

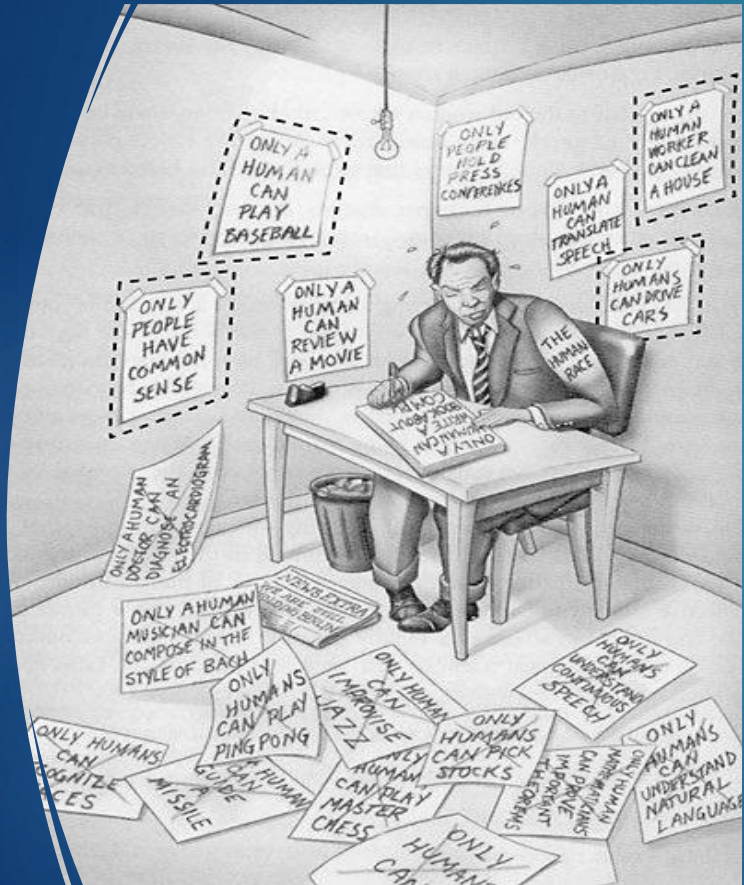


# Executive Certificate Program in Machine Learning & Artificial Intelligence

NEENA PANDEY

IIM VISAKHAPATNAM

# AI & ML



The 50 most attractive nationalities revealed: India is No.1, USA comes second and Britain has the most handsome men (while AI images show 'beautiful people' in each country)

By Ted Thornhill, Mailonline Travel Editor  
12:42 03 Mar 2023, updated 12:13 04 Mar 2023



