

# **Effective Big Data Management for Business analytics and decision**



# Topics (1)

Introduction to Big Data  
Basic understanding of business analytic  
Data format (structured vs unstructured)  
Fundamental of Data management



# Topics (2)

Data analysis techniques

Statistical analysis

Working with data analysis tool

Business intelligence

Workshop with Visualization data tool



# Introduction

Big Data

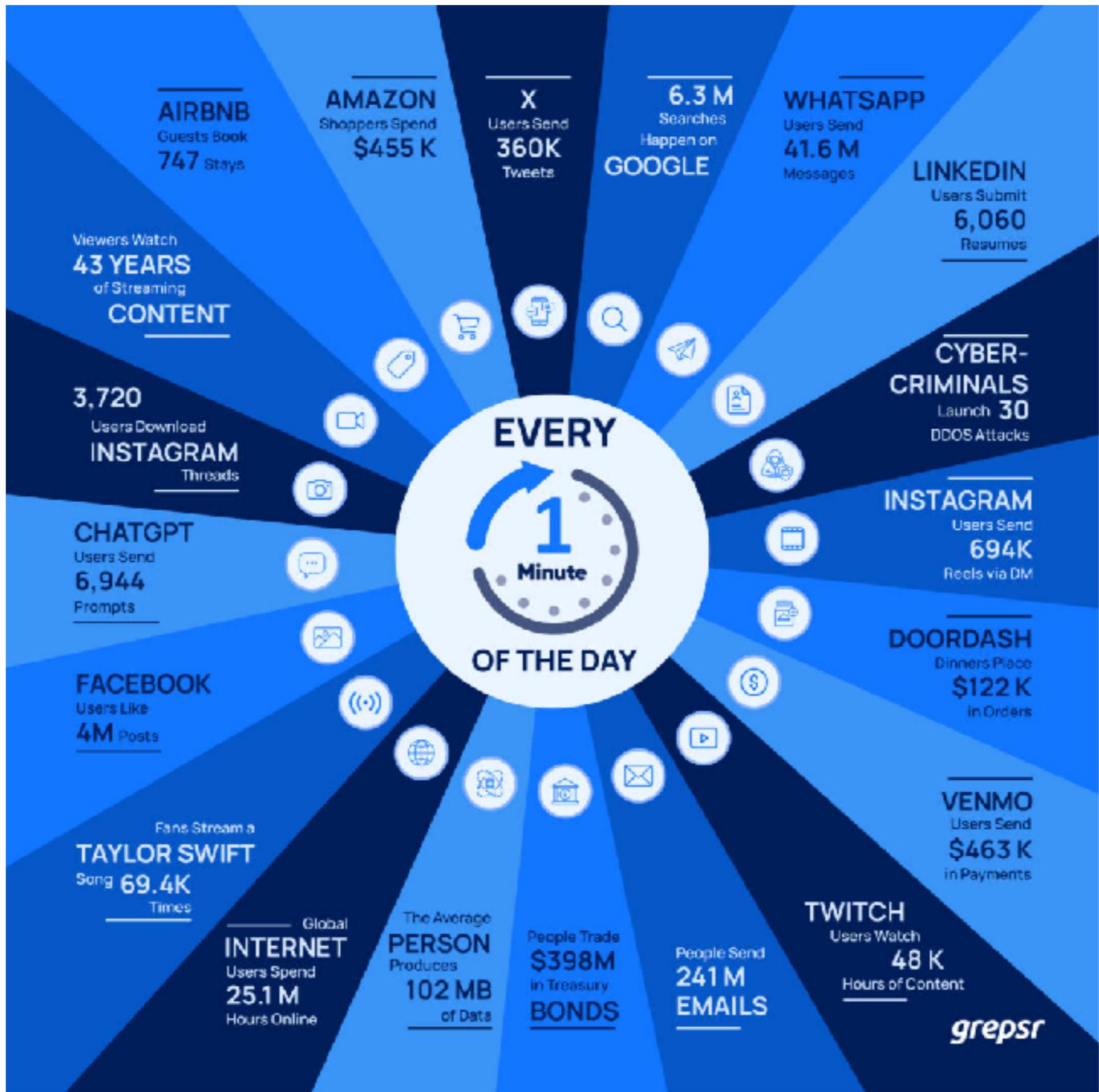
Business  
Analytics (BA)

Business  
Intelligence (BI)

