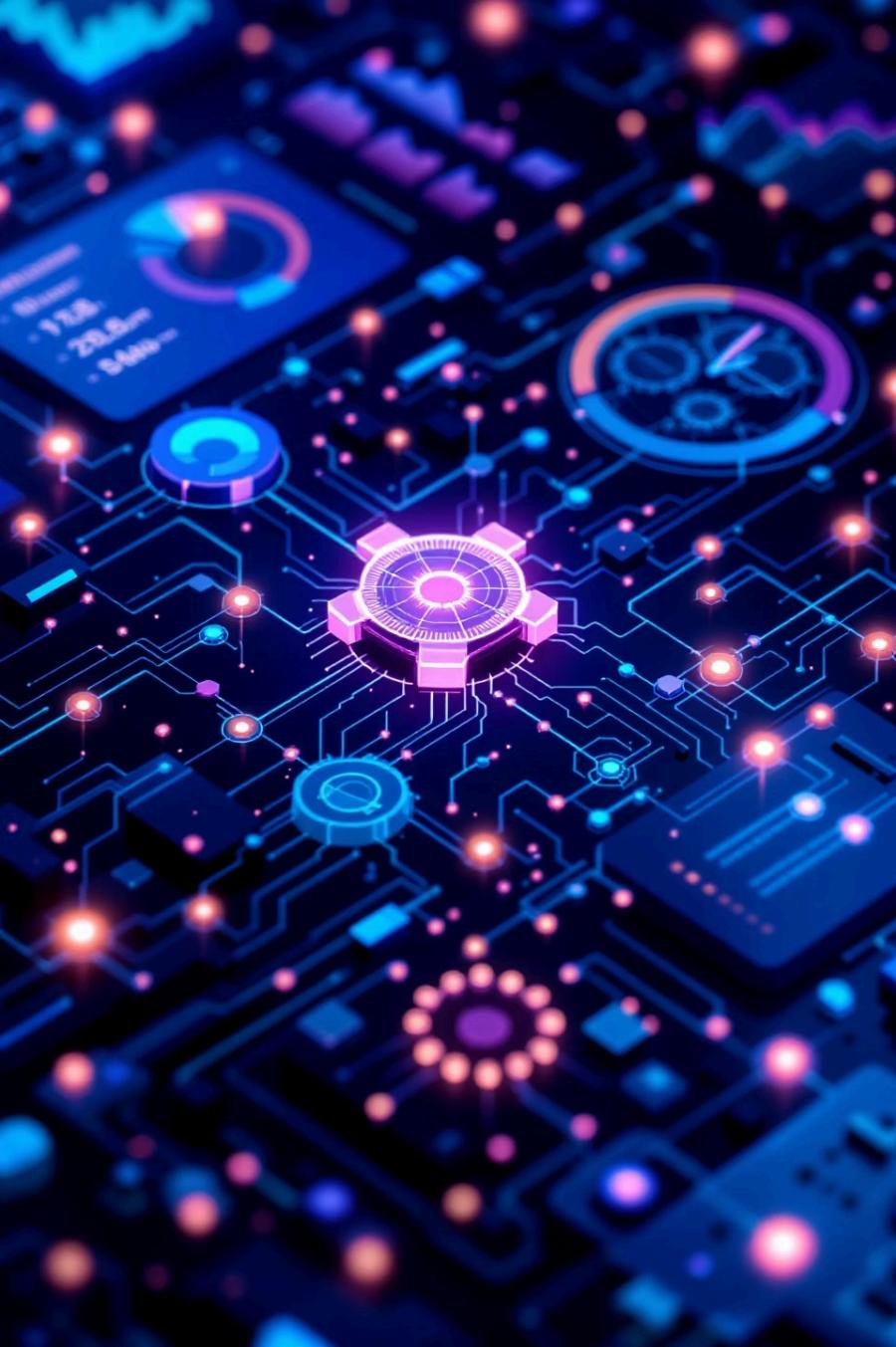




Data Science and Statistical Analysis

Welcome to this comprehensive presentation on data science fundamentals. We'll explore the core concepts of data types, statistical measures, probability, and artificial intelligence. This presentation provides both theoretical foundations and practical examples to help you understand these critical concepts in data analysis and machine learning.

Muhammad Saeed



Understanding Data Types

Structured Data

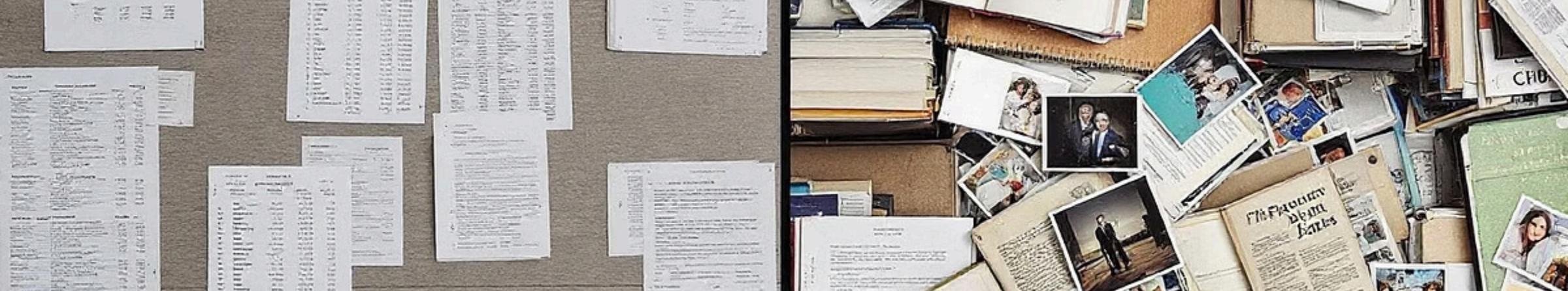
Organized and easily searchable in tabular formats. Usually stored in relational databases (SQL) and fits neatly into a predefined model.

Unstructured Data

Unorganized data that lacks a clear structure and doesn't fit into traditional databases. More difficult to analyze without special tools.

Semi-structured Data

Contains tags or markers to separate elements but doesn't conform to the structure of standard databases. Examples include XML and JSON files.




 Search
 Student for hmp air studennts

Larck None

Name

Name	ID	Major	GPA
D:	\$:45,000		D
MeA	\$:25,000		1
GPA	\$:25,000		D
Name	\$:45,000		D
Laine	\$:55,000		1
Btannda	\$:25,000		D
Cannel	\$:35,000		V
Drickien	\$:90,000		U
WPA	\$:65,000		NO
GPA	\$:20,000		D
GPA	\$:25,000		D
GPA	\$:55,000		D
GPA	\$:25,000		D
Nane	\$5,000		D
New	\$:17,000		D
Haney	\$:35,000		D
Cale	\$:25,000		D
GPA	\$:45,000		D

Structured Data Examples

Student ID	Name	Age	Grade
001	John Smith	18	A
002	Sarah Johnson	17	B+
003	Michael Brown	18	A-

Structured data is organized in a predictable way, making it easy to search and analyze. This student database example shows how data is arranged in rows and columns with clear relationships. Each student has a unique identifier and associated attributes like name, age, and grade.

Other examples include bank transactions with date, amount, and transaction type, or inventory systems with product codes, quantities, and prices.

Unstructured Data Examples



Social Media

Posts, comments, and interactions that contain text, images, videos, and other media without a standardized format.

Unstructured data makes up approximately 80-90% of all data generated today. It requires specialized tools like natural language processing, computer vision, and machine learning to extract meaningful insights.



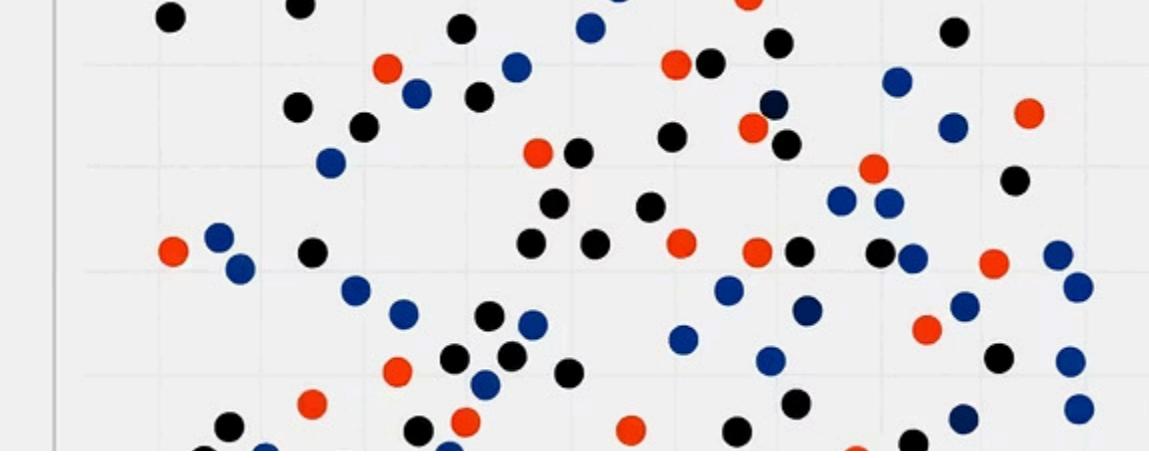
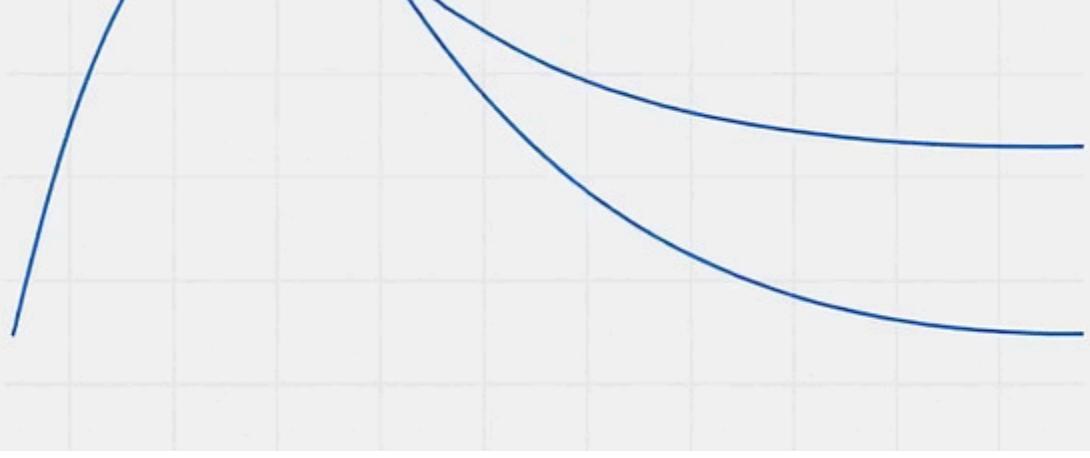
Emails

Messages with varying formats, attachments, and content that don't follow a consistent structure.



Medical Images

X-rays, MRIs, and other diagnostic images that contain valuable information but aren't organized in a tabular format.



Quantitative Data



Numerical Data

Data that represents measurable quantities expressed as numbers.



Continuous Variables

Can take any value within a range, including fractions or decimals.



Discrete Variables

Can take only specific values, often whole numbers that are countable and finite.

Continuous vs. Discrete Variables

Continuous Variables

Can take any value within a range, including fractions or decimals. There are infinite possible values.

- Temperature (e.g., 36.6°C)
- Height (e.g., 175.4 cm)
- Weight (e.g., 60.5 kg)
- Time (e.g., 3.45 seconds)

Discrete Variables

Can take only specific values, often whole numbers. They are countable and finite.

- Number of students in a class (e.g., 25)
- Number of cars in a parking lot
- Count of errors in a program
- Number of children in a family



Qualitative Data



Categorical Variables

Data divided into groups or categories that describe attributes or qualities.



Nominal Variables

Categories without a specific order or ranking (e.g., eye color, gender).



Ordinal Variables

Categories with a meaningful order but no fixed interval (e.g., satisfaction levels).



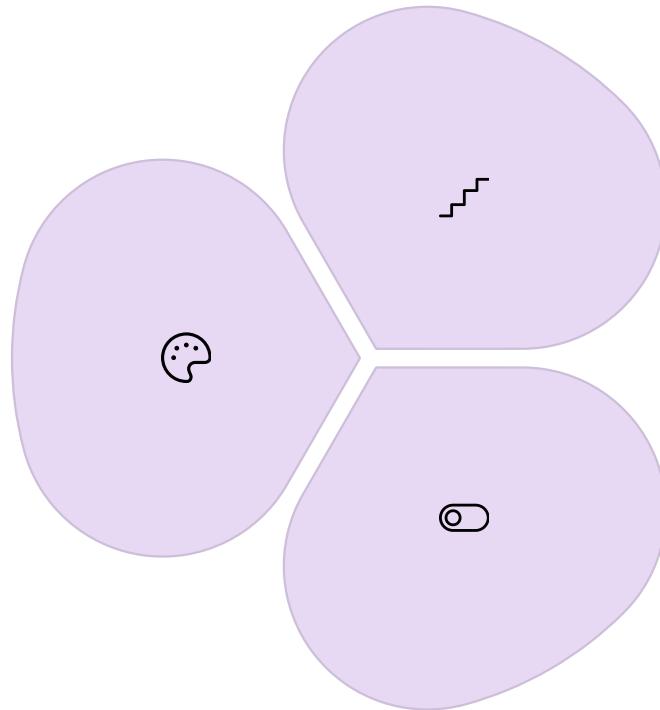
Binary Variables

Only two possible outcomes (e.g., yes/no, true/false, pass/fail).