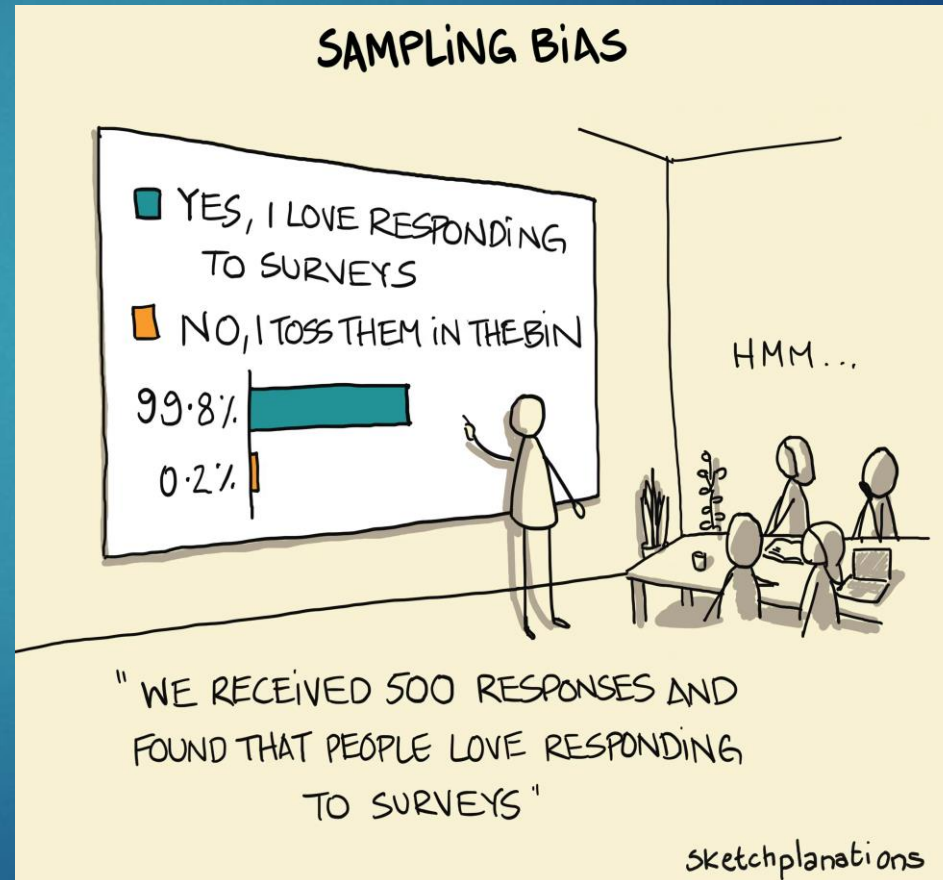
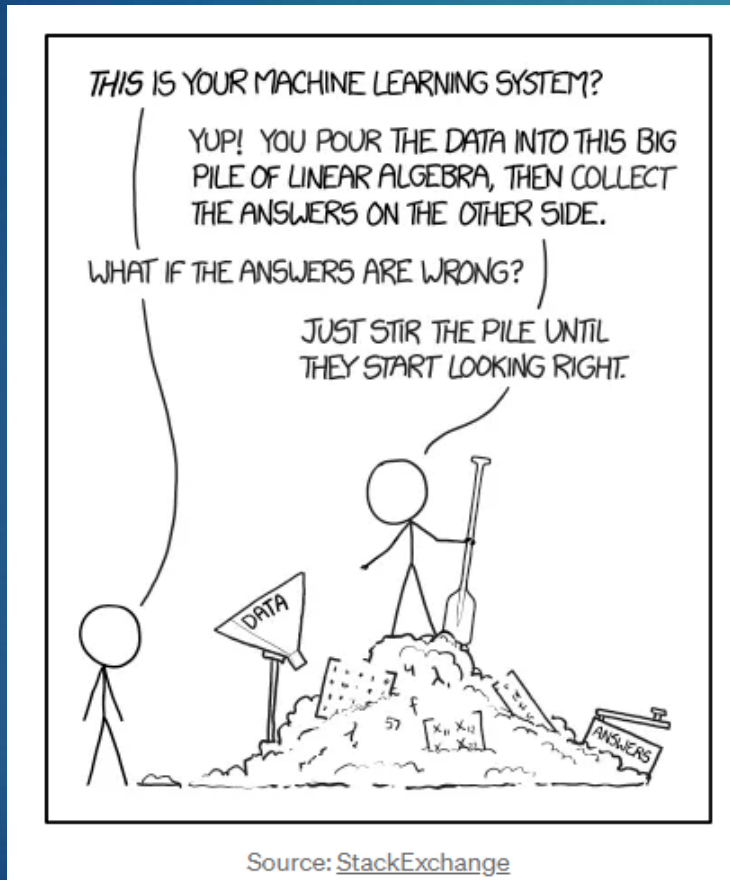
pretty wild that the first job openai took was sam's


2:22 AM · Nov 18, 2023

27.8K Reply Share

Read 268 replies

# AI/ML & Statistics





# Module 1: Fundamentals of Statistical Learning

# Over the course of 4 sessions

- ▶ Descriptive Statistics
  - ▶ Types of measurement, Tabular & Pictorial representation
  - ▶ Numerical measures – Central tendency (Location) and Dispersion
- ▶ Probability
  - ▶ Probability (Basic & Conditional), Bayes' Theorem
  - ▶ Probability distributions
- ▶ Hypothesis Testing
  - ▶ Point & Interval Estimation
  - ▶ Hypothesis Testing

# Statistics

- ▶ Refers to *numerical facts* such as averages, medians, percentages, and maximums that help us understand various business and economic situations.
- ▶ Also refers to the *art and science* of collecting, analyzing, presenting, and interpreting data.
- ▶ Business Applications
  - ▶ Statistical sampling procedures used by public auditing firms
  - ▶ Financial advisors may use PE ratio and dividend yields to guide their investment advice
  - ▶ Production: Statistical control charts to monitor production output
  - ▶ Information Systems: To monitor performance and anomaly in computer networks

# Data

- ▶ Elements, Variables and Observations
- ▶ Scales of measurement
  - ▶ Nominal, Ordinal, Interval, and Ratio
  - ▶ Categorical & Quantitative data
  - ▶ Cross-sectional vs. Time Series Data
- ▶ Data Sources
  - ▶ Internal company records – across functional departments
  - ▶ Govt. agencies
  - ▶ Business Database services
  - ▶ Industry associations
  - ▶ Internet etc.