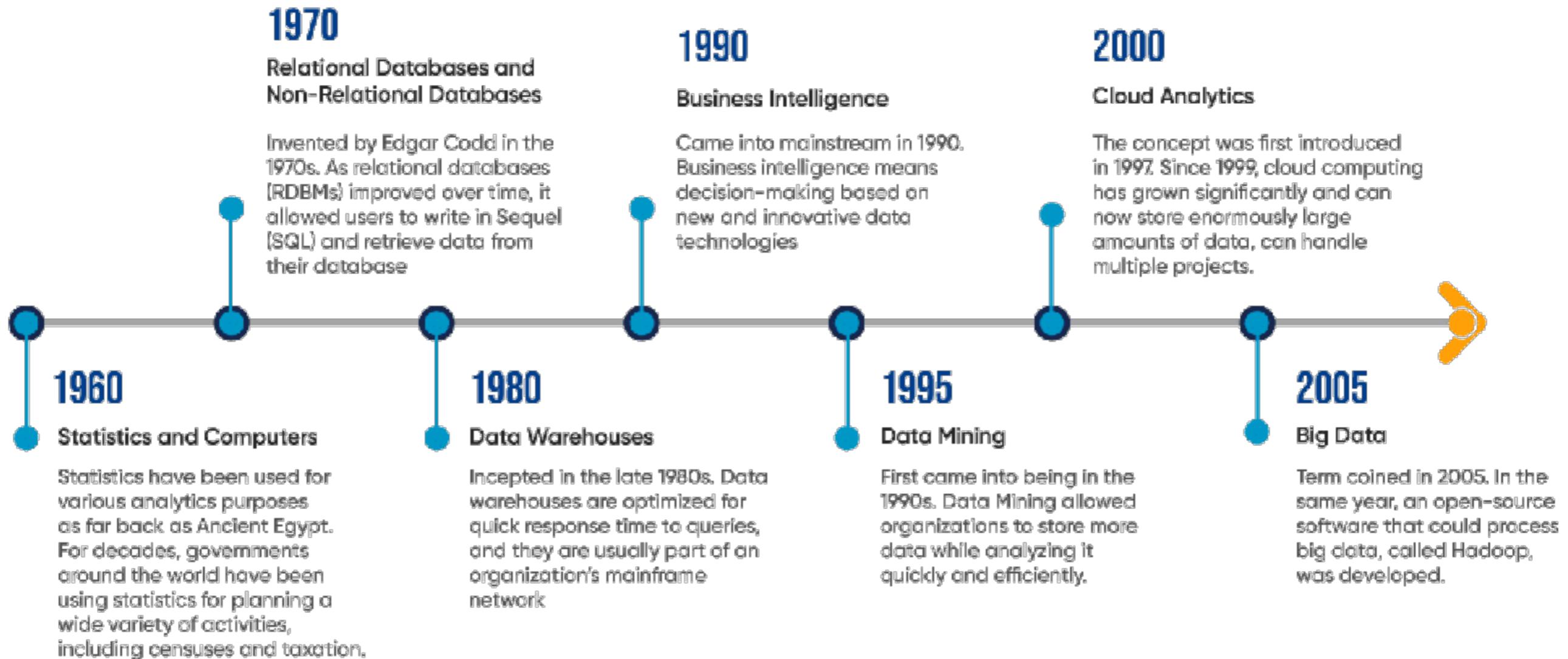


# Introduction to Big Data



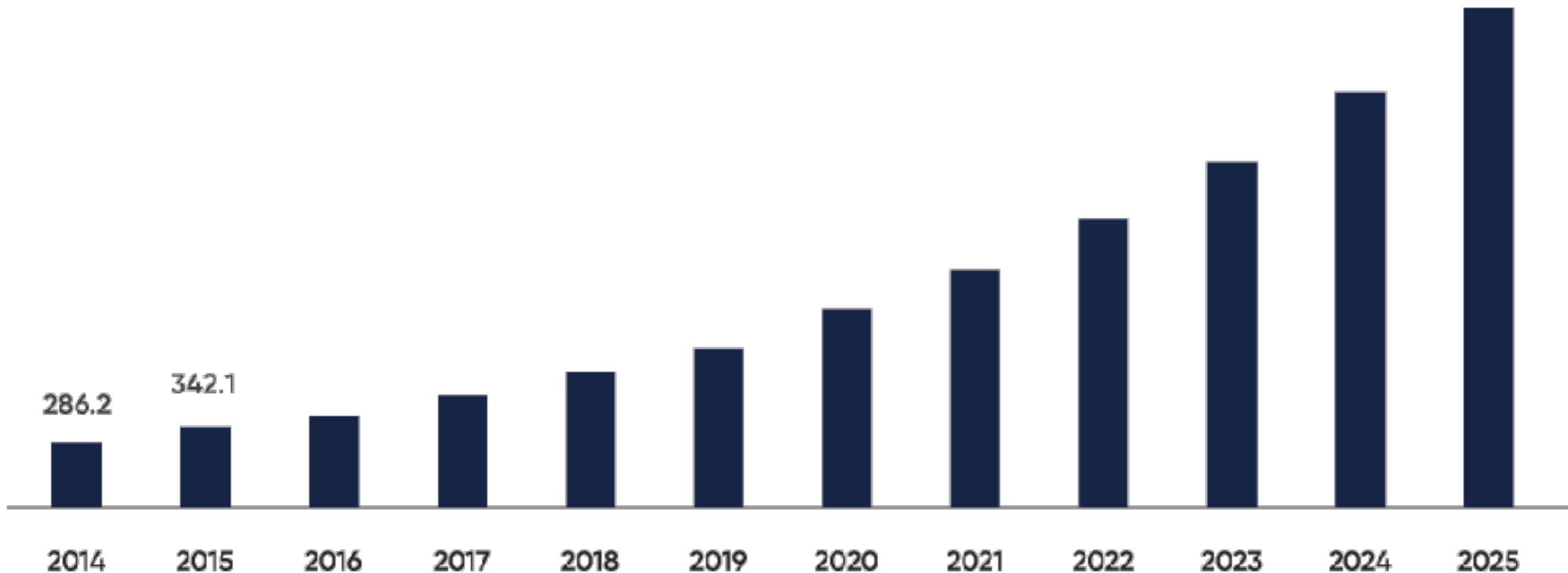


# Timeline



# Timeline

U.S. Data Analytics outsourcing market size, 2014 - 2025 (USD Million)



# Big Data

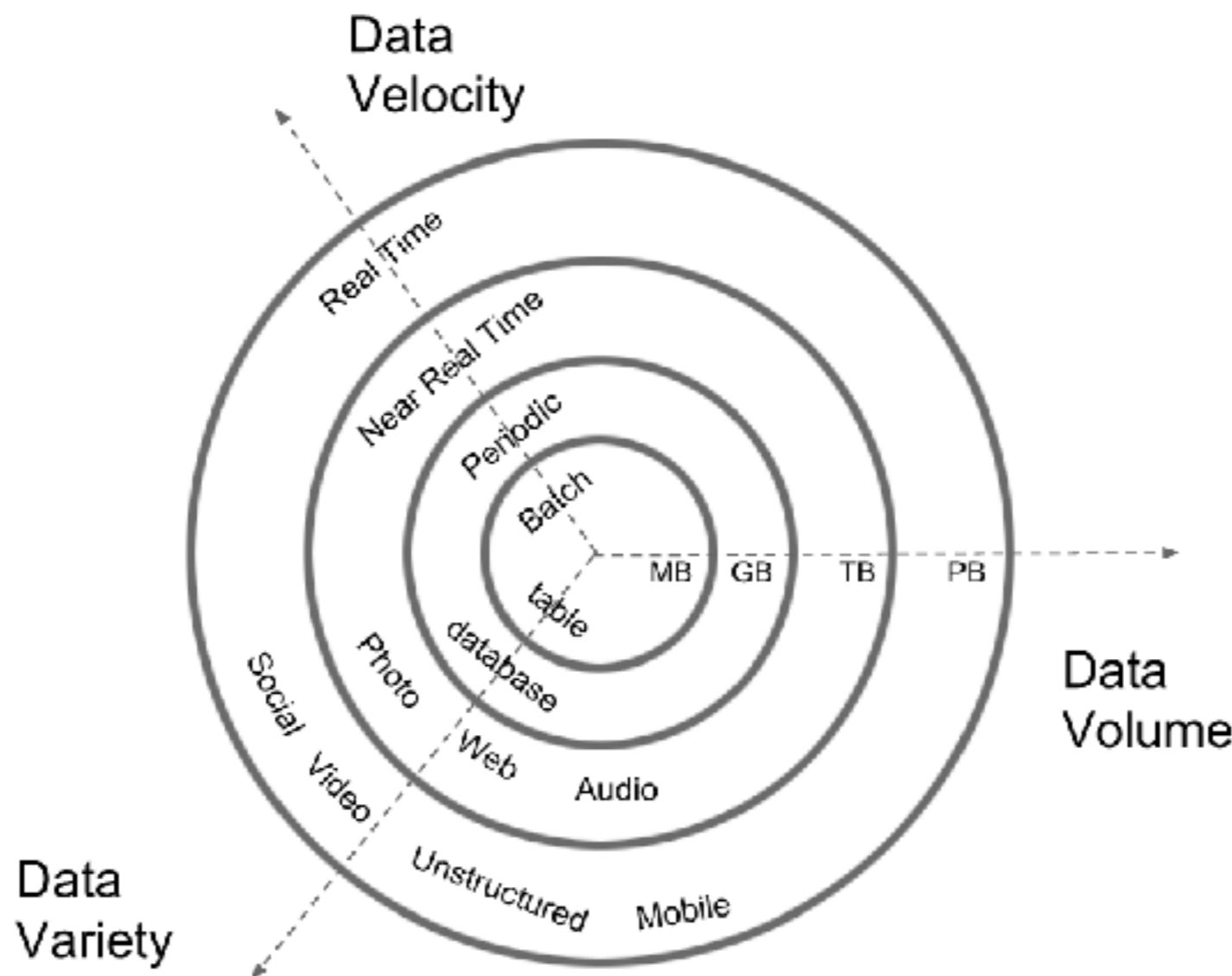
Refer to the vast amount of structured and unstructured data generated at high velocity, variety, and volume