Examples of Qualitative Data Types

Nominal Variables

- Eye color (blue, green, brown)
- Gender (male, female)
- Blood type (A, B, AB, O)



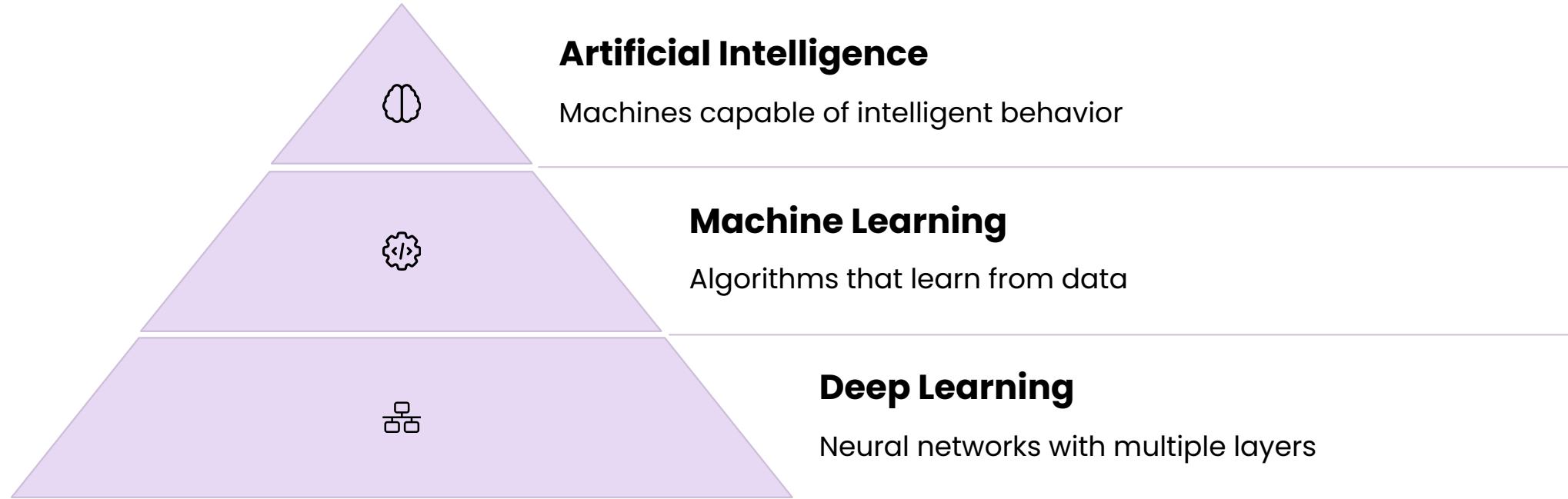
Ordinal Variables

- Satisfaction level (very satisfied, satisfied, neutral, dissatisfied)
- Education level (High School, Bachelor, Master, PhD)
- Movie ratings (1-5 stars)

Binary Variables

- Yes/No responses
- True/False values
- Pass/Fail results

Artificial Intelligence Overview



Artificial Intelligence is the broader field that aims to create machines capable of intelligent behavior, mimicking human decision-making and problem-solving. Machine Learning is a subset of AI that involves algorithms that learn from data and make predictions without explicit programming. Deep Learning is a specialized subset of ML using neural networks with multiple layers to learn from large amounts of data.



Artificial Intelligence Applications



Self-driving Cars

Autonomous vehicles that use AI to navigate roads, detect obstacles, and make driving decisions.



Virtual Assistants

AI-powered tools like Siri, Alexa, and Google Assistant that respond to voice commands and questions.



Healthcare Diagnostics

AI systems that analyze medical images and patient data to assist in disease diagnosis and treatment planning.



Fraud Detection

AI algorithms that identify unusual patterns in financial transactions to prevent fraud and enhance security.

Machine Learning Types



Supervised Learning

Algorithms learn from labeled training data



Unsupervised Learning

Algorithms find patterns in unlabeled data

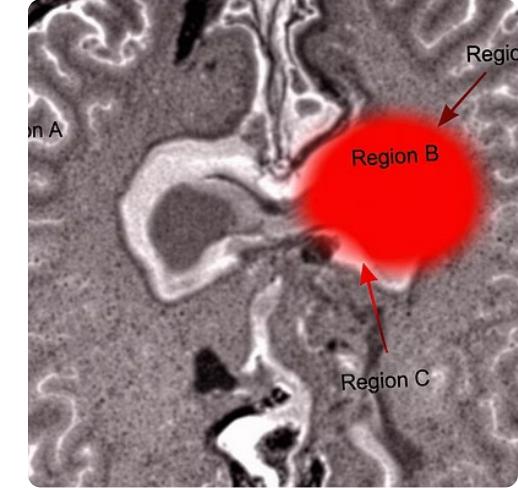
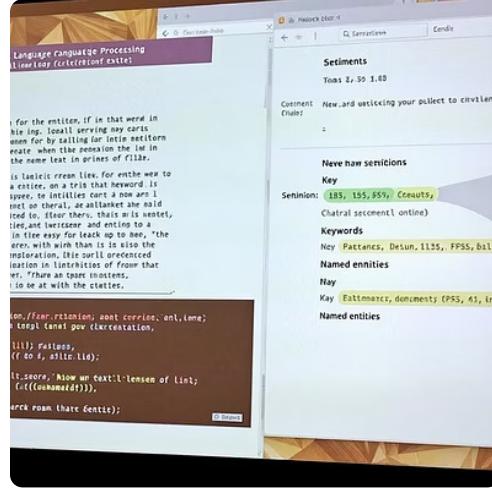


Reinforcement Learning

Algorithms learn through trial and error

Supervised learning uses labeled data to train models for prediction tasks like house price prediction or email classification. Unsupervised learning discovers hidden patterns in unlabeled data, useful for customer segmentation or anomaly detection. Reinforcement learning enables agents to learn optimal behaviors through interaction with an environment, as seen in game-playing AI like AlphaGo.

Deep Learning Applications



Deep learning has revolutionized many fields through its ability to process and learn from unstructured data. Image recognition systems can identify faces and objects with human-level accuracy. Natural language processing powers chatbots and translation services. Computer vision enables autonomous vehicles to "see" and interpret their surroundings. Medical applications include analyzing X-rays and MRIs to detect diseases earlier and more accurately than human doctors in some cases.

Measures of Central Tendency



Mean (Arithmetic Average)

The sum of all values divided by the number of values. Sensitive to extreme values or outliers.



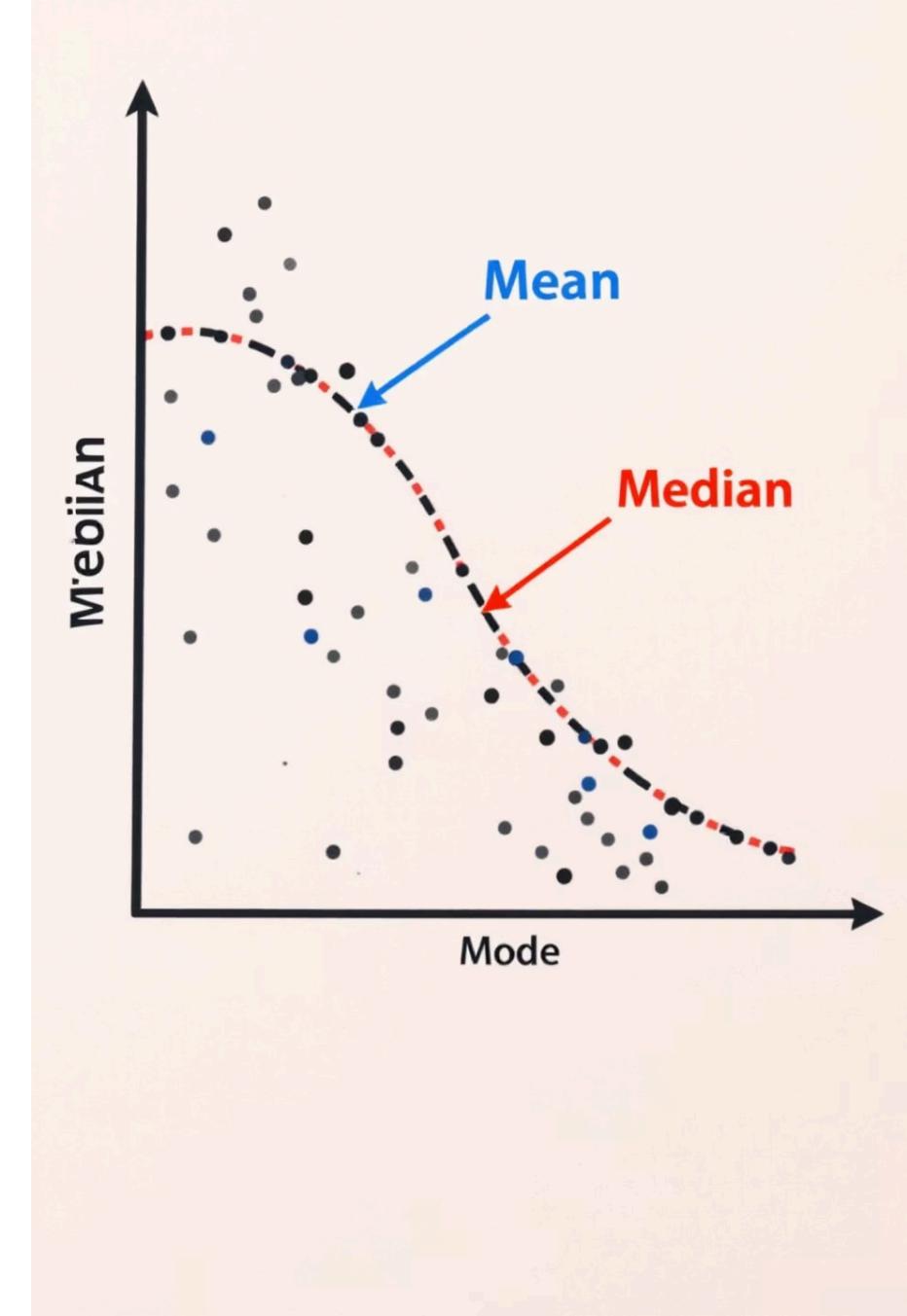
Median

The middle value when data is arranged in order. Less affected by outliers than the mean.



Mode

The value that occurs most frequently in a dataset. A dataset can have no mode, one mode, or multiple modes.



he example are calculate of the mean all noo betying, meses,
lately for the cdata poins of the mean, thare, and i the can th



Mean: The Arithmetic Average

$$\Sigma X$$

Sum of Values

Add all data points together

$$N$$

Number of Values

Count the total data points

$$\Sigma X / N$$

Mean Formula

Sum divided by count

The mean is calculated by adding all values in a dataset and dividing by the total number of values. For example, if five students scored 60, 70, 80, 90, and 100 on an exam, the mean would be $(60+70+80+90+100)/5 = 400/5 = 80$. The mean is useful for calculating average income, temperature, or age, but can be skewed by extreme values or outliers.

2, 4, 6, 8, 10
 $= (6$

1, 3) 5 7
4 & 5

median

4 - 5

Median: The Middle Value



Sort the Data

Arrange all values in ascending or descending order

Count the Values

Determine if there is an odd or even number of values

Find the Middle

For odd count: select the middle value

For even count: average the two middle values

Median Examples

Odd Number of Elements

Data: 45, 55, 60, 70, 85

Sorted: 45, 55, 60, 70, 85

Median = 60 (middle value)

The median represents the central value that is unaffected by extreme values at either end of the distribution.

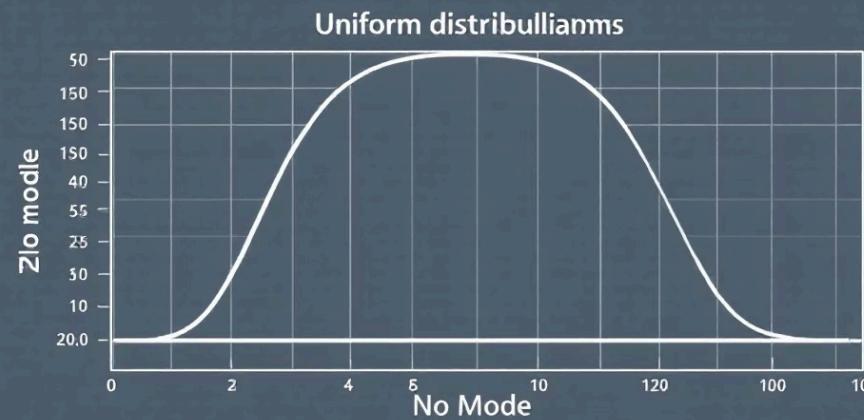
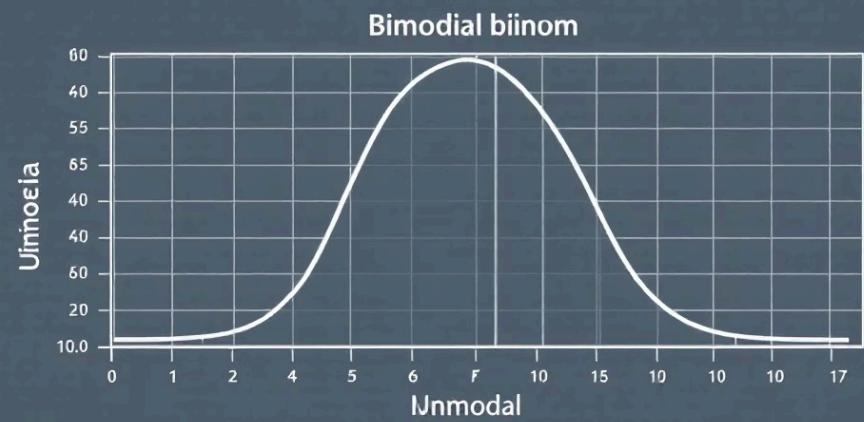
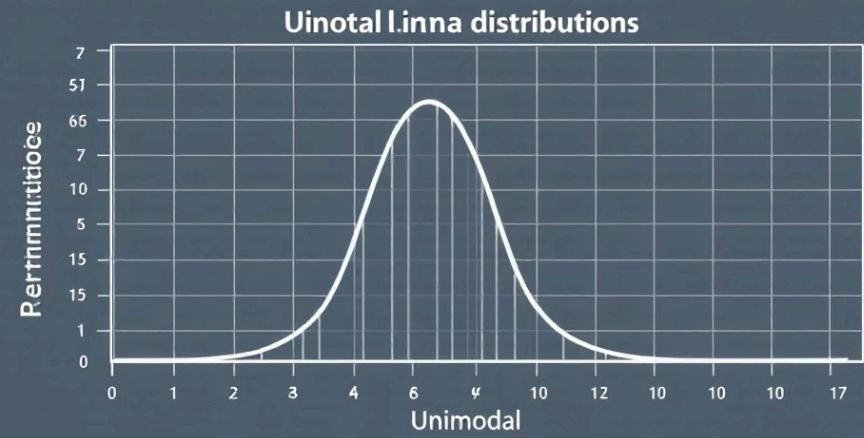
Even Number of Elements

Data: 20, 30, 40, 50

Sorted: 20, 30, 40, 50

Median = $(30 + 40) / 2 = 35$

With an even number of values, we take the average of the two middle values to find the median.