# Types of Statistical Studies

- ▶ Observational
  - ▶ No attempt is made to control or influence the variables of interest
  - ▶ Example: Survey
- ▶ Experimental
  - ▶ Identification of the variable of interest; Performed under controlled conditions
  - ▶ Other variables are identified and controlled or changed to study their impact on the variable of interest
  - ▶ Example: Agriculture, Medical experiments, A/B testing for websites or other digital products

# Statistical Inference

- ▶ Keywords: Census, Sample survey, population, Statistics inference
- ▶ Statistical inference
  - ▶ Estimation of the population characteristic of interest
  - ▶ Provide a statement of quality or precision associated with the estimate
  - ▶ Point and Interval estimation
  - ▶ Confidence interval

# Analytics & Business Use Cases

- ▶ Analytics
  - ▶ Transforming data into insights for better decision-making
- ▶ Descriptive – Our focus in this module
  - ▶ Data summarized and presented in an understandable form; Describing what happened in the past
  - ▶ Customer segmentation; Sentiment Analysis
- ▶ Predictive
  - ▶ Process Automation – using Chatbots; Customer Life Time Value
- ▶ Prescriptive
  - ▶ Optimization with constraints

# Predictive vs. Prescriptive

## Predictive

- Models certain aspects of business
- Forecasts what's likely to happen
- Outputs are non-actionable
- Tends to optimize one function at the expense of others
- Based on hypotheses using pre-determined scenario with finite options

## Prescriptive

- Models the entire business
- Recommends specific business decisions
- Considers inter-dependencies
- Supports what-if scenarios
- Accounts for all inputs, variables, and outputs
- Models that truly reflect how the business operates

# Predictive vs. Prescriptive – Use Cases & Methods

## Predictive

- Demand planning, Inventory Control, Maintenance requirements, Customer churn
- Regression – Linear, logistic; Neural networks, Naïve Bayesian Classifiers

## Prescriptive

- Optimizing inventory over all stores to meet customer requirements and increase overall profitability, Optimizing manufacturing and inventory strategy
- Heuristics, Optimization – Linear Programming, Transportation/Assignment Problems