Data is too complex to processed by traditional databases and tools

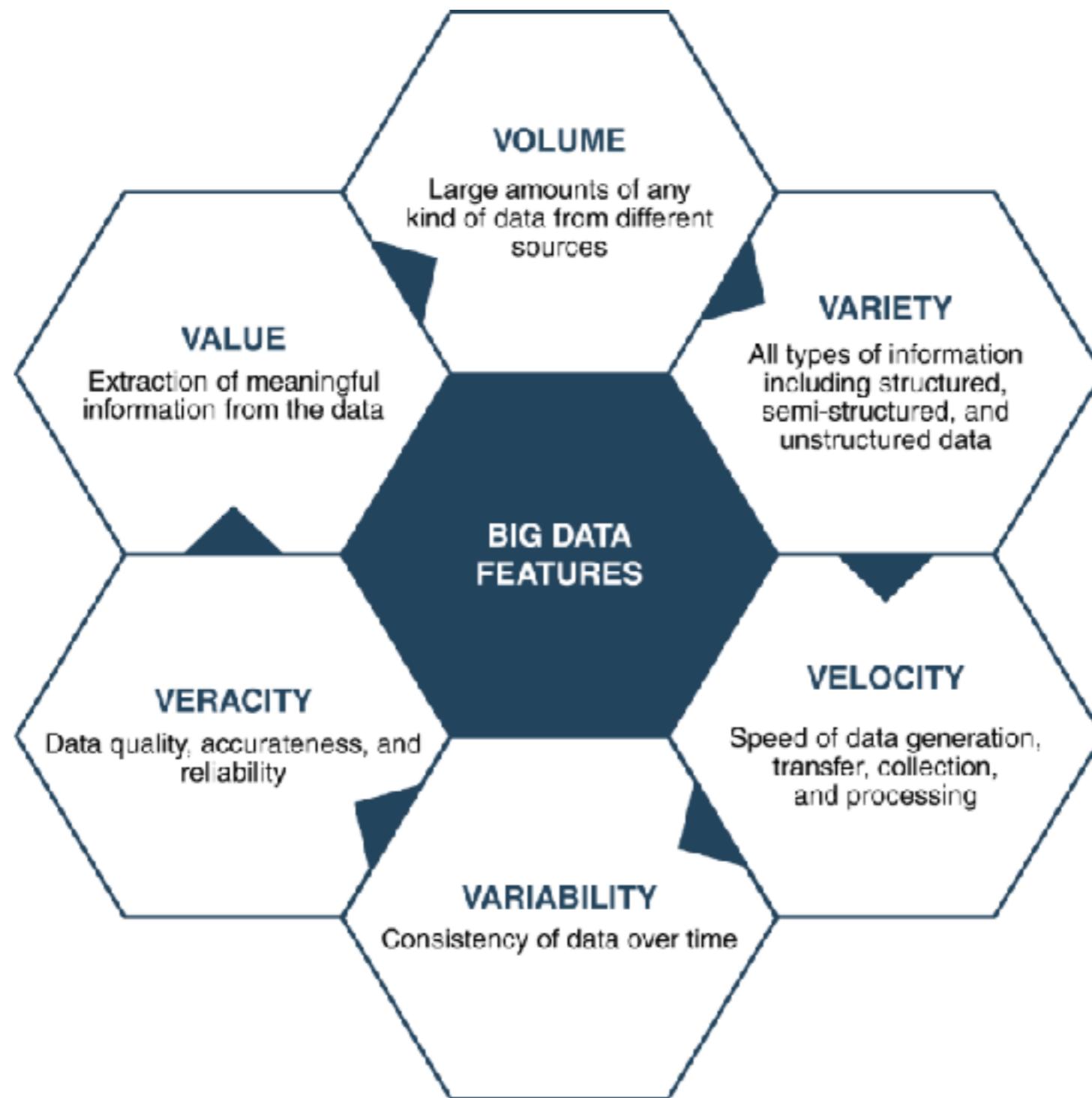


# Characteristics for Big Data

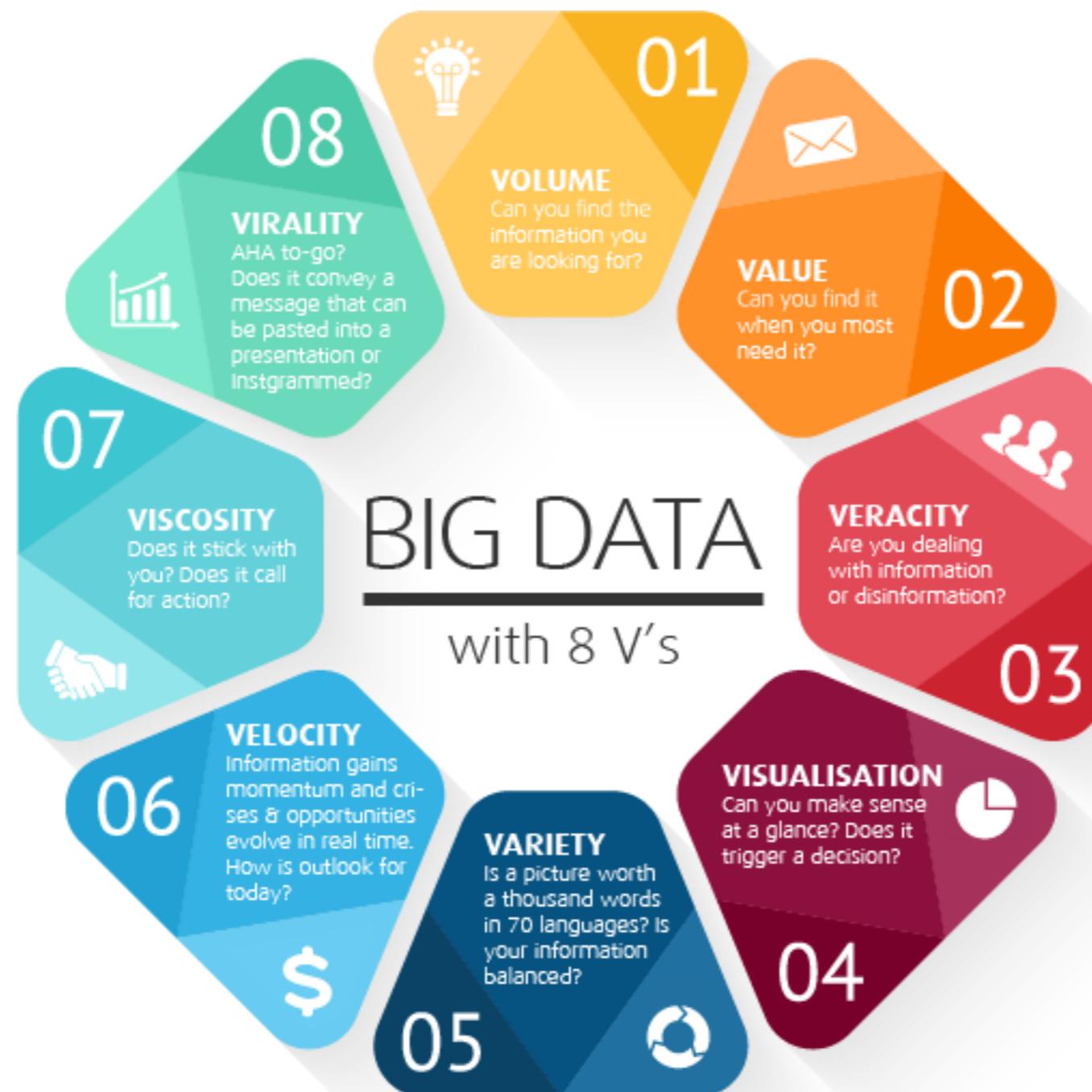
3 V's !!!