Mode: The Most Frequent Value

Unimodal

One value occurs most frequently

Example: 2, 4, 4, 6, 7

Mode = 4 (occurs twice)

Bimodal

Two values occur with equal highest frequency

Example: 3, 3, 5, 5, 7, 9

Modes = 3 and 5 (each occurs twice)

No Mode

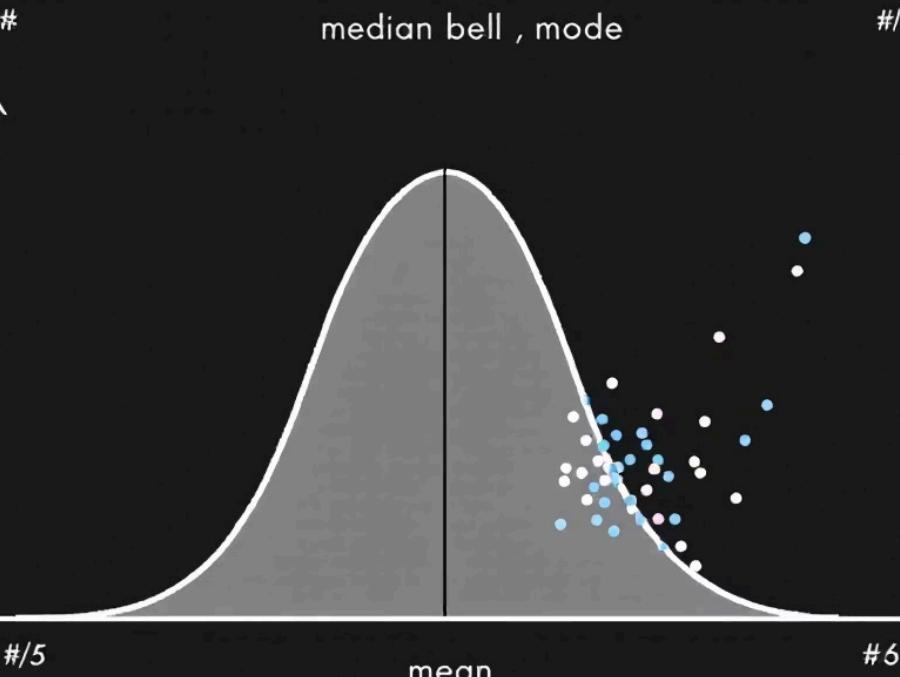
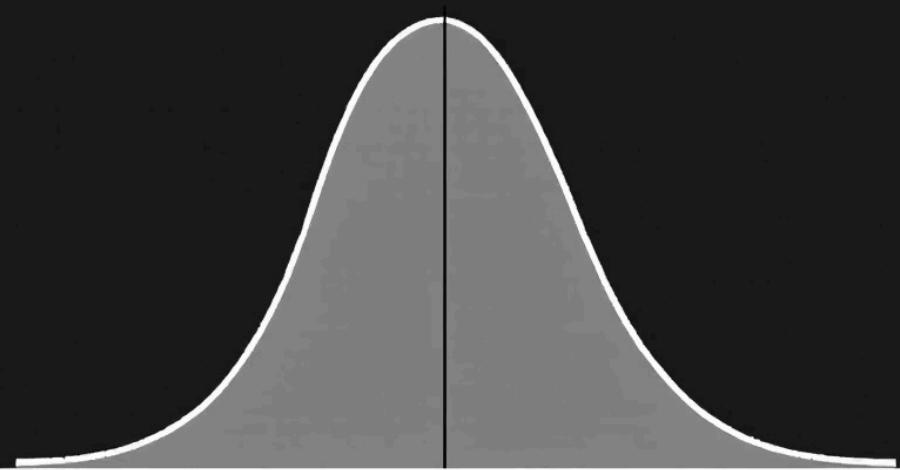
All values appear with equal frequency

Example: 1, 2, 3, 4, 5

No mode (all values appear once)

IMPACT OF OUTLIERS

Median mediutions and nor median, and Mode



Mean, pulled'stret by outliers, out whiens aufbyed "tail
out pulled out effect y out.lins by in the mean.



Comparing Measures of Central Tendency

Measure	Strengths	Weaknesses	Best Use Cases
Mean	Uses all data points	Affected by outliers	Normally distributed data
Median	Not affected by outliers	Ignores actual data values	Skewed data (e.g., income)
Mode	Shows most common value	May not exist or be unique	Categorical data

The choice between mean, median, and mode depends on the data distribution and the specific analysis goals. For normally distributed data, all three measures tend to be similar. For skewed distributions, the median often provides a better representation of the "typical" value than the mean.

Measures of Dispersion

Variance

Average squared deviation from
the mean

Kurtosis

Tailedness of distribution



Standard Deviation

Square root of variance

Coefficient of Variation

Relative variability as
percentage

Skewness

Asymmetry of distribution



Variance: Measuring Data Spread

Population Variance (σ^2)

$$\sigma^2 = \frac{\sum (X_i - \mu)^2}{N}$$

Where:

- X_i : Each value
- μ : Population mean
- N: Number of values

Sample Variance (s^2)

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

Where:

- X_i : Each value
- \bar{X} : Sample mean
- n: Number of values
- n-1: Degrees of freedom

Variance Example

Given Dataset

Data: 4, 6, 8, 10, 12

$$\text{Mean} = (4+6+8+10+12)/5 = 40/5 = 8$$

Calculate Squared Deviations

$$(4-8)^2 = 16, (6-8)^2 = 4, (8-8)^2 = 0, (10-8)^2 = 4, (12-8)^2 = 16$$

Find Sample Variance

$$s^2 = (16+4+0+4+16)/(5-1) = 40/4 = 10$$

A variance of 10 means the data values, on average, deviate by 10 squared units from the mean. The higher the variance, the more spread out the data. Variance is useful for comparing the spread of different datasets, but its squared units can make interpretation challenging.

How the variance?

$$\begin{aligned}\text{Step B: variance: } & X \times (+ \pm (X \times 1)) = 111 \\ & \times 25 - 60 \\ & x = 2 - 14\end{aligned}$$

Step B: numbers:

$$x - 2 \cdot x^4 (= 3.0^*)$$

Calculatt:

$$\begin{aligned}2 \cdot x^3 & (x (1, + \sqrt{= 6, 5}) = \\ & (1 \times 1,)^{\frac{1}{2}} = (5, <)\end{aligned}$$

$$\begin{aligned}\text{Step B: nussice } & \equiv 1 + + 1 \\ & = 1 \times + 5, * x\end{aligned}$$

Step 7: semout: $(2, 3.5) + 4$ in, the Vowy)

$$= (\bar{x} - 4^t) \left[\begin{array}{l} x \\ x = 2 \end{array} \right] = (x, 20, 15) \\ r = 1 \cdot 10$$

Step 4: imoneneng is $(1.2 + \bar{x} (n(1, h))$

$$5. (5 v x) \left[\begin{array}{l} 4,5 \\ x = 3 \end{array} \right] = < 1, x (6) + (1)$$

Calculate:

$$+ 1b (62) + 5 (7) in 1 \times 7) + = 0.55)$$

"Calculate is abelowing variance":

$$+ T, r) + xample - (2 (14)) = * (0) = 00, (67))$$

Standard Deviation



Definition

The square root of variance, expressing dispersion in the same unit as the original data.



Formula

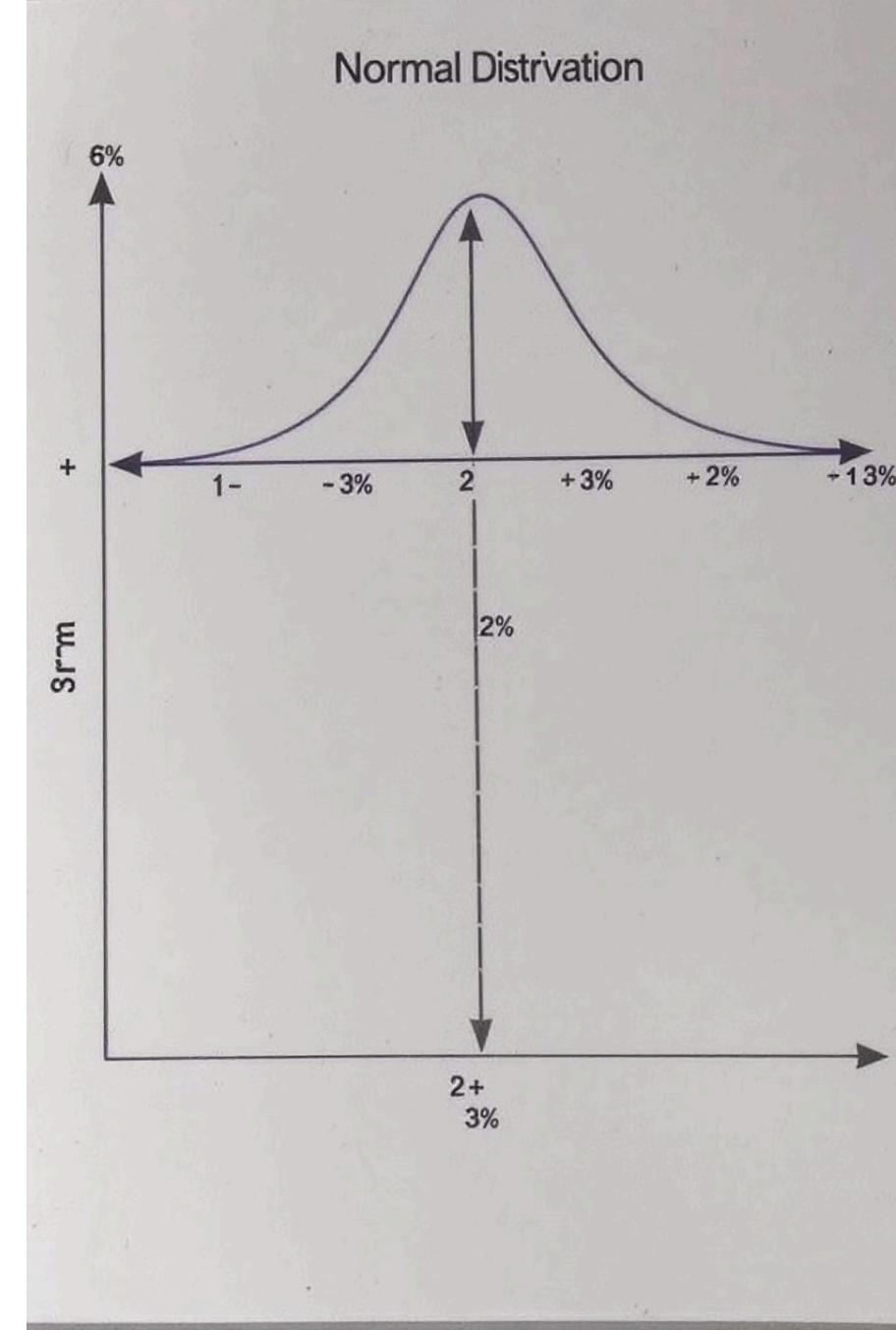
$$\sigma = \sqrt{\sigma^2} \text{ (population)} \text{ or } s = \sqrt{s^2} \text{ (sample)}$$



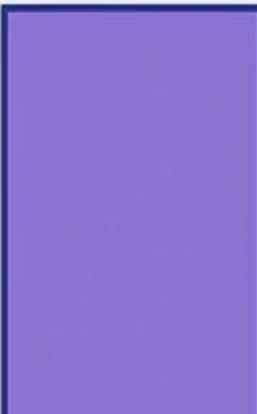
Interpretation

Smaller SD means data points are closer to the mean; larger SD indicates more variability.

From our previous example with a variance of 10, the standard deviation would be $\sqrt{10} \approx 3.16$. This means that, on average, data points deviate from the mean by about 3.16 units. Standard deviation is widely used because it's easier to interpret than variance, as it's expressed in the same units as the original data.



Coefficient of Variation



Dataset R

Coefficient of Variation (cv)

SD

Standard Deviation

Measure of absolute dispersion

Mean

Mean Value

Average of the dataset

SD/Mean

Relative Variation

Expressed as percentage

The Coefficient of Variation (CV) expresses standard deviation as a percentage of the mean, allowing comparison between datasets with different units or means. For example, if Machine A produces widgets with mean weight = 100g and SD = 2g, its CV = $(2/100) \times 100 = 2\%$. If Machine B produces widgets with mean weight = 10g and SD = 1g, its CV = $(1/10) \times 100 = 10\%$. This shows Machine A is more consistent in production despite having a larger standard deviation.

Skewness: Measuring Asymmetry

Positive Skew

Long tail on the right side of the distribution.

Mean > Median > Mode

Example: Income distribution where most people earn average amounts but a few earn very high salaries.

Negative Skew

Long tail on the left side of the distribution.

Mean < Median < Mode

Example: Exam scores where most students perform well but a few score very poorly.

Skewness helps determine if the average (mean) is higher or lower than the majority's actual value. A perfectly symmetrical distribution has a skewness of zero. Skewness is important for understanding the shape of your data distribution and can influence which statistical methods are appropriate to use.