# Descriptive Statistics: Tabular & Graphical Representation

# Categorical Data Summarization: Graphs & Tables

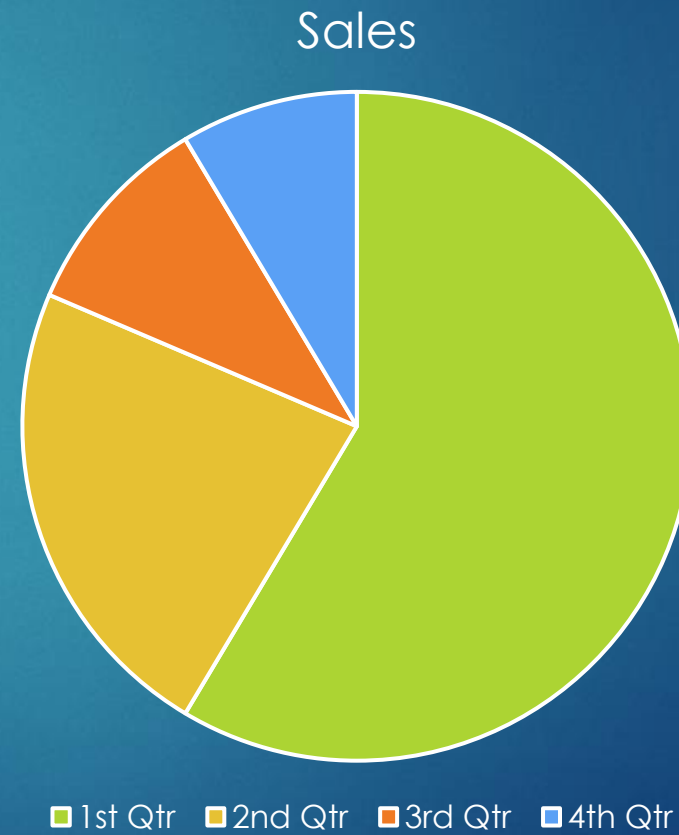
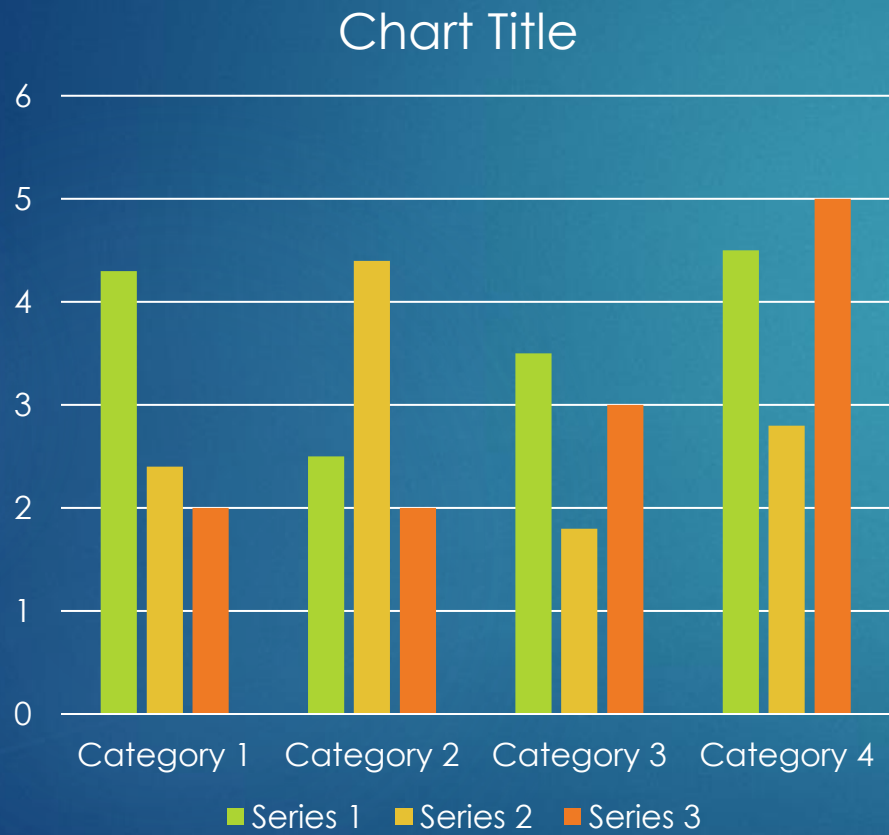
- ▶ Frequency Distribution
  - ▶ Tabular summary of data showing the number (frequency) of observations in each of several *non-overlapping categories* or classes.
- ▶ Relative Frequency distribution/ Percentage frequency distribution
- ▶ Bar Chart
  - ▶ A graphical display for depicting qualitative data
  - ▶ Using a bar of fixed width drawn above each class label
- ▶ Pie Chart
  - ▶ Relative frequency and percent frequency distributions for categorical data.

# Frequency Distribution: Tabular



<b>Ratings</b>	<b>Frequency</b>	<b>Relative Frequency</b>	<b>Percent Frequency</b>
Poor	2	0.10	10%
Below Average	3	0.15	15%
Average	5	0.25	25%
Above Average	9	0.45	45%
Excellent	1	0.05	5%
Total	20	1.00	100%

# Bar Chart & Pie Charts



# Quantitative Data Summarization: Graphs & Tables

- ▶ Only one variable of interest
  - ▶ Frequency, Relative Frequency, and Percent Frequency Distributions
  - ▶ Determine the number of nonoverlapping classes, width of class, class limits
  - ▶ Histogram
    - ▶ Horizontal axis – variable of interest; Vertical axis – frequency, relative frequency, percent frequency
  - ▶ Cumulative Distributions – Facebook adoption over the years
- ▶ Relationship between two variables
  - ▶ For comparisons between the two:
    - ▶ Side-by-side Bar Chart – for comparing two variables
    - ▶ Stacked Bar Chart – to compare relative frequency
  - ▶ For relationships between the two:
    - ▶ Scatter plot
    - ▶ Trendline



# Descriptive Statistics: Numerical Measures

# Quantitative Data Summary: Numerical Measures

- ▶ Measures of Central Tendency (Location)
- ▶ Measures of Variability (Dispersion)
- ▶ Measures of Distribution Shape, Relative Location, Detecting Outliers
- ▶ Measures of Association

# Measures of Central Tendency

- ▶ Mean, Median, Mode, Weighted Mean, Percentiles, Quartiles
- ▶ Mean
  - ▶ Sample mean  $\bar{x}$  is the point estimator of the population mean,  $\mu$ .
- ▶ Median

$$\bar{x} = \frac{\sum x_i}{n}$$

  - ▶ Value in the middle of a data set when the data items are arranged in ascending order.
  - ▶ When is it preferred? – When a data set has extreme value
- ▶ Mode
  - ▶ Value that occurs with the greatest frequency

# Measures of Central Tendency

- ▶ Weighted Mean

- ▶ Mean where each observation has a weight that reflects its reflective importance
- ▶ Choice of weights depends upon the application – for e.g., number of credits earned for each grade

$$\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$$

where:  $x_i$  = value of observation  $i$   
 $w_i$  = weight for observation  $i$

- ▶ Another example?