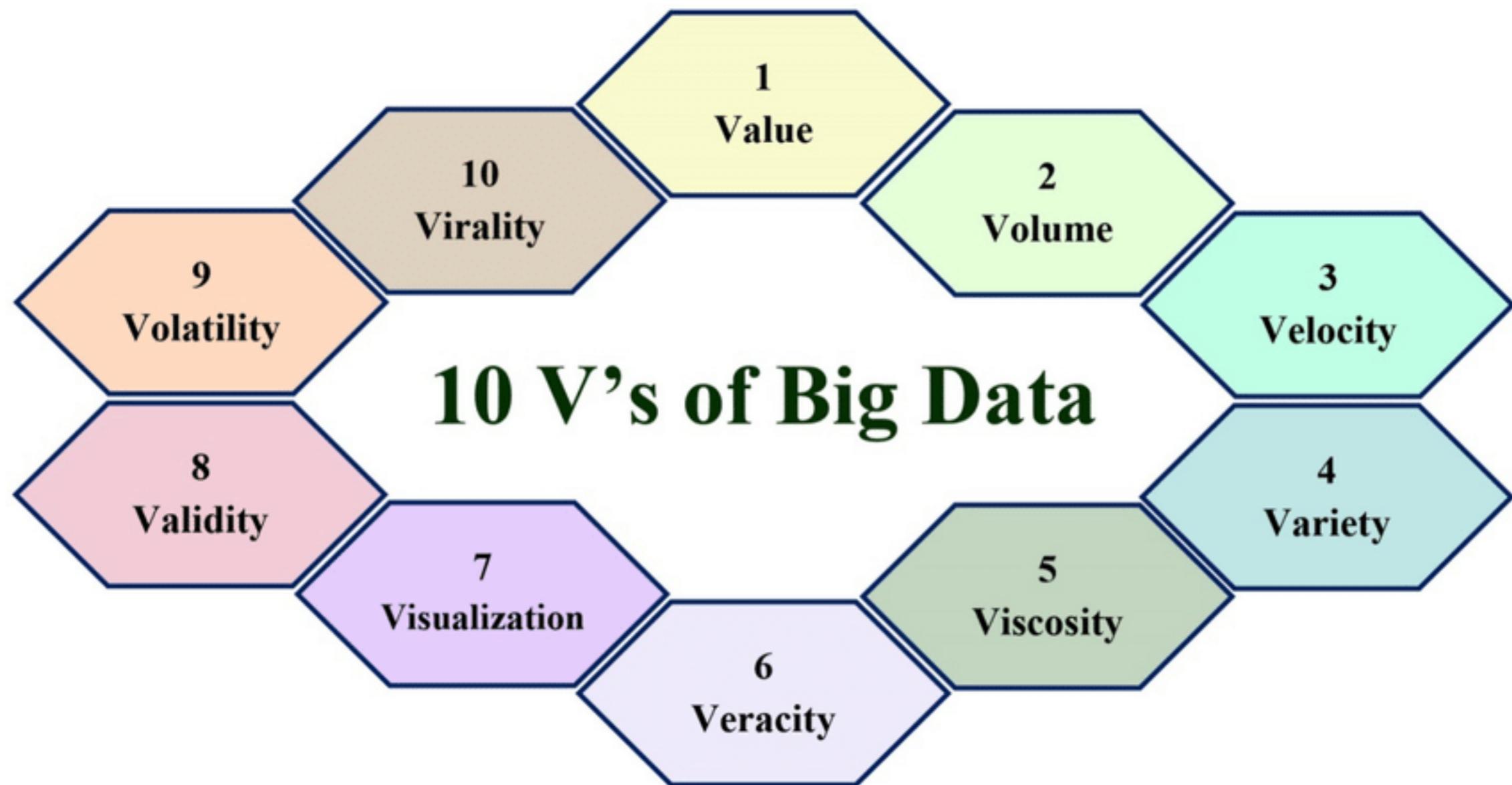
# Characteristics for Big Data



# Characteristics for Big Data



# Characteristics for Big Data

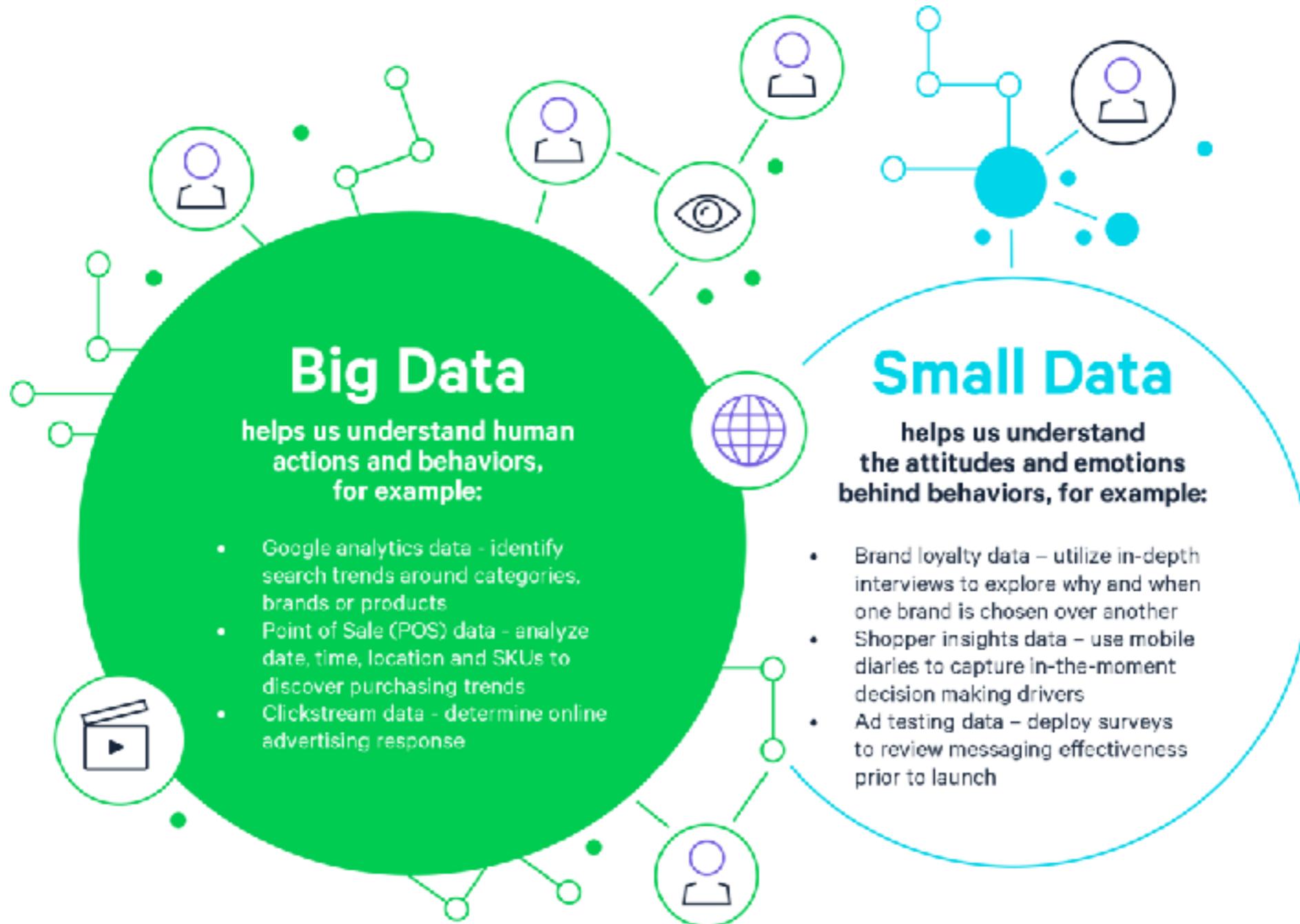


# Purpose of Big Data

To extract **meaningful** information from massive datasets that conventional tools can't manage, using **advanced analytics** and data processing methods.