Kurtosis: Measuring Tailedness

Mesokurtic ($K \approx 3$)

Normal distribution
(benchmark)

Standard bell curve shape

Leptokurtic ($K > 3$)

Sharper peak and heavier tails

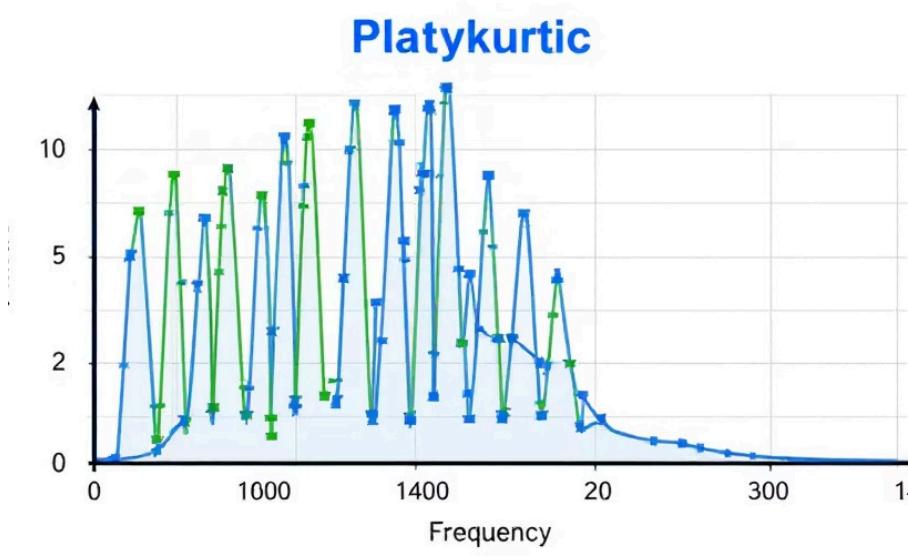
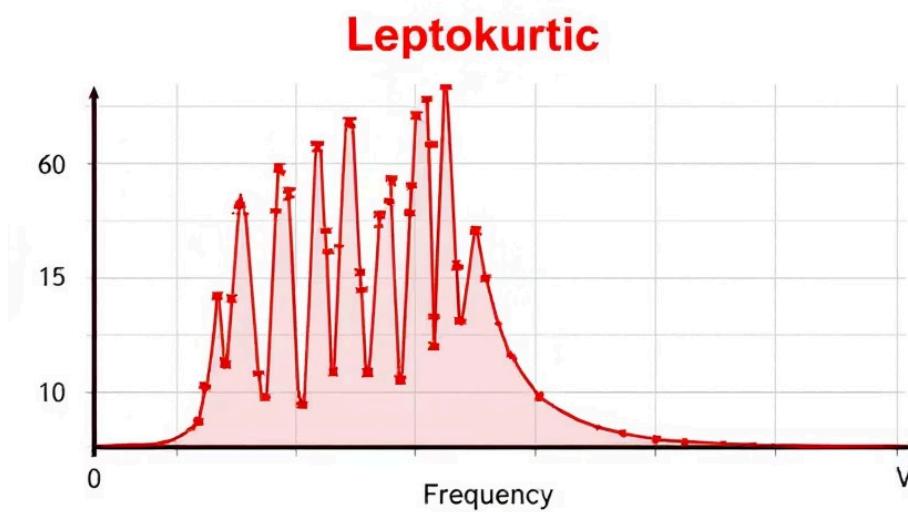
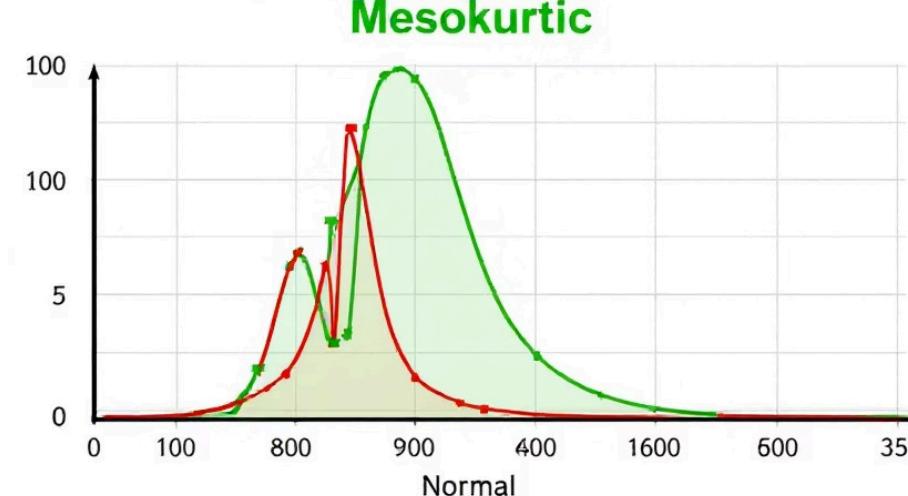
More outliers than normal distribution

Platykurtic ($K < 3$)

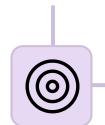
Flatter peak and lighter tails

Fewer outliers than normal distribution

Kurtosis indicates the "tailedness" or peakedness of a distribution. High kurtosis (leptokurtic) distributions have more frequent extreme outcomes, making them riskier. Financial returns often show high kurtosis. Low kurtosis (platykurtic) distributions have more evenly distributed data with fewer outliers, like uniformly spread exam scores.



Measures of Position



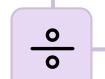
Z-Score

Standardized score showing distance from mean in standard deviation units



Percentile

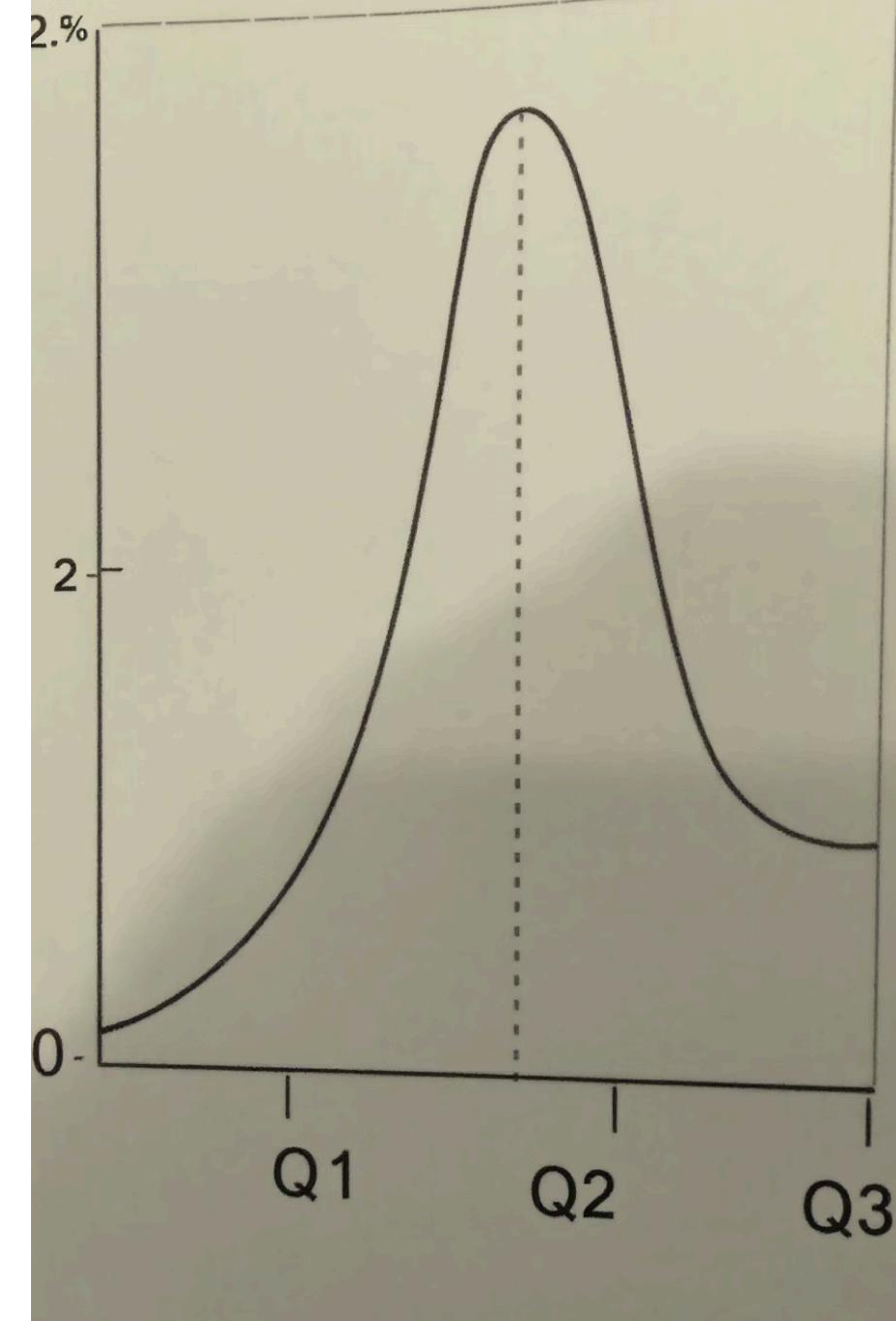
Percentage of values below a certain point in a dataset



Quartile

Values that divide a ranked dataset into four equal parts

Measures of position help identify the relative standing of a data point within a dataset. They are essential for comparing values across different distributions and for understanding how data is distributed across its range.



Z-Score (Standard Score)

