Mean}} \times 100$$

$$\sqrt{V} = \text{S.D}$$



$$0 < CV < 1$$

$$\begin{aligned} & (20)^2 + (10)^2 + (-10)^2 + (-40)^2 + (-90)^2 \\ & \hline & = \text{Variance} \Rightarrow \sqrt{12050} \\ & \hline & = \text{S.D} \end{aligned}$$

Pen : — (24) CY2022 (24) CY2023 CY2024
 J F M . . . D J F M . . . J F M

Mumbai (5m) = [Tk 2k X X X 5k . . .]
 Chennai (1m) = 12 12
 Delhi (8m) = 24
 Kolkata (2m) =
 Bangalore (6m) = ✓ ✓ ✓ ✓ ✓

$M_m \Rightarrow$