# Big Data vs Small Data

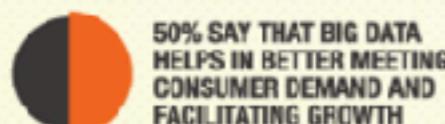
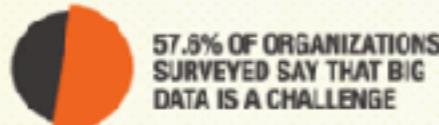


# BIG DATA



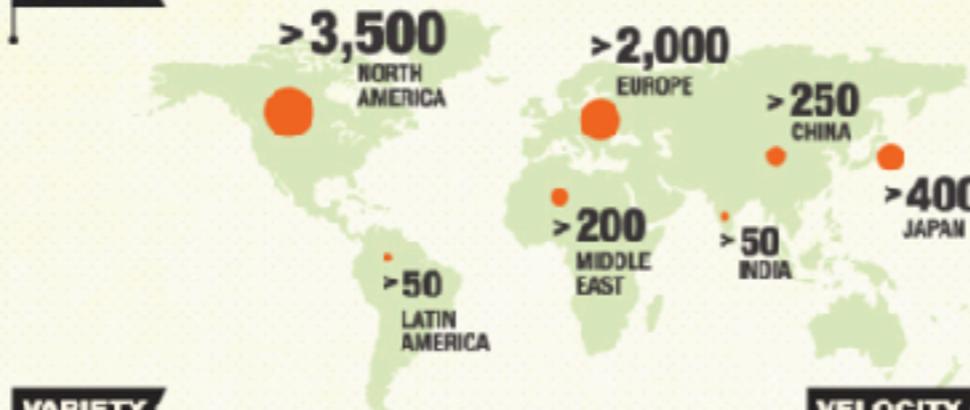
*Big Data* is data that is too large, complex and dynamic for any conventional data tools to capture, store, manage and analyze.

The right use of Big Data allows analysts to spot trends and gives niche insights that help create value and innovation much faster than conventional methods.



The “three V’s”, i.e the Volume, Variety and Velocity of the data coming in is what creates the challenge.

## VOLUME



## VARIETY



### PEOPLE TO PEOPLE

NETIZENS, VIRTUAL COMMUNITIES, SOCIAL NETWORKS, WEB LOGS...



### PEOPLE TO MACHINE

ARCHIVES, MEDICAL DEVICES, DIGITAL TV, E-COMMERCE, SMART CARDS, BANK CARDS, COMPUTERS, MOBILES...



### MACHINE TO MACHINE

SENSORS, GPS DEVICES, BAR CODE SCANNERS, SURVEILLANCE CAMERAS, SCIENTIFIC RESEARCH...

## VELOCITY



### 2.9 MILLION EMAILS SENT EVERY SECOND



### YOU TUBE

20 HOURS OF VIDEO UPLOADED EVERY MIN



50 MILLION TWEETS PER DAY

\$300 billion is the potential annual value to Healthcare

TRANSPARENCY IN CLINICAL DATA AND CLINICAL DECISION SUPPORT

**\$165B CLINICAL**

**\$9B PUBLIC HEALTH**

PUBLIC HEALTH SURVEILLANCE AND RESPONSE SYSTEMS

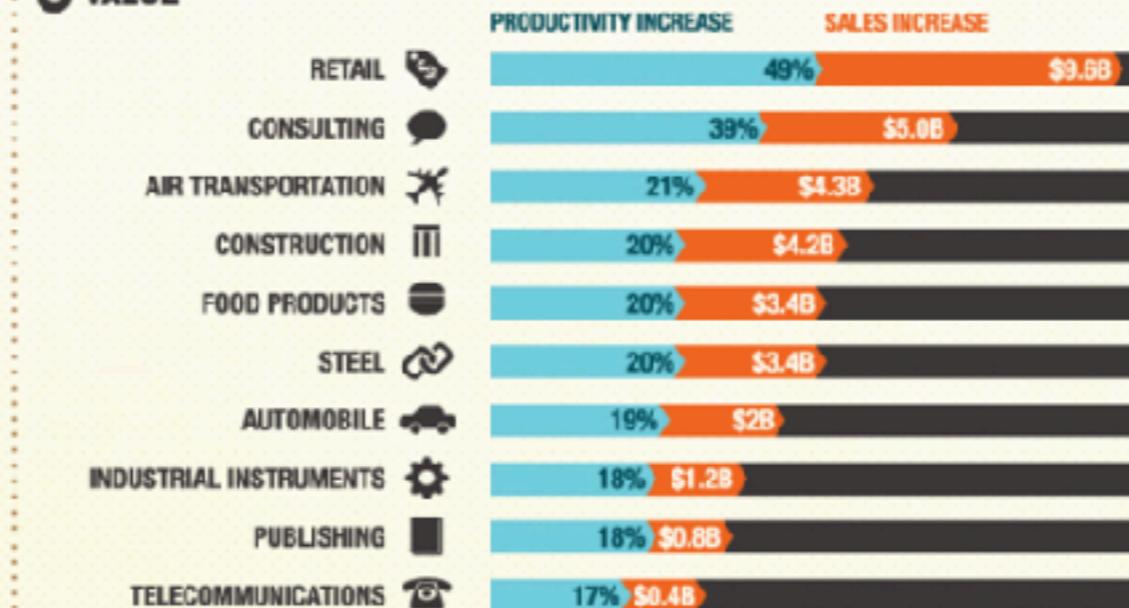
**\$108B R&D**

**\$5B BUSINESS MODEL**

**\$47B ACCOUNTS**

RESEARCH AND DEVELOPMENT; PERSONALIZED MEDICINE; CLINICAL TRIAL DESIGN  
ADVANCED FRAUD DETECTION; PERFORMANCE BASED DRUG PRICING

## VALUE



40% PROJECTED GROWTH IN GLOBAL DATA CREATED PER YEAR



5% PROJECTED GROWTH IN GLOBAL IT SPENDING PER YEAR



The estimated size of the digital universe in 2011 was 1.8 zettabytes. It is predicted that between 2009 and 2020, this will grow 44 fold to 35 zettabytes per year. What is the right data management strategy for you, to successfully utilize Big Data?

Sources - ① Realizing the Rewards of Big Data - Wipro Report ② Big Data: The Next Frontier for Innovation, Competition and Productivity - McKinsey Global Institute Report ③ ComScore, Radicati Group ④ Measuring the Business Impacts of Effective Data - study by University of Texas/Austin ⑤ U.S. Department of Labor.