$$X - \mu$$

Deviation

Distance from the mean

$$\sigma$$

Standard Deviation

Measure of spread

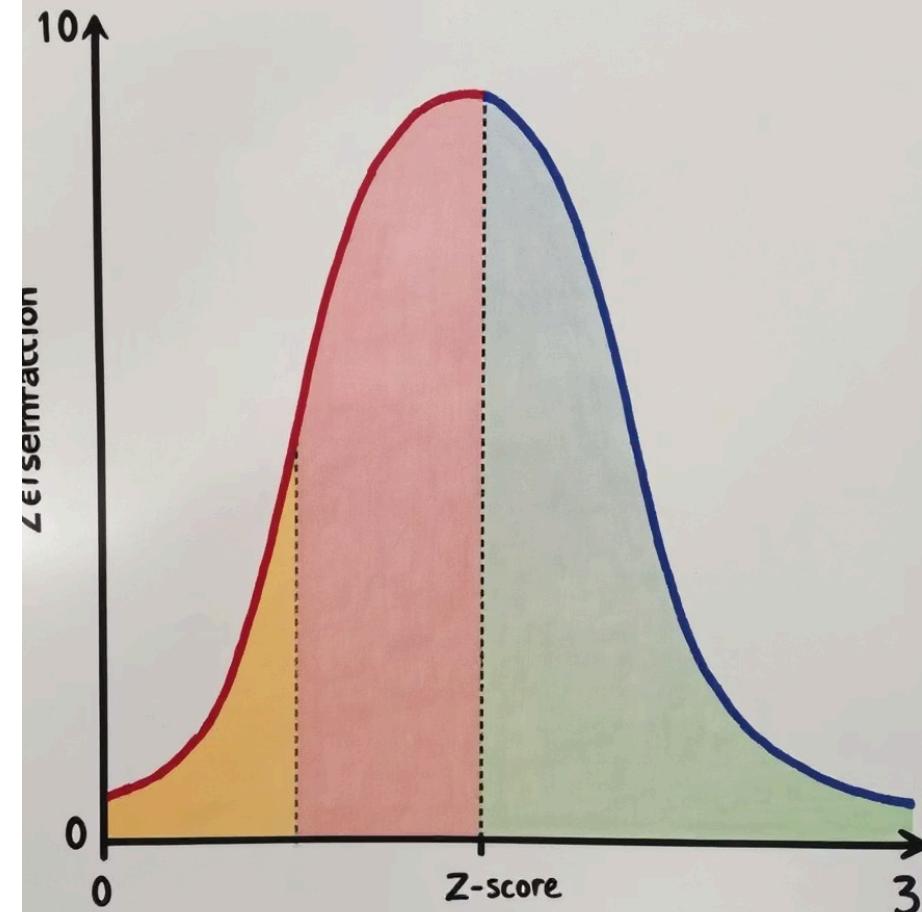
$$Z$$

Standardized Score

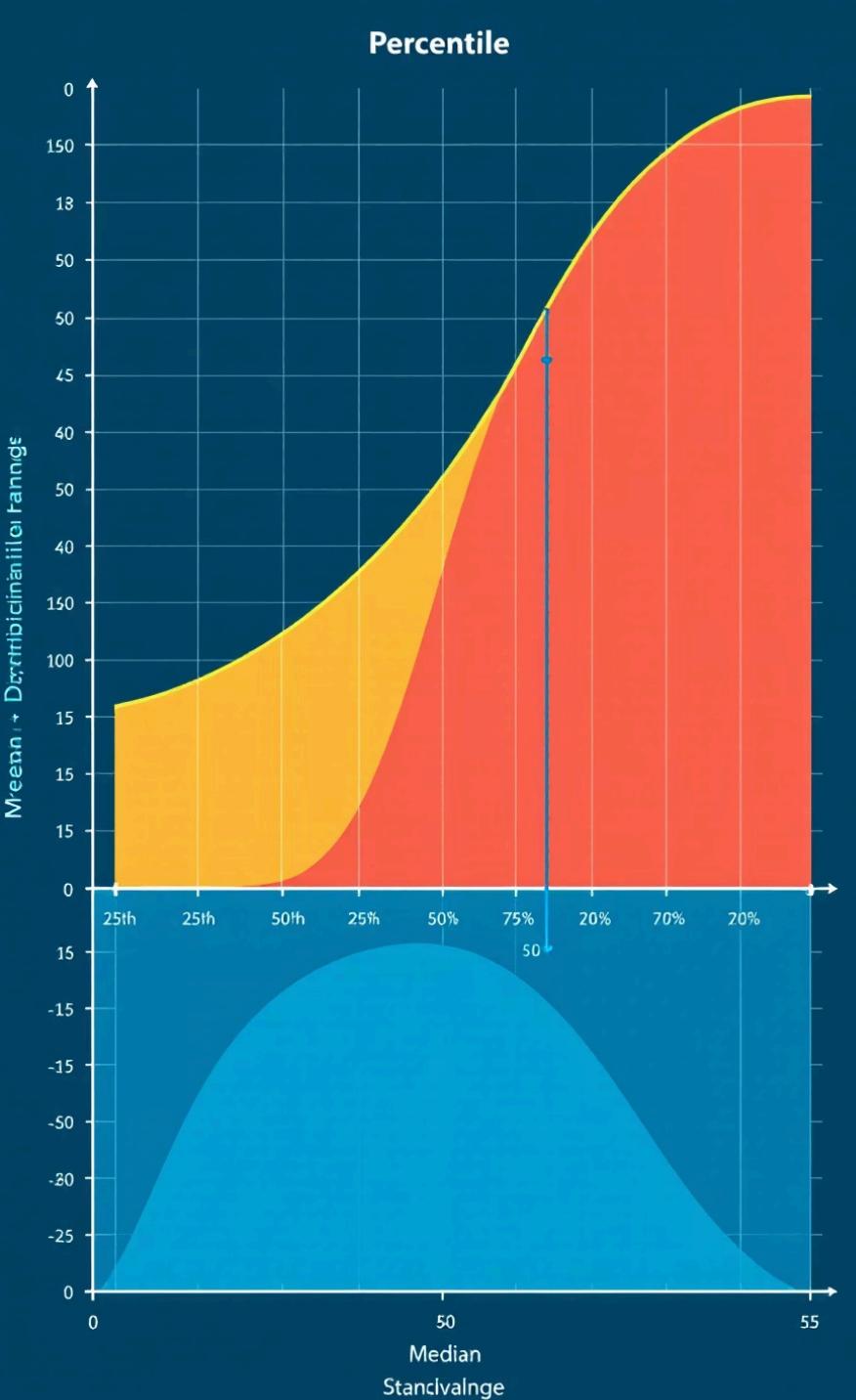
Deviation in SD units

A Z-score tells you how many standard deviations a data point is from the mean. For example, if a test has a mean score of 70 with a standard deviation of 10, and John scores 85, his Z-score would be $Z = (85-70)/10 = 1.5$. This means John scored 1.5 standard deviations above the mean. Z-scores are useful for comparing values from different distributions and identifying outliers (typically $Z > 2$ or $Z < -2$).

Z-Score Distribution



68% within ± 1 standard deviations
95% within ± 2 standard deviations
99.7% within ± 3 standard deviations



Percentiles



Rank Data

Sort values from lowest to highest

Calculate Position

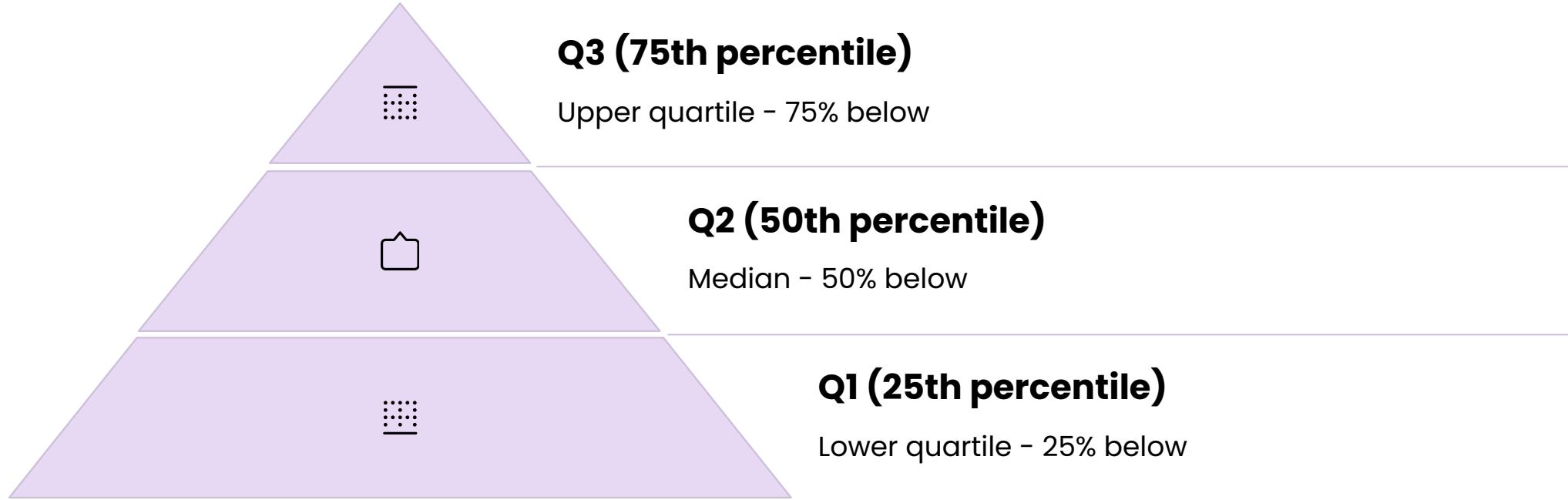
Determine percentage below each value

Interpret Results

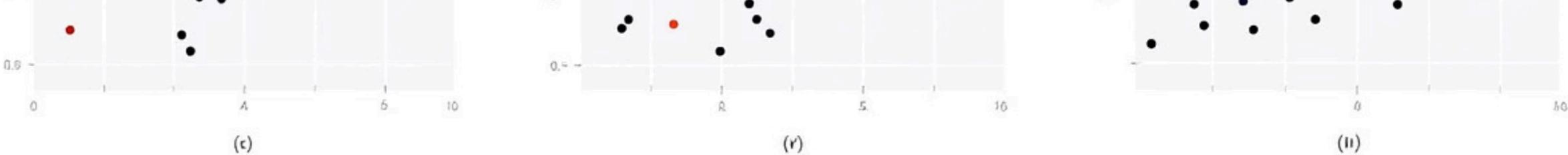
Understand relative standing

A percentile indicates the percentage of values below a certain point in a dataset. If you are in the 90th percentile in a test, it means 90% of the scores are below yours. For example, in a class with scores [40, 50, 55, 60, 65, 70, 75, 80, 85, 90], the 90th percentile would be close to the 9th score: 85. Percentiles are widely used in entrance exams, growth charts, and performance evaluations.

Quartiles



Quartiles divide a ranked dataset into four equal parts. For the dataset [10, 20, 30, 40, 50, 60, 70, 80, 90], $Q1 = 30$ (25% point), $Q2 = 50$ (Median), and $Q3 = 70$ (75% point). The Interquartile Range (IQR) = $Q3 - Q1 = 70 - 30 = 40$, which is used to detect outliers. Quartiles give insights into data spread and central grouping, and are essential for creating box plots.



Correlation Coefficient

$r=1$

Perfect Positive

Variables increase together

$r=0$

No Correlation

No linear relationship

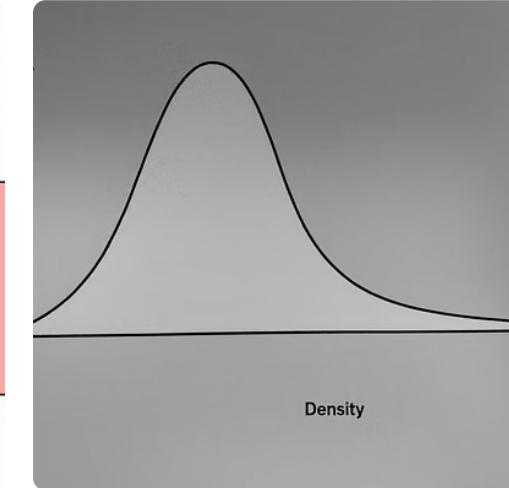
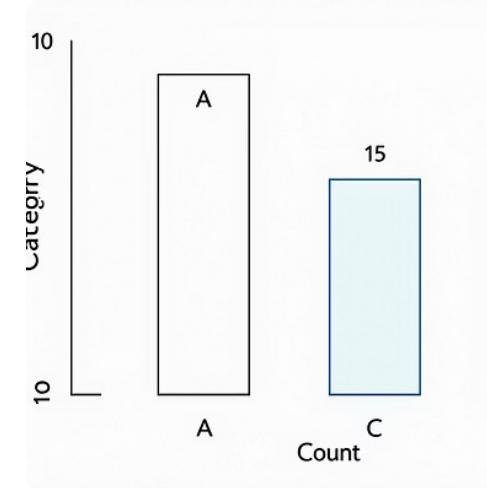
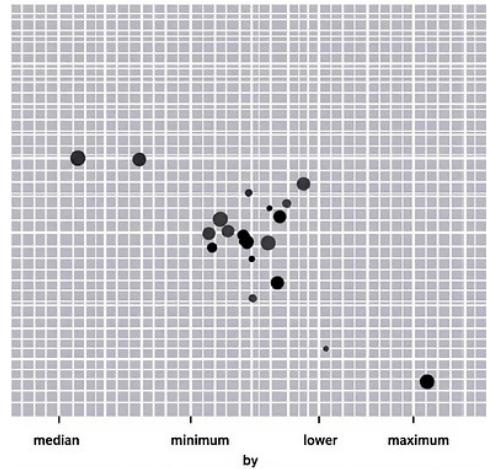
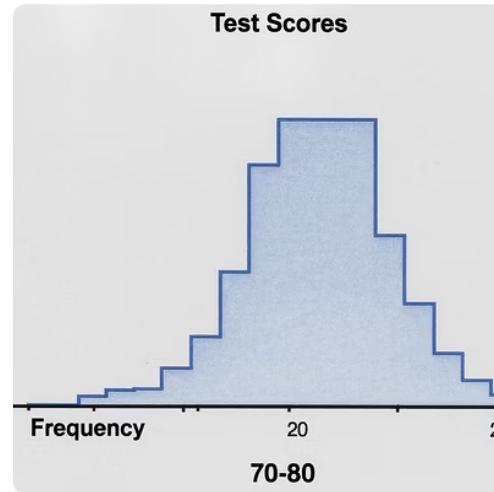
$r=-1$

Perfect Negative

One increases as other decreases

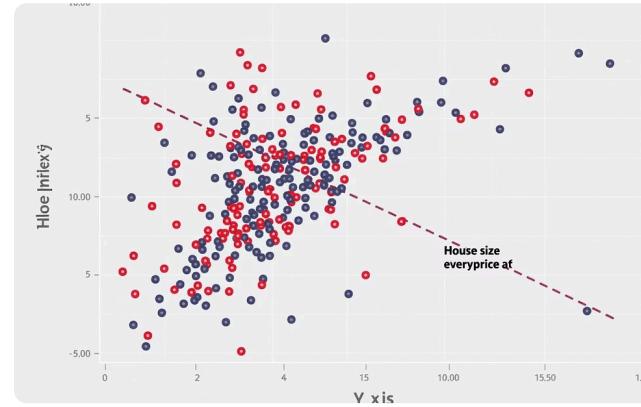
The correlation coefficient (r) measures the strength and direction of the linear relationship between two variables. It ranges from -1 to 1. For example, study time vs exam scores might show a strong positive correlation ($r \approx 0.8$), meaning more study time generally leads to higher scores. The correlation formula accounts for how variables move together relative to their means.

Data Visualization: Univariate Plots



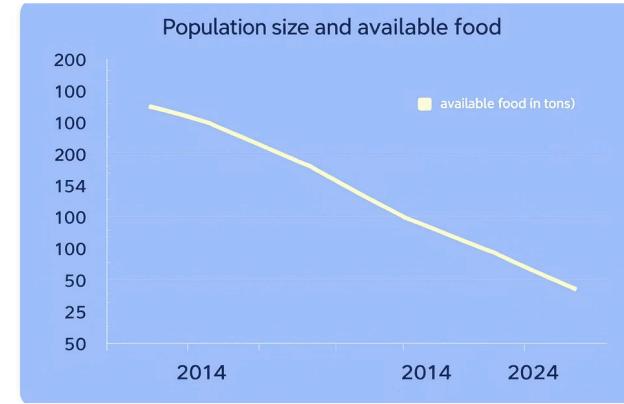
Univariate plots visualize the distribution of a single variable. Histograms show frequency distributions by dividing data into bins. Box plots display the median, quartiles, and potential outliers. Bar charts are used for categorical data to show counts or frequencies. Density plots provide a smooth estimate of the distribution. These visualizations help identify central tendency, spread, and shape of the data distribution.

Data Visualization: Bivariate Plots



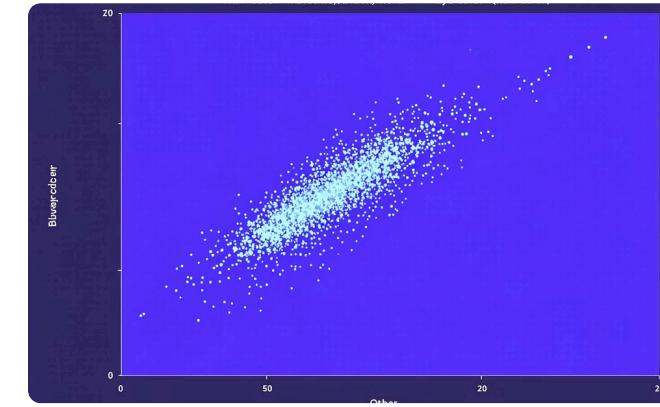
Scatter Plot

Shows relationship between two continuous variables, revealing patterns, trends, and potential correlations.



Line Chart

Displays trends between two variables, often used when one variable is time-based.



Heatmap

Visualizes the relationship between two categorical variables using color intensity to represent frequency or another metric.

Bivariate plots help identify relationships between two variables. They're essential for correlation analysis, trend identification, and pattern recognition in data science and statistical analysis.

Data Visualization: Multivariate Plots

3D Plots

Visualize relationships between three variables simultaneously, allowing for more complex pattern identification.