# Measures of Central Tendency

## ▶ Percentiles

- ▶ Provides information about how the data are spread over the interval from the smallest value to the largest value.
- ▶  $p^{\text{th}}$  percentile of a data set is a value such that at least  $p$  percent of the items take on this value or less and at least  $(100 - p)$  percent of the items take on this value or more.
- ▶ Arrange the data in ascending order.
- ▶ Compute  $L_p$ , the location of the  $p^{\text{th}}$  percentile.

$$L_p = \left( \frac{p}{100} \right) (n + 1)$$

# Measures of Central Tendency

- ▶ 80<sup>th</sup> percentile for 70 observations

$$L_p = \left(\frac{p}{100}\right)(n + 1) = \left(\frac{80}{100}\right)(70 + 1) = 56.8$$

- ▶ For e.g., NEET Score; Percentile =  $635 + 0.8(649 - 635) = 646.2$

525	530	530	535	535	535	535	535	540	540
540	540	540	545	545	545	545	545	550	550
550	550	550	550	550	560	560	560	565	565
565	570	570	572	575	575	575	580	580	580
580	585	590	590	590	600	600	600	600	610
610	615	625	625	625	635	649	650	670	670
675	675	680	690	700	700	700	700	715	715

# Measures of Central Tendency: Quartiles

- ▶ First Quartile = 25<sup>th</sup> Percentile
- ▶ Second Quartile = 50<sup>th</sup> Percentile
- ▶ Third Quartile = 75<sup>th</sup> Percentile
  
- ▶ Second Quartile is same as ??

# Measures of Variability

- ▶ Think of a business case where even if you know measure of location, measure of variability add significantly high value!
- ▶ Examples
  - ▶ Suppliers' average delivery time and variability in delivery time
  - ▶ Average return on investment and variation in investment return
- ▶ Common measures
  - ▶ Range, Interquartile Range, Variance, Standard Deviation, Coefficient of variation

# Measures of Variability

- ▶ Range
  - ▶ Difference between the largest and smallest data value.
  - ▶ Very sensitive to two specific data points
- ▶ Interquartile Range
  - ▶ Fixes earlier problem
  - ▶ Range for the middle 50% of the data – the difference between the third and first quartile
  - ▶  $IQR = Q_3 - Q_1$  – Try with the given data
- ▶ Variance
  - ▶ Measure of variability that utilizes all the data
  - ▶ average of the squared deviations between each data value and the mean

$$s^2 = \frac{\sum(x_i - \bar{x})^2}{n - 1}$$

# Measures of Variability

- ▶ Standard Deviation

- ▶ More easily interpretable

- ▶ Coefficient of Variation

- ▶ How large the deviation is in relation to the mean
  - ▶ What's a good value for Coefficient of Variation
  - ▶ When can this be useful?

$$\left[ \frac{s}{\bar{x}} \times 100 \right] \%$$

# Standardized Value

- ▶ Number of standard deviations a data value  $x_i$  is from the mean.
- ▶ A measure of the relative location of the observation in a data set.
- ▶ Think of a Use-case

$$Z_i = \frac{x_i - \bar{x}}{s}$$

# Detecting Outliers – Why & How to detect

## ▶ Why

- ▶ unusually small or unusually large value in a data set.
- ▶ data value with a z-score less than  $-3$  or greater than  $+3$  might be considered an outlier.
- ▶ Reasons: Might be an incorrectly recorded value, incorrect inclusion into the dataset, or correct value but an outlier data point

## ▶ Detection

- ▶  $|z| > 3$

# Five-Number Summaries and Box Plots

- ▶ Two tools that accomplish this are five-number summaries and box plots.

- ▶ Smallest Value
- ▶ First Quartile
- ▶ Median
- ▶ Third Quartile
- ▶ Largest Value

525	530	530	535	535	535	535	535	540	540
540	540	540	545	545	545	545	545	550	550
550	550	550	550	550	560	560	560	565	565
565	570	570	572	575	575	575	580	580	580
580	585	590	590	590	600	600	600	600	610
610	615	625	625	625	635	649	650	670	670
675	675	680	690	700	700	700	700	715	715

- ▶ Box Plots

# Measure of Association

- ▶ Till now, summarized data for one variable at a time
- ▶ What's the relationship between two variables?
  - ▶ Covariance & Correlation Coefficient
- ▶ Covariance
  - ▶ Measure of linear association between two variables

$$s_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$