DO BUSINESS BETTER

WIPRO | OVER 100,000 EMPLOYEES | 54 COUNTRIES | CONSULTING | SYSTEM INTEGRATION | OUTSOURCING

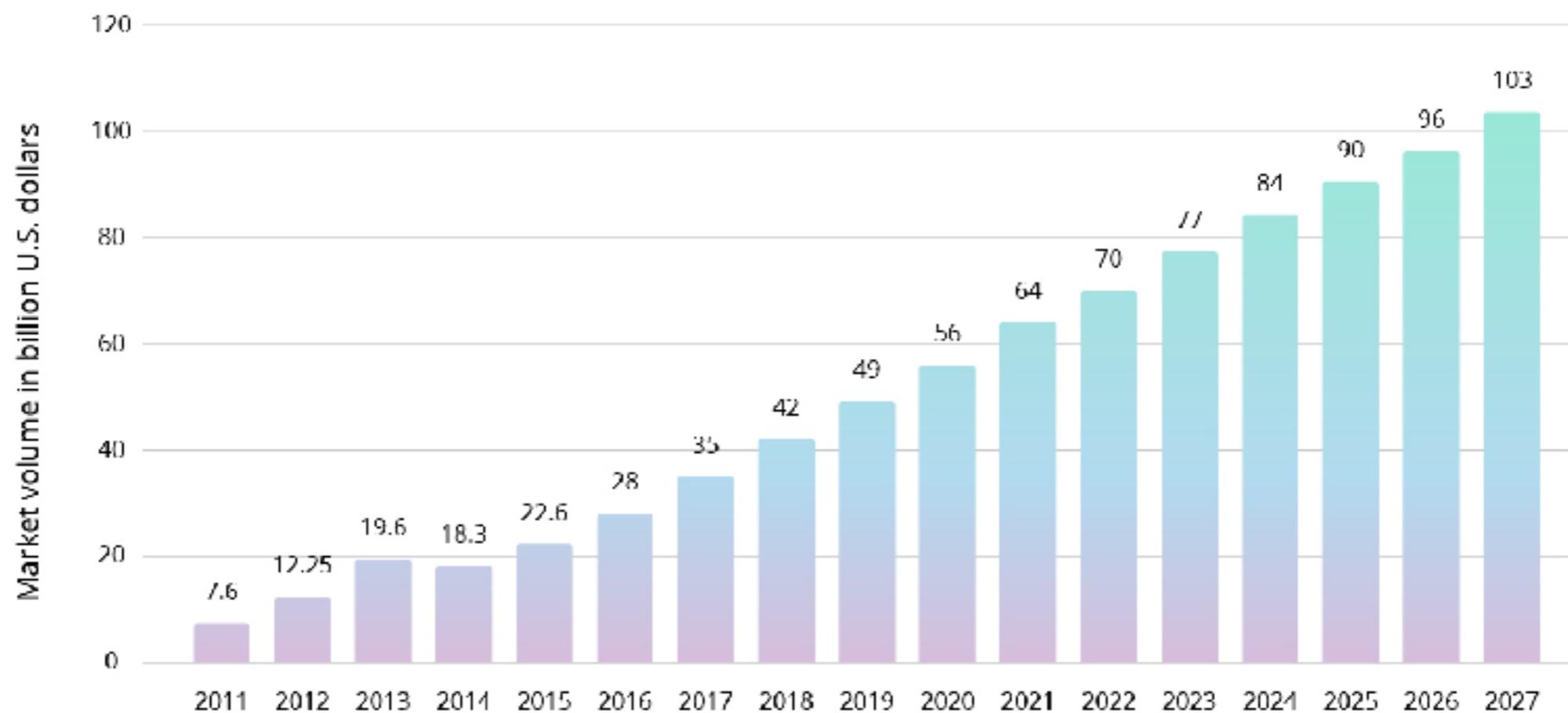


Sharing

16

# Big Data Trend

**Big data market size revenue forecast worldwide from 2011 to 2027**  
(in billion U.S. dollars)



<https://www.statista.com/statistics/254266/global-big-data-market-forecast/>



# Big Data Challenges

Storage

Data Quality

Analysis

Cost vs Value