Example: 3D scatter plot showing height, weight, and age of individuals.

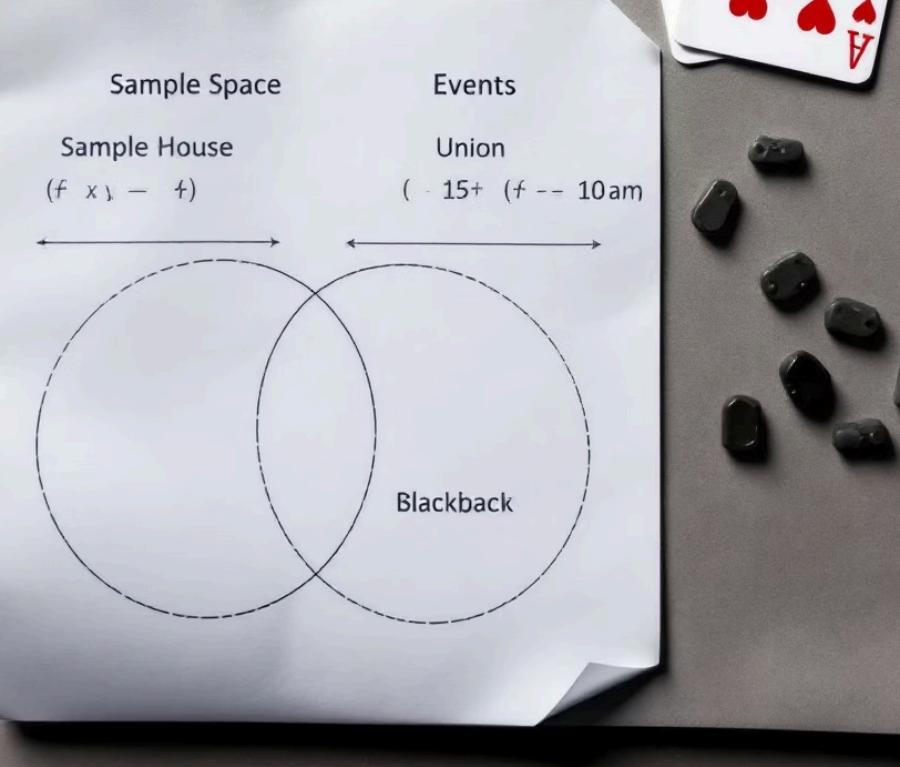
Multivariate plots are crucial for understanding complex relationships in datasets with many variables. They help data scientists identify patterns that might not be visible when examining variables in isolation or pairs. These visualizations are commonly used in exploratory data analysis before building predictive models.

Pair Plots

Matrix of scatter plots showing relationships between multiple pairs of variables, with histograms on the diagonal.

Example: Visualizing relationships between multiple features in a machine learning dataset.

Basic Probability Concepts



Sample Space

Set of all possible outcomes of an experiment.

Example: When rolling a die, $S = \{1, 2, 3, 4, 5, 6\}$

Event

Subset of the sample space.

Example: Rolling an even number, $E = \{2, 4, 6\}$

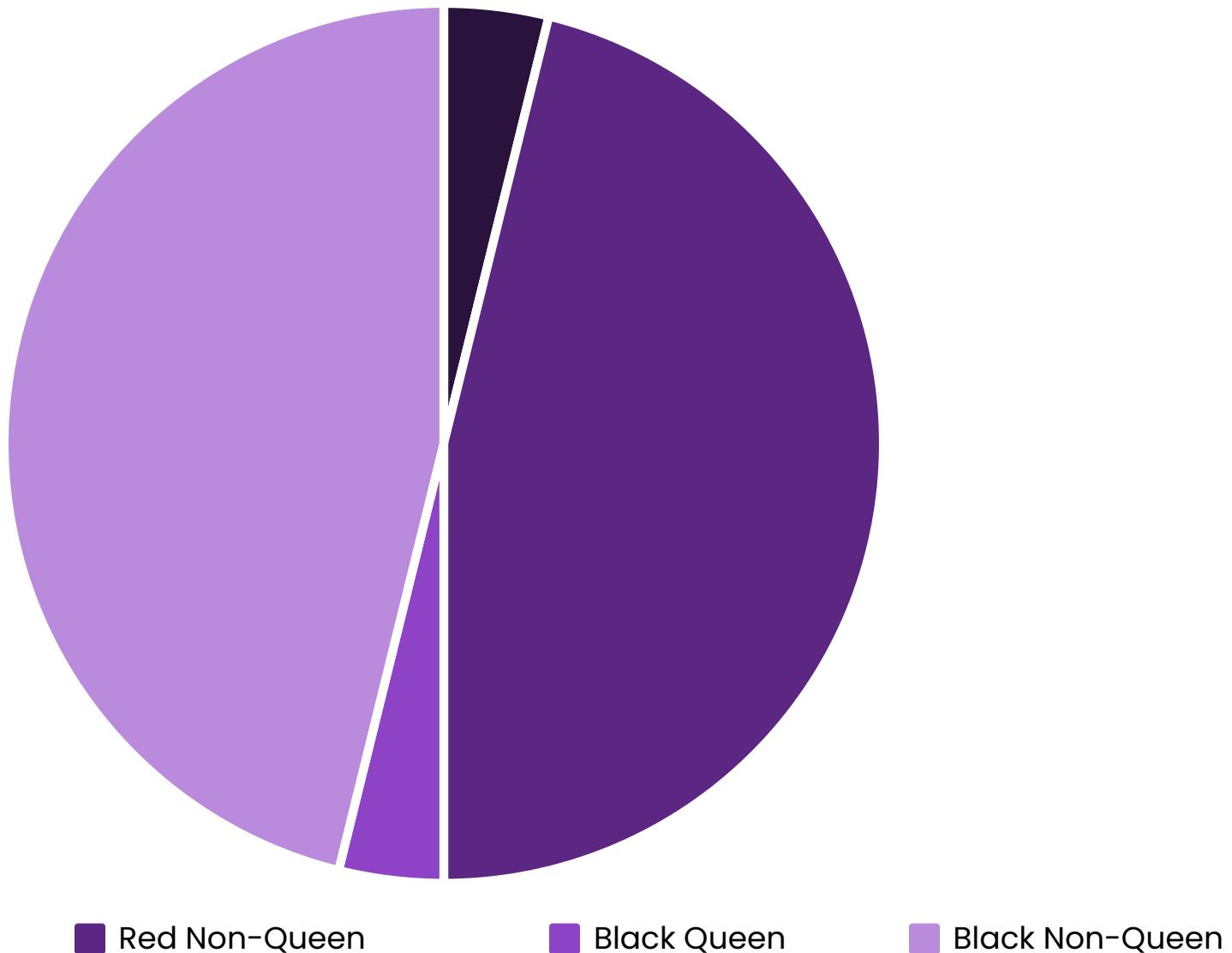
Probability

$P(E) = \text{Number of favorable outcomes} / \text{Total number of outcomes}$

Example: $P(\text{even number}) = 3/6 = 0.5$

Probability quantifies the likelihood of events occurring. It ranges from 0 (impossible) to 1 (certain). Understanding basic probability is essential for statistical inference, risk assessment, and decision-making under uncertainty.

Joint Probability



Joint probability $P(A \cap B)$ is the probability of two events occurring together. For example, when drawing a card from a standard deck, the probability of getting a card that is both red and a queen is $P(\text{Red} \cap \text{Queen}) = 2/52 = 1/26$. This represents the intersection of the set of red cards (26 cards) and the set of queens (4 cards), which contains 2 cards (the queen of hearts and the queen of diamonds).

Joint probability is essential for understanding how events interact and is the foundation for more complex probability concepts.

Marginal Probability

	Queen	Non-Queen	Marginal (Color)
Red	2/52	24/52	$26/52 = 1/2$
Black	2/52	24/52	$26/52 = 1/2$
Marginal (Card Type)	$4/52 = 1/13$	$48/52 = 12/13$	1

Marginal probability $P(A)$ is the probability of an event regardless of other events. In a deck of cards, the marginal probability of drawing a red card is $P(\text{Red}) = 26/52 = 1/2$, regardless of whether it's a queen or not. Similarly, the marginal probability of drawing a queen is $P(\text{Queen}) = 4/52 = 1/13$, regardless of its color.

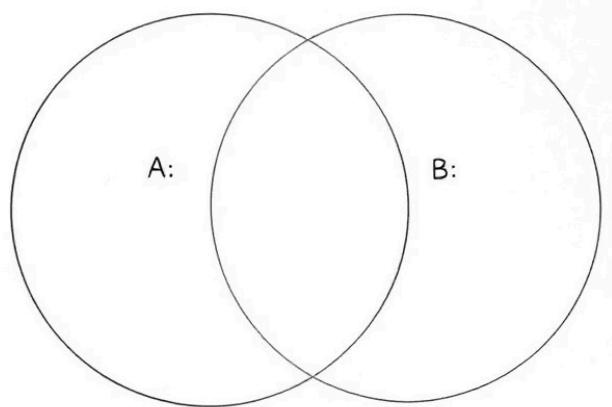
Marginal probabilities can be calculated by summing joint probabilities across all possible values of the other variable.

Marginal Probability



Conditional Probability

D: IN. Hearts



und Playing card Chances

DIAGRAM 2:
Conditional probability

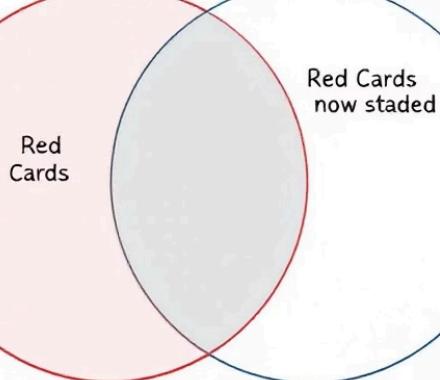


Diagram
the shape
for the a
condition
Hearts give
Red Car



Question

What is the probability of event A given that event B has occurred?

Formula

$$P(A|B) = P(A \cap B) / P(B)$$

Example

$$P(\text{Queen}|\text{Red}) = P(\text{Red} \cap \text{Queen}) / P(\text{Red}) = (2/52) / (26/52) = 2/26 = 1/13$$

Conditional probability measures the likelihood of an event occurring given that another event has already occurred. For example, if we know a card drawn from a deck is red, the probability it's a queen is $P(\text{Queen}|\text{Red}) = 1/13$. This differs from the unconditional probability of drawing a queen, which is $1/13$ regardless of color.

Probability Distributions Overview

Discrete Probability Distributions

Used for countable outcomes like the number of successes or occurrences.

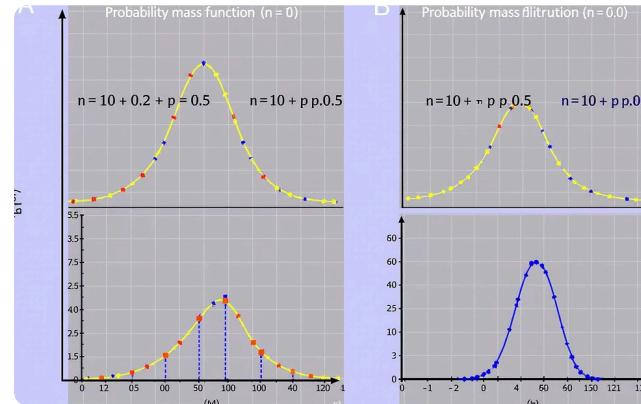
- Binomial: Number of successes in fixed trials
- Poisson: Number of events in fixed time/space
- Geometric: Trials until first success
- Negative Binomial: Trials until r successes

Continuous Probability Distributions

Used for measurable variables that can take any value within a range.

- Normal: Bell-shaped, symmetric distribution
- Uniform: Equal probability across range
- Exponential: Time between events
- Beta: Proportions and probabilities

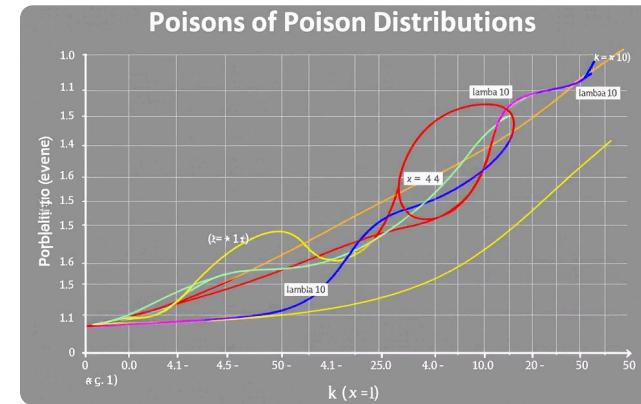
Discrete Probability Distributions



Binomial Distribution

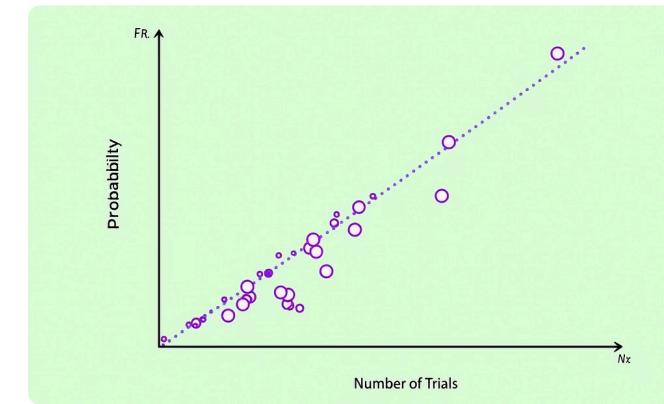
Models the number of successes in a fixed number of independent trials, each with the same probability of success.

Discrete probability distributions assign probabilities to distinct, countable outcomes. They are essential for modeling events like the number of defects in manufacturing, customer arrivals, or success/failure scenarios.



