



# Statistical Thinking for Data Science

① Descriptive <sup>Statistics</sup> ~~Data~~ Sumit Kumar Yadav

Department of Management Studies

② Inferential Statistics

May 15, 2022





- ❏ Descriptive Statistics
- ❏ Concept of Probability
- ❏ Random Variable, Discrete and Continuous
- ❏ Expected Value, Variance, Correlation



# Definitions of Statistics

SRT

1.)  
2.)  
:  
:  
500.)

HS

1.)  
2.)  
:  
:  
400.)



- ❑ Art of learning from data – Sheldon M. Ross, Introduction to Probability and Statistics for Engineers and Statistics
- ❑ Statistics is a branch of mathematics working with data collection, organization, analysis, interpretation and presentation. - Wikipedia
- ❑ Statistics may be regarded as (i) the study of populations, (ii) as the study of variation, (iii) as the study of methods of the reduction of data. – Fisher, 1925

- ❑ Describing the data, Summarize the data, etc
- ❑ Using numbers, pictures etc. → ~ Box Plot ~
- ❑ [https://www.ted.com/talks/hans\\_rosling\\_the\\_best\\_stats\\_you\\_ve\\_ever\\_seen](https://www.ted.com/talks/hans_rosling_the_best_stats_you_ve_ever_seen)
- ❑ [https://www.ted.com/talks/hans\\_rosling\\_asia\\_s\\_rise\\_how\\_and\\_when](https://www.ted.com/talks/hans_rosling_asia_s_rise_how_and_when)

# Summary Statistics



$$X_1, X_2, \dots, X_n$$

$$\bar{\mu} = \frac{X_1 + X_2 + \dots + X_n}{n}$$

$$\begin{array}{r|l} \sum & \sum \\ 100 & 80 \\ 150 & 90 \\ 200 & 1000 \end{array}$$

~Average~

- ❑ Measures of Central Tendency (mean, median, mode)
- ❑ Measures of Dispersion (Range, Variance)
- ❑ Chebyshev Inequality  $\rightarrow (\text{Max Value} - \text{Min Value})$

<u>D1</u>		<u>D2</u>	<u>D3</u>
8	(-10)	0	9
9	(-5)	5	9
10	0	10	9
11	5	15	9
12	10	20	14



Co-variance and Correlation



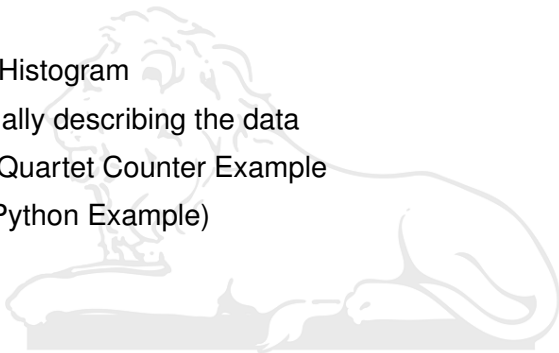


Scatter Plot, Histogram

Need for visually describing the data

Anscombe's Quartet Counter Example

Box-Plot (in Python Example)



# A Few More terminologies



- ☐ Cross-sectional Data
- ☐ Time Series Data
- ☐ Panel Data

- ☐ Qualitative Data

1. Nominal
2. Ordinal

- ☐ Quantitative Data

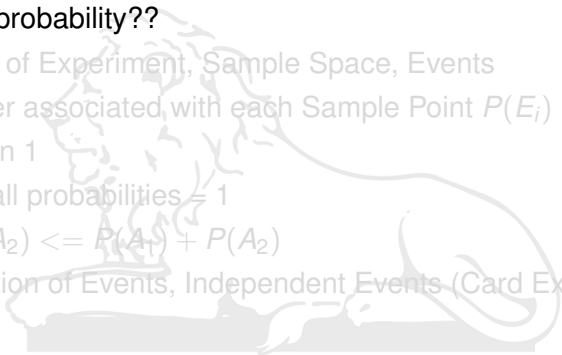
1. Interval
2. Ratio







- ❑ What is probability??
- ❑ Concept of Experiment, Sample Space, Events
- ❑ A number associated with each Sample Point  $P(E_i)$
- ❑ Less than 1
- ❑ Sum of all probabilities = 1
- ❑  $P(A_1 \cup A_2) \leq P(A_1) + P(A_2)$
- ❑ Intersection of Events, Independent Events (Card Example)



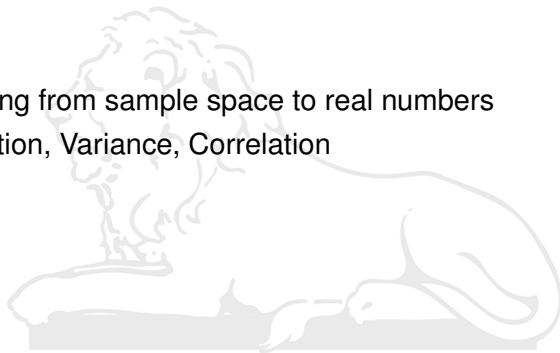


- ❑ What is probability??
- ❑ Concept of Experiment, Sample Space, Events
- ❑ A number associated with each Sample Point  $P(E_i)$
- ❑ Less than 1
- ❑ Sum of all probabilities = 1
- ❑  $P(A_1 \cup A_2) \leq P(A_1) + P(A_2)$
- ❑ Intersection of Events, Independent Events (Card Example)

# GodBole's Problem



- ❑ A mapping from sample space to real numbers
- ❑ Expectation, Variance, Correlation



# Pooled Testing Example



A LAB doing Covid testing gets 1000 samples to test everyday. However, due to the positivity rate drop in cases of Covid samples, the LAB is contemplating if it is better to mix the samples to get the result in lesser number of tests. Assuming 3% positivity rate, what is the number of samples that should be pooled together?