Accessibility

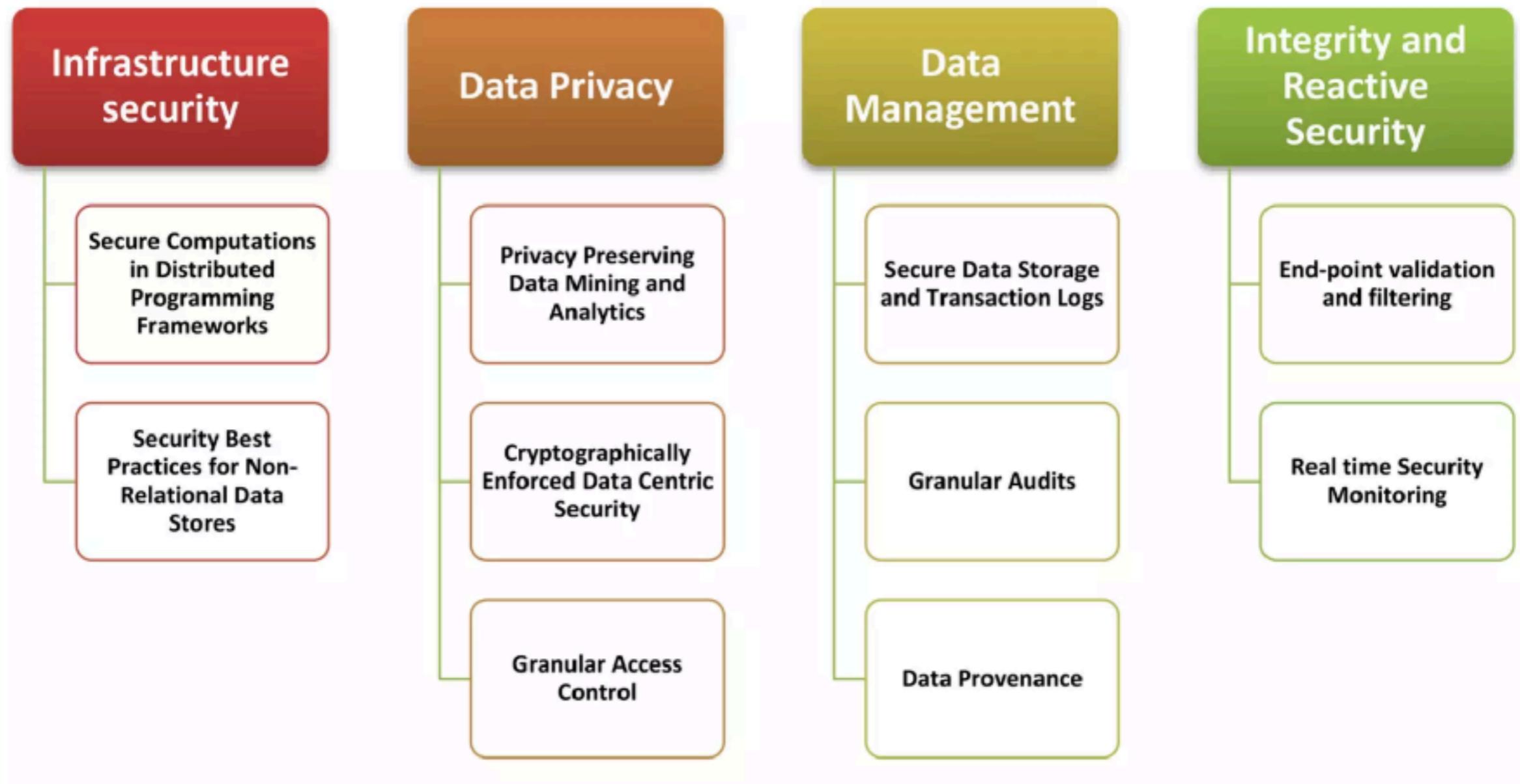
Security



# Big Data Challenges



# Big Data Security Challenges



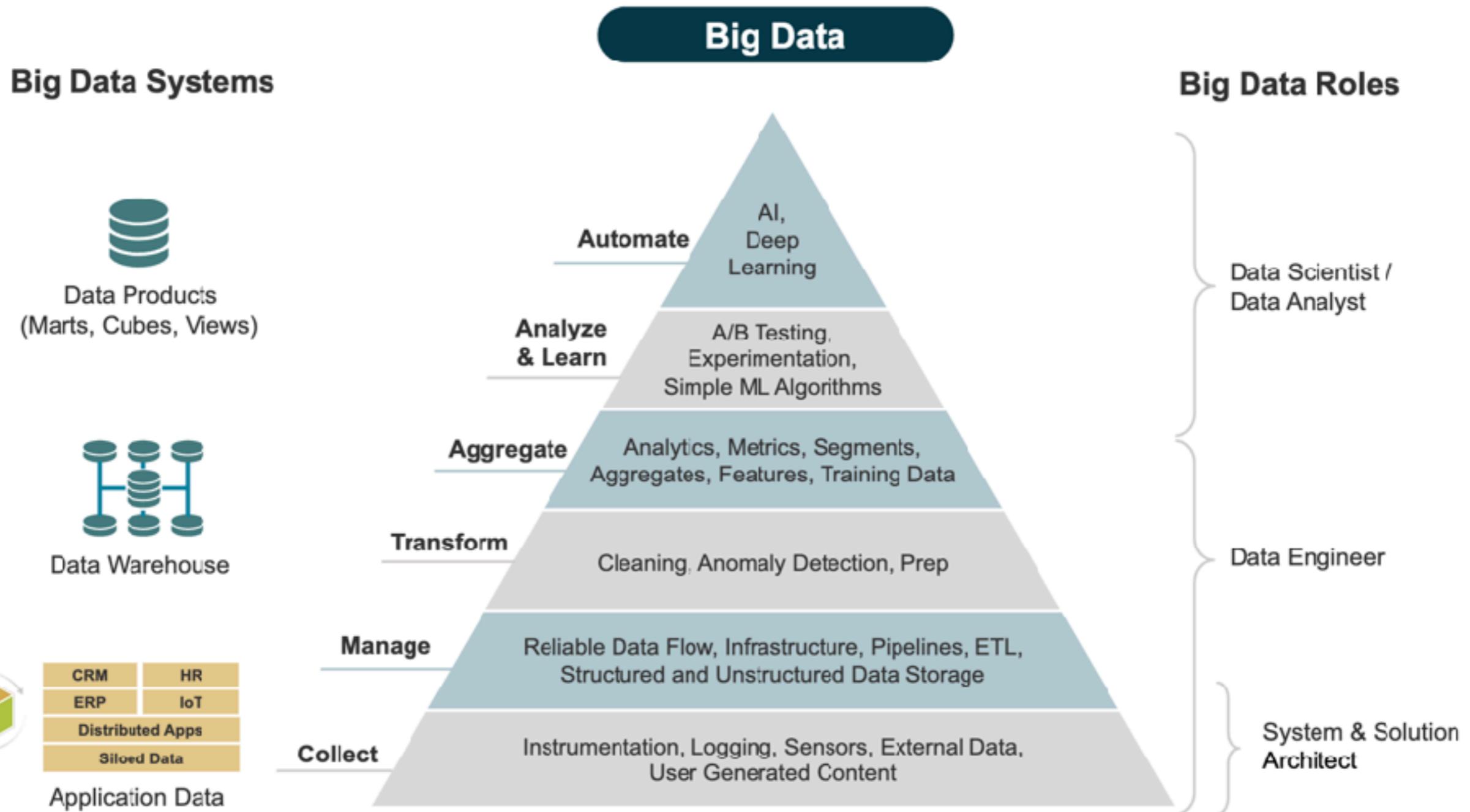
# Data Governance



<https://intellias.com/future-big-data-trends/>



# All about Big Data



© Scaled Agile, Inc.

<https://scaledagileframework.com/big-data/>

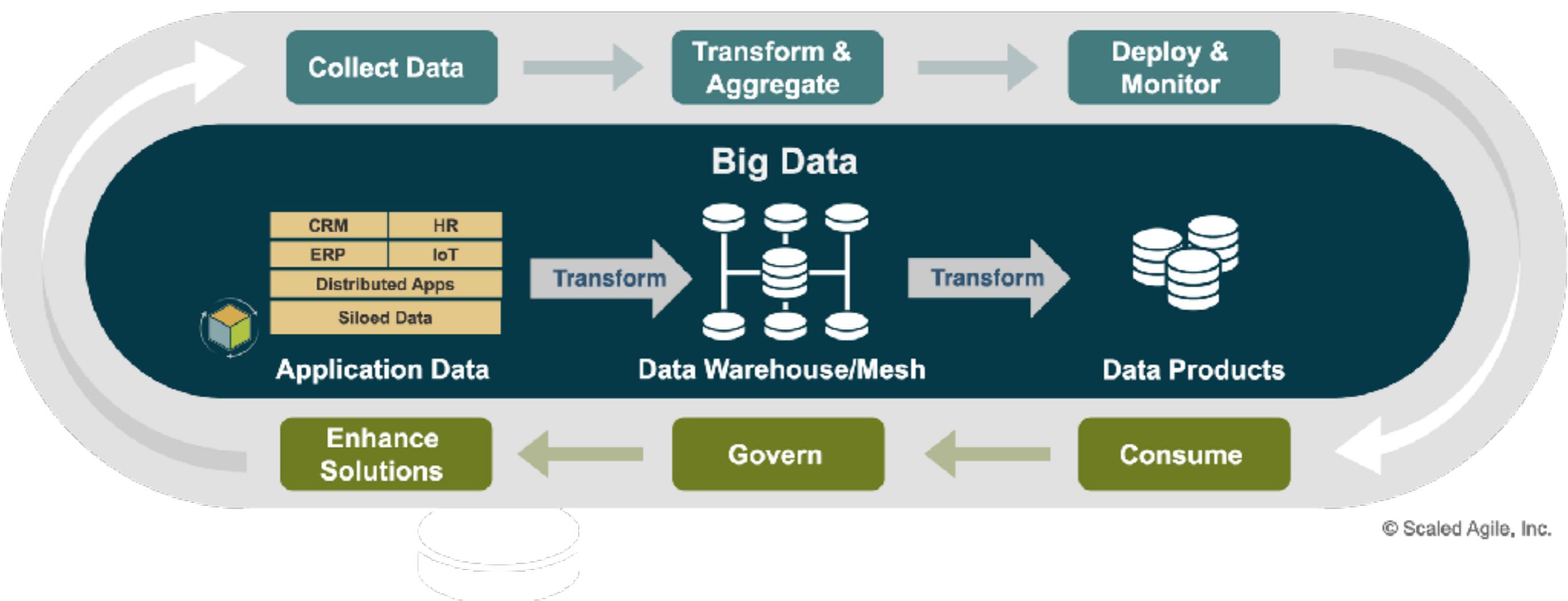


Sharing

© 2020 - 2024 Siam Chamnankit Company Limited. All rights reserved.