Poisson Distribution

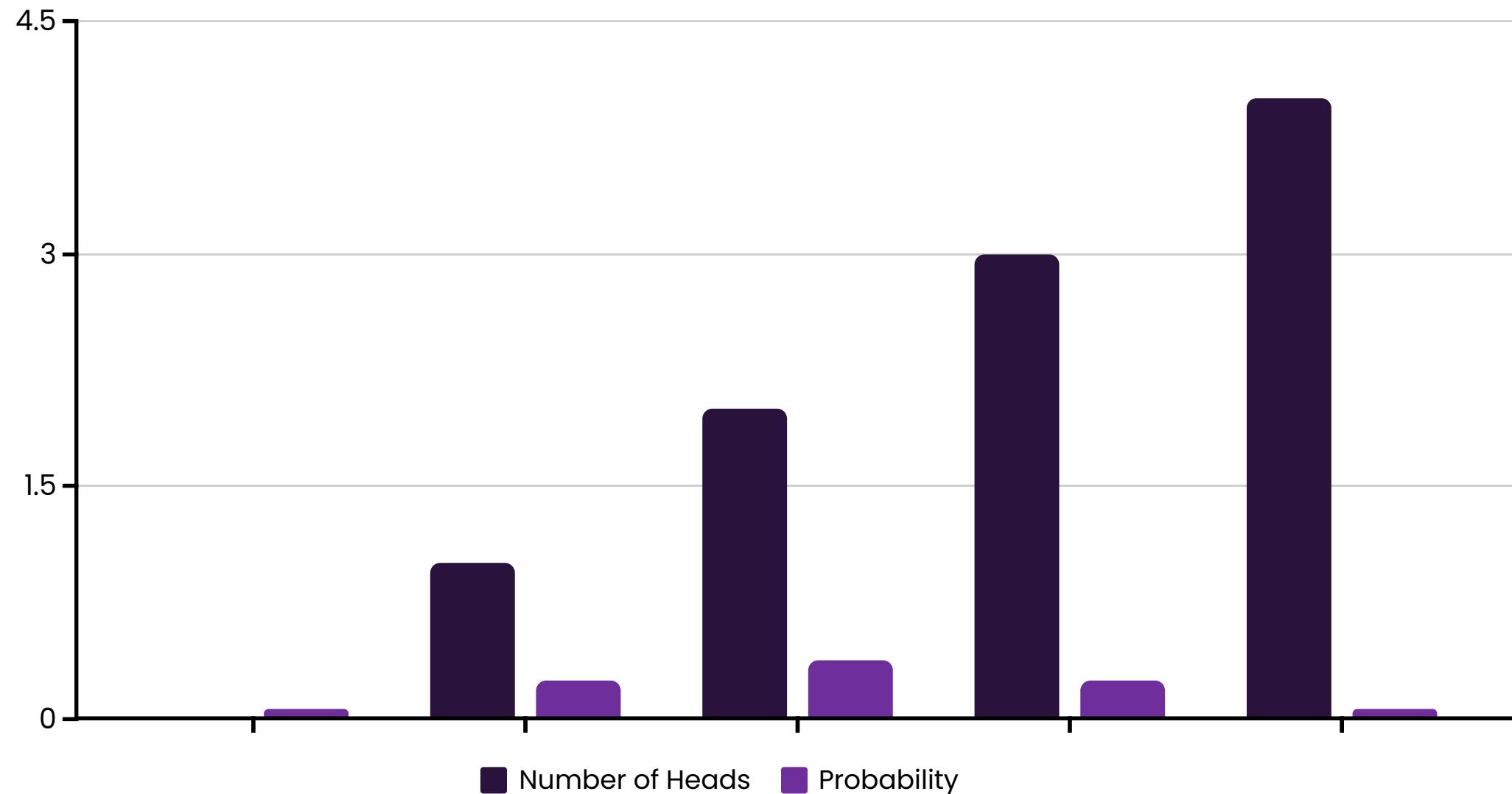
Models the number of events occurring in a fixed interval of time or space, assuming events occur independently.



Geometric Distribution

Models the number of trials needed to achieve the first success in a sequence of independent trials.

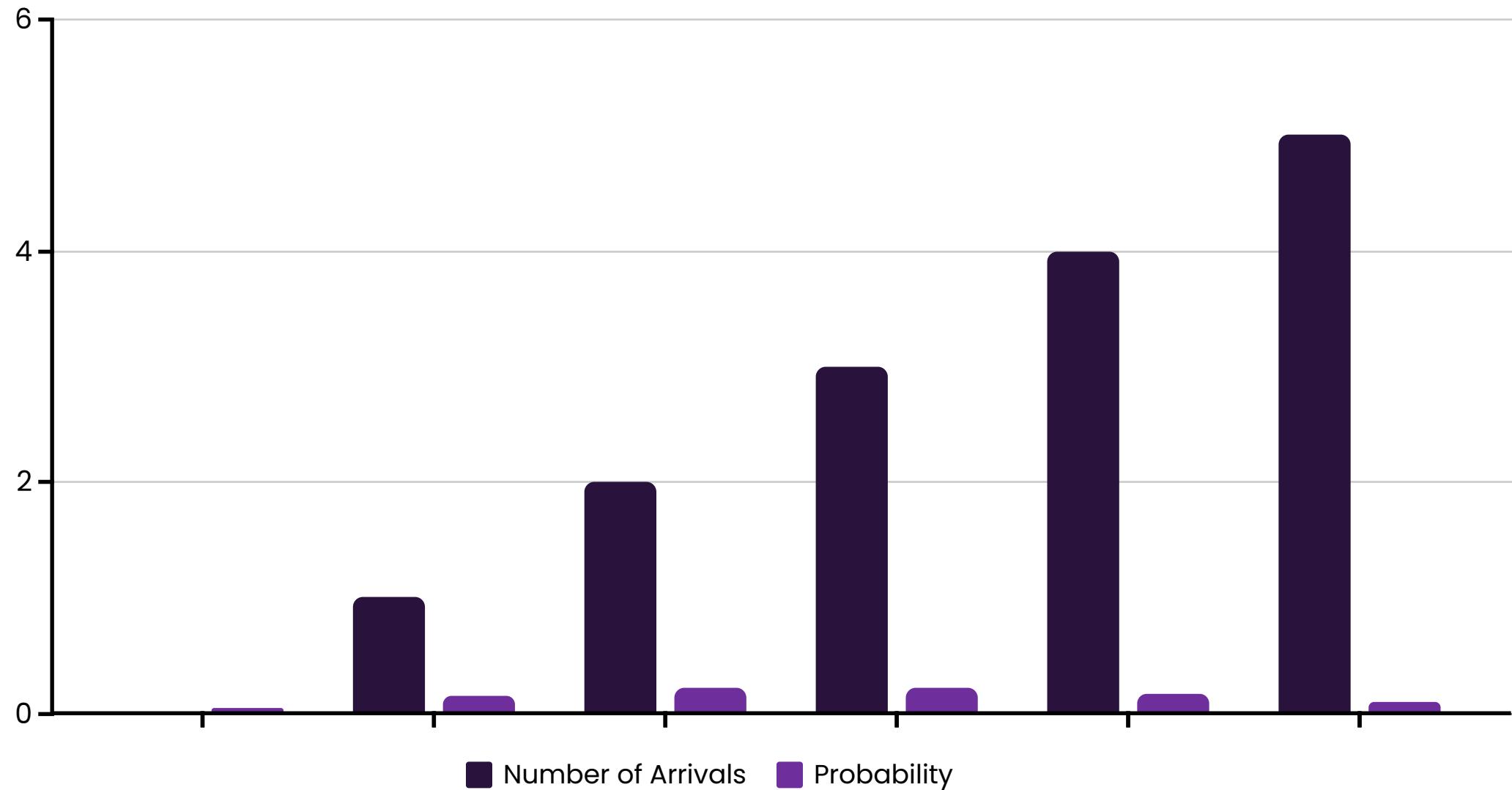
Binomial Distribution Example



The binomial distribution models the number of successes in a fixed number of independent trials. For example, when flipping a fair coin 4 times, the probability of getting exactly 2 heads is given by the binomial probability mass function: $P(X=2) = (4 \text{ choose } 2) \times (0.5)^2 \times (0.5)^2 = 6 \times 0.25 \times 0.25 = 0.375$.

The binomial distribution is widely used in quality control, polling, and any scenario involving success/failure outcomes over multiple trials.

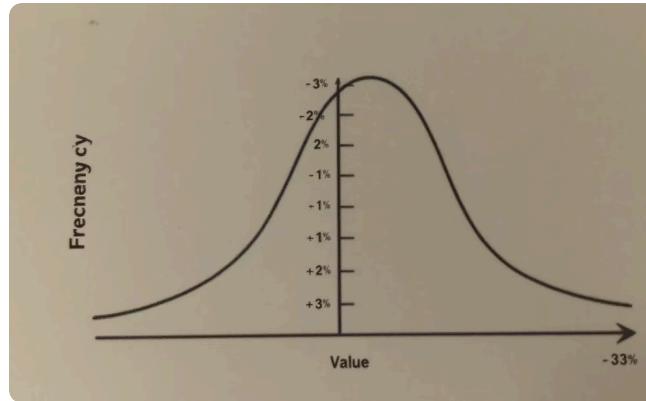
Poisson Distribution Example



The Poisson distribution models the number of events occurring in a fixed interval, assuming events occur independently. For example, if a call center receives an average of 3 calls per hour, the probability of receiving exactly 2 calls in the next hour is: $P(X=2) = e^{-3} \times 3^2 / 2! \approx 0.224$.

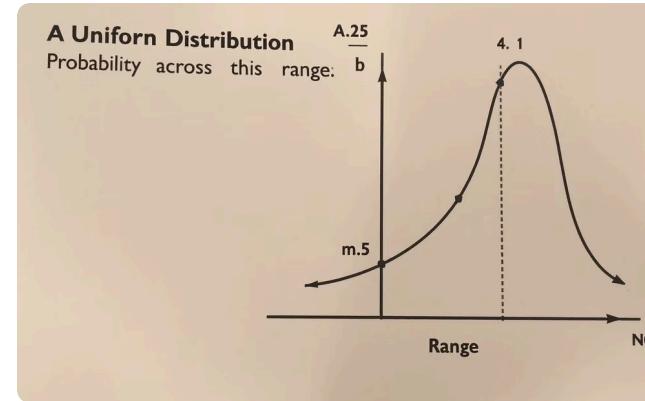
This distribution is commonly used for modeling rare events like equipment failures, website traffic spikes, or customer arrivals.

Continuous Probability Distributions



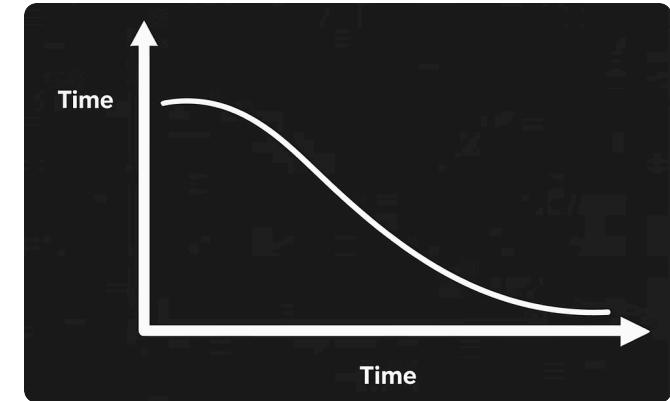
Normal Distribution

Bell-shaped, symmetric distribution defined by mean and standard deviation. Central to statistical theory due to the Central Limit Theorem.



Uniform Distribution

Equal probability across all values in a range. Used when any value in an interval is equally likely to occur.



Exponential Distribution

Models the time between events in a Poisson process. Characterized by a constant failure/arrival rate.

Continuous probability distributions assign probabilities to ranges of values rather than individual points. They are essential for modeling measurements like height, weight, time, and other variables that can take any value within a range.

Normal Distribution Properties

68%

Within 1σ

Percentage within one standard deviation of mean

95%

Within 2σ

Percentage within two standard deviations of mean

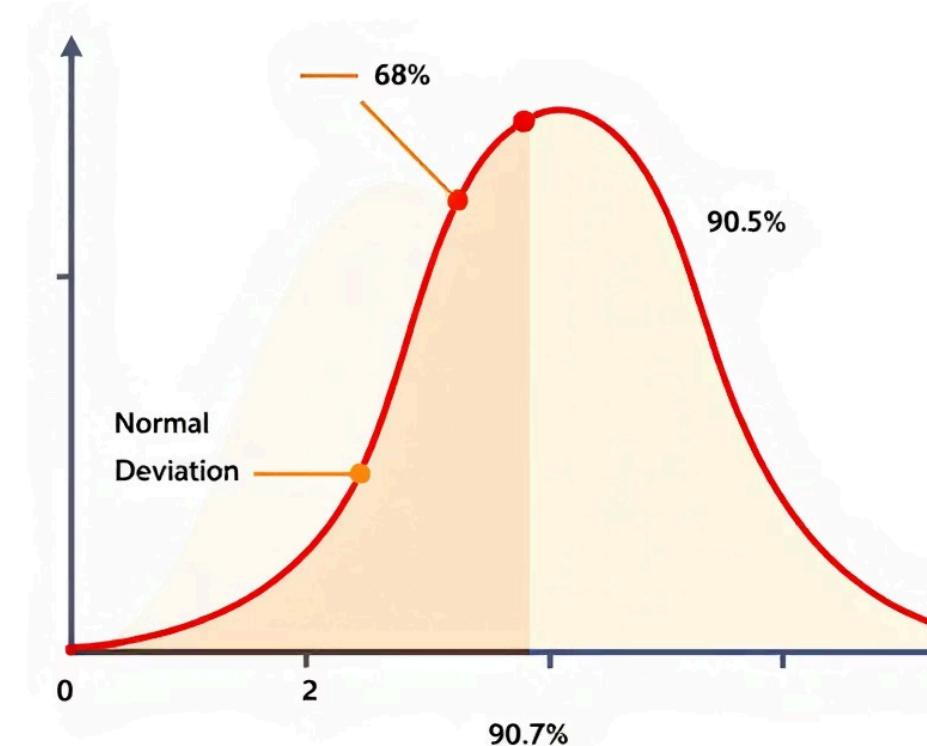
99.7%

Within 3σ

Percentage within three standard deviations of mean

The normal distribution is defined by its mean (μ) and standard deviation (σ). It's symmetric around the mean, with the familiar bell-shaped curve. The empirical rule (68-95-99.7 rule) helps interpret normal distributions: approximately 68% of values fall within 1σ of the mean, 95% within 2σ , and 99.7% within 3σ . This distribution is foundational in statistics due to the Central Limit Theorem.

A Normal Distribution Curve



Bayesian Probability

Prior Probability

Initial belief before evidence

New Evidence

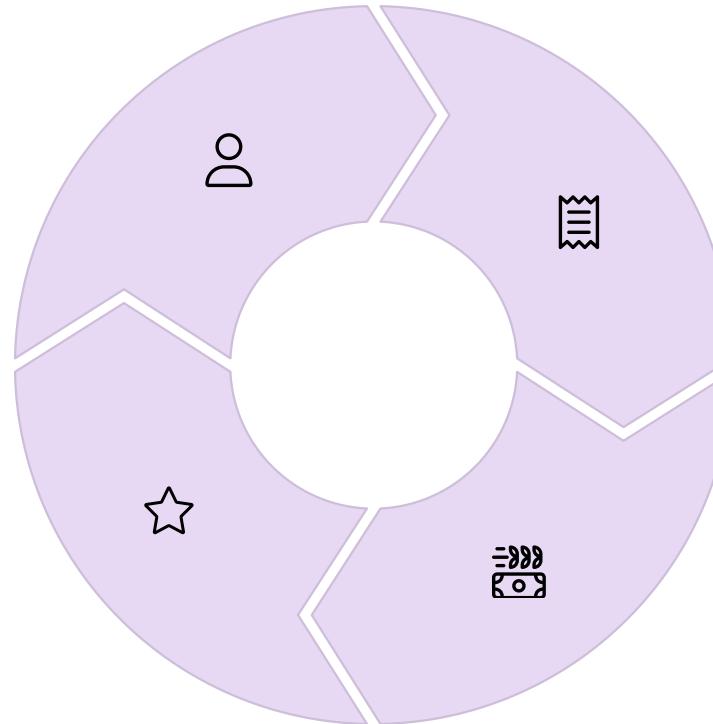
Data that updates our belief

Posterior Probability

Updated belief after evidence

Likelihood

Probability of evidence given hypothesis



Bayesian probability updates the likelihood of an event based on new evidence. It's expressed through Bayes' Theorem: $P(A|B) = P(B|A) \times P(A) / P(B)$. This approach allows for incorporating prior knowledge and updating beliefs as new information becomes available, making it powerful for decision-making under uncertainty.

Bayes' Theorem Example: Medical Testing

Prior Probability

$P(\text{Disease}) = 0.01$ (1% of population has the disease)

Test Accuracy

$P(\text{Positive}|\text{Disease}) = 0.9$ (90% sensitivity)

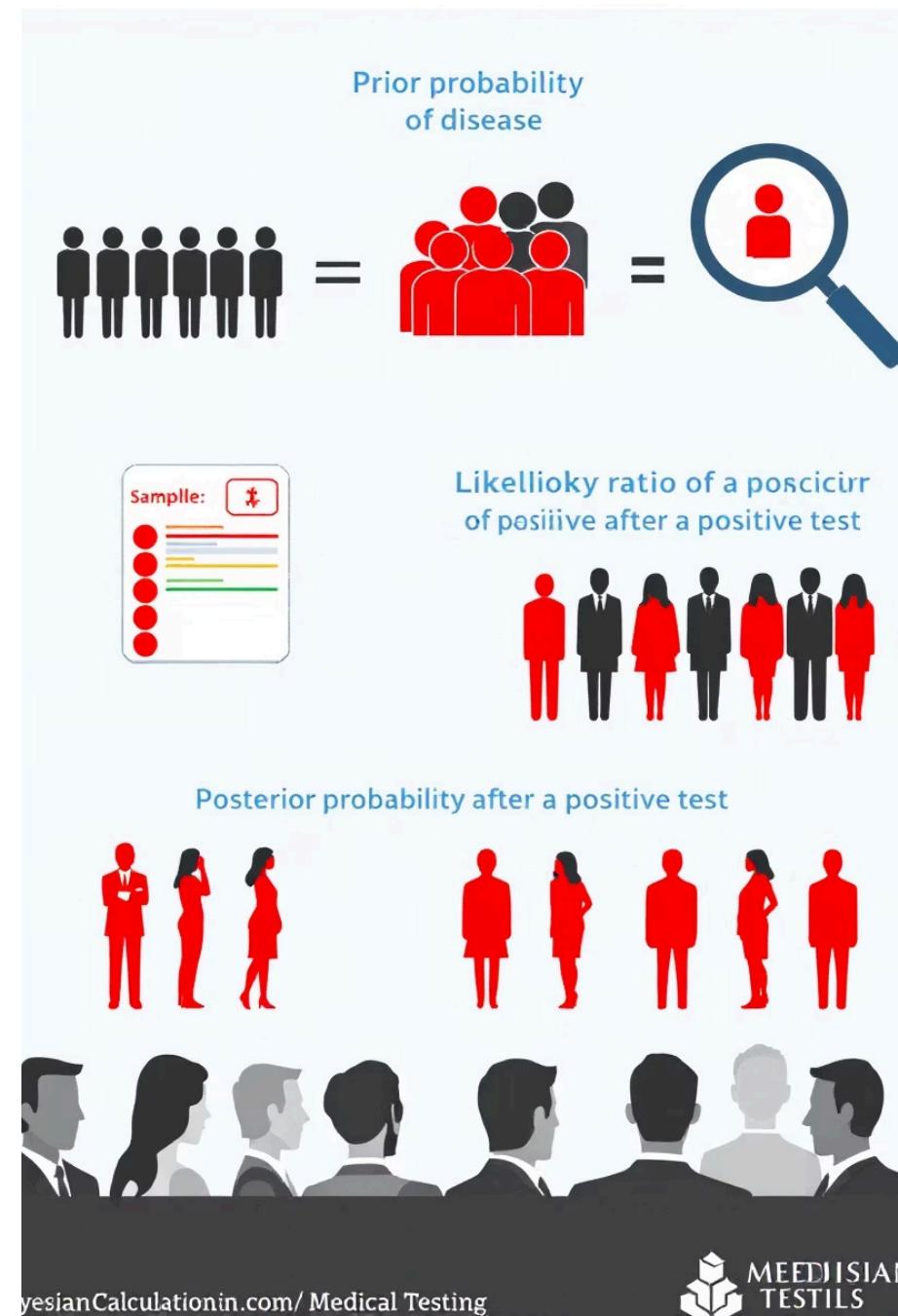
$P(\text{Positive}|\text{No Disease}) = 0.05$ (5% false positive rate)

Posterior Probability

$$\begin{aligned} P(\text{Disease}|\text{Positive}) &= P(\text{Positive}|\text{Disease}) \times P(\text{Disease}) / \\ &\quad P(\text{Positive}) \\ &= 0.9 \times 0.01 / [0.9 \times 0.01 + 0.05 \times 0.99] \approx 0.153 \end{aligned}$$

Despite a positive test result, there's only about a 15.3% chance the person actually has the disease. This counterintuitive result demonstrates the importance of considering base rates (prior probabilities) when interpreting test results, especially for rare conditions. This concept is crucial in medical diagnostics, spam filtering, and many other applications.

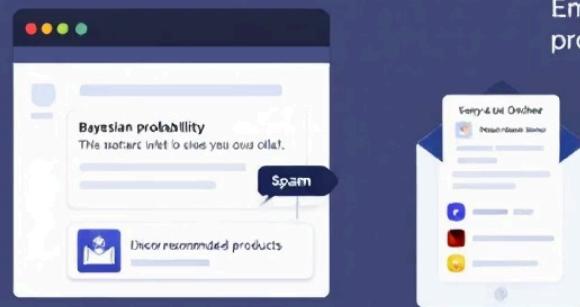
Bayesian Calculation in Medical Testing





Due to
applican email only
receivers and of the
cancer, and startling,
forecaster of the
Bayesian probability.

Email inbox



Email inbox and Bayesian probability

Shop recommended products



Gives upictions



Applications of Bayesian Probability



Medical Diagnosis

Updating disease probability based on test results and symptoms.



Spam Filtering

Classifying emails based on word probabilities and updating with user feedback.



Recommendation Systems

Predicting user preferences based on past behavior and similar users.

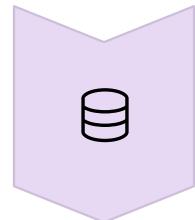


Weather Forecasting

Updating precipitation probabilities as new atmospheric data arrives.

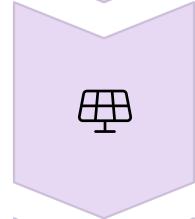


Data Science Workflow



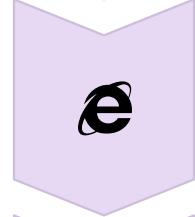
Data Collection

Gathering structured and unstructured data from various sources



Data Cleaning

Handling missing values, outliers, and inconsistencies



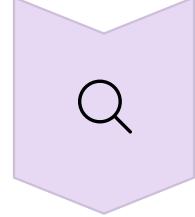
Exploratory Analysis

Understanding distributions, relationships, and patterns



Modeling

Building predictive or descriptive models using statistical techniques



Evaluation & Deployment

Testing model performance and implementing in production

Statistical Analysis in Machine Learning

Descriptive Statistics

Understanding data characteristics through measures of central tendency and dispersion.

- Feature engineering decisions
- Outlier detection
- Data transformation choices

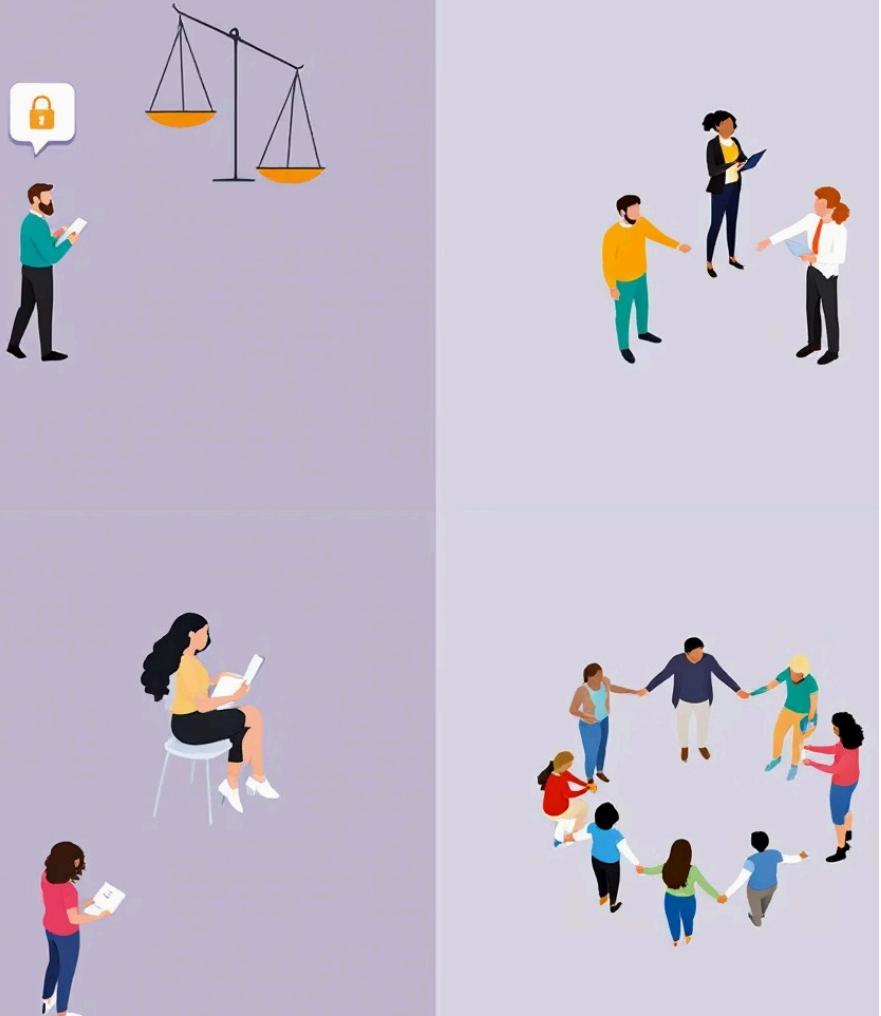
Statistical analysis forms the foundation of machine learning. Understanding data distributions helps in choosing appropriate algorithms. Correlation analysis identifies relevant features. Probability theory underlies many ML algorithms, from Naive Bayes to Neural Networks. Statistical tests help validate model performance and ensure results are significant.

Inferential Statistics

Drawing conclusions about populations from samples.

- Hypothesis testing for feature selection
- Confidence intervals for predictions
- A/B testing for model comparison

Ethical Considerations in Data Science



Ethical Considerations in Data Science



Data Privacy

Ensuring personal information is protected and used responsibly, with proper consent and anonymization.



Algorithmic Fairness

Preventing models from perpetuating or amplifying biases present in training data.



Transparency

Making algorithms and decision processes understandable to stakeholders and end users.



Social Impact

Considering the broader implications of data-driven systems on society and individuals.



Future Trends in Data Science



Automated Machine Learning

Democratizing AI through tools that automate model selection and hyperparameter tuning



Explainable AI

Making complex models more transparent and interpretable



Edge Computing

Moving data processing closer to the source for faster insights and reduced bandwidth

4

Quantum Computing

Leveraging quantum mechanics to solve previously intractable problems

The future of data science will be shaped by advances in computing power, algorithm development, and increasing data volumes. As AI becomes more integrated into critical systems, the focus on explainability, fairness, and robustness will intensify. Interdisciplinary approaches combining domain expertise with data science will lead to breakthroughs in fields from healthcare to climate science.