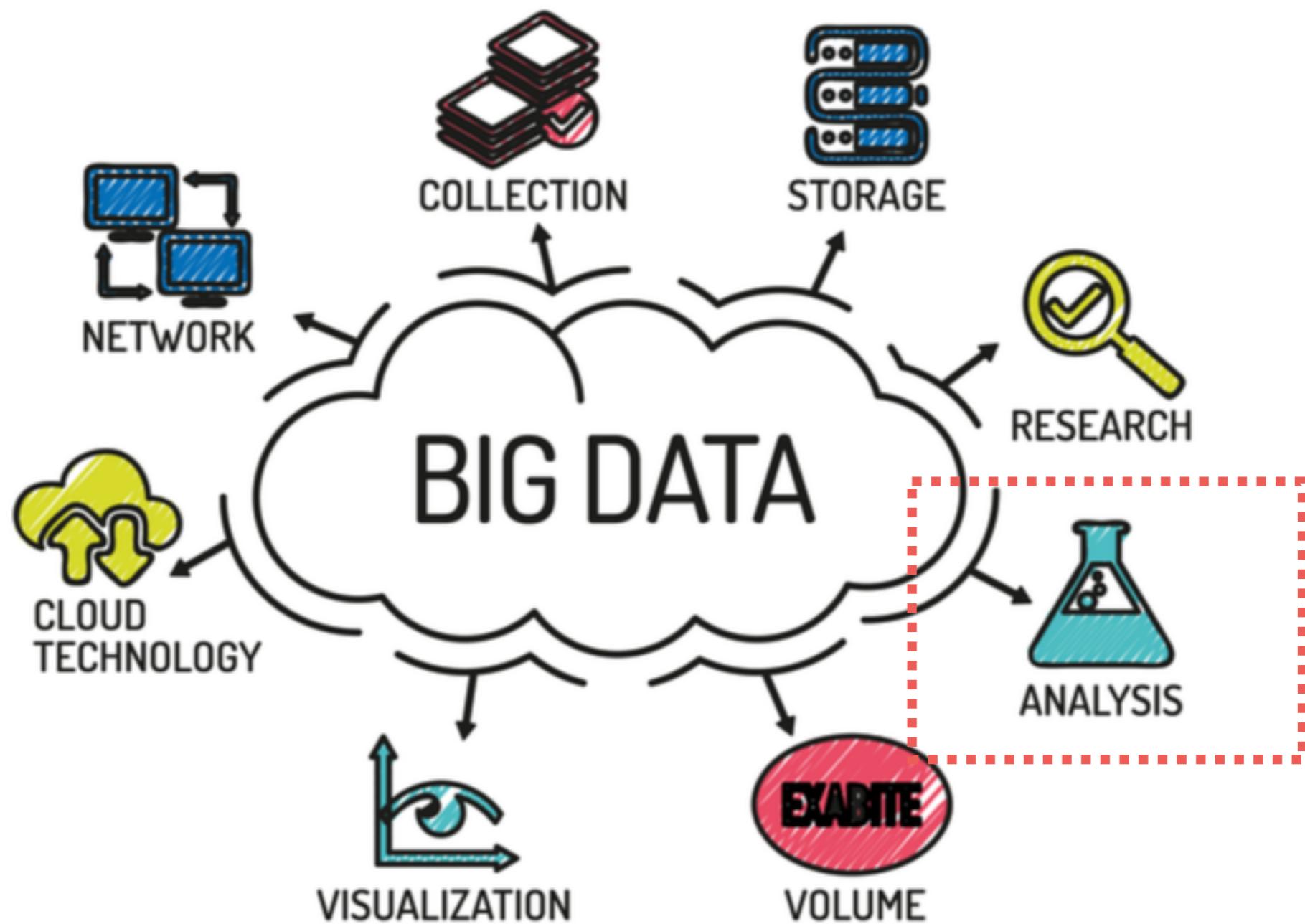
# DataOps

## Collaboration data management across teams

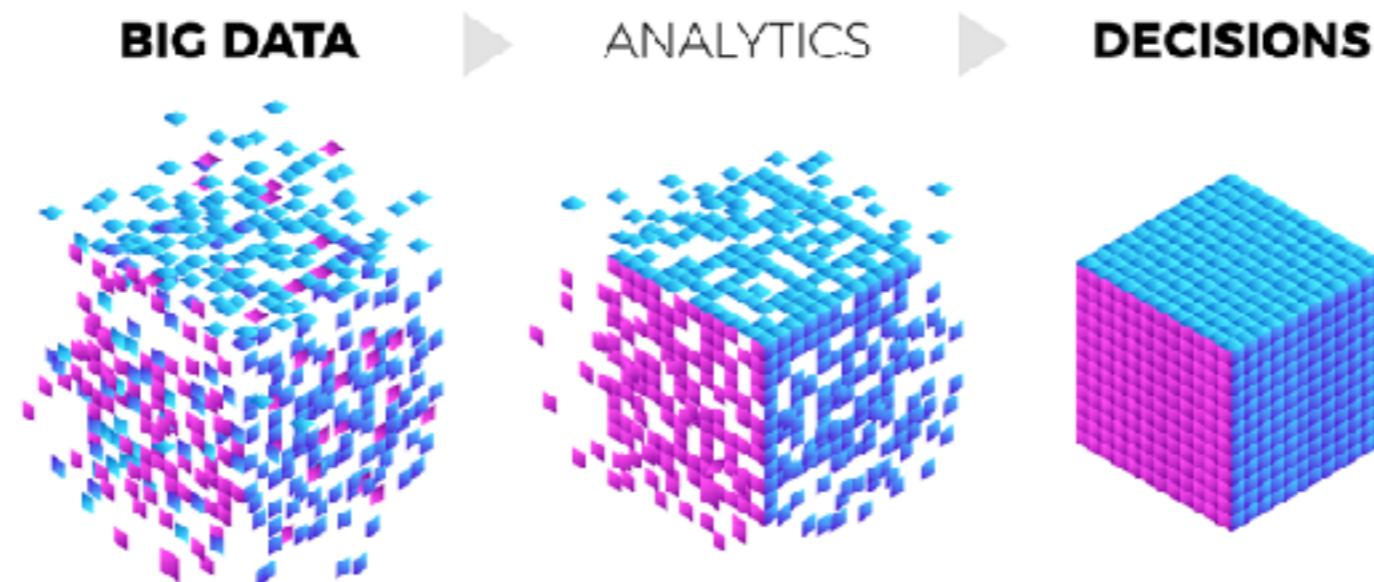


<https://scaledagileframework.com/big-data/>





# Rise of Big Data Analytics



stargazr



# **Introduction to Business Analytics (BA)**



# **Business Analytics (BA)**

**Practice** of using data, statistical analysis, and quantitative methods to drive **decision-making** and **improve business outcomes**

It often involves predictive modeling, data mining, and machine learning to **forecast trends** and derive **insights**.



# Purpose of BA

The goal of BA is  
to help businesses make informed,  
**data-driven decisions** by predicting  
future outcomes or trends



# Data-driven Decision



Make confident decisions



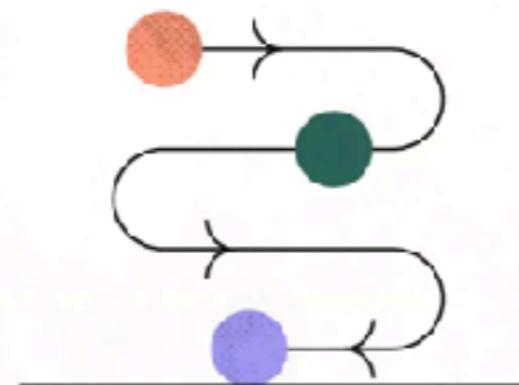
Guard against biases



Find unresolved questions



Set measurable goals



Improve company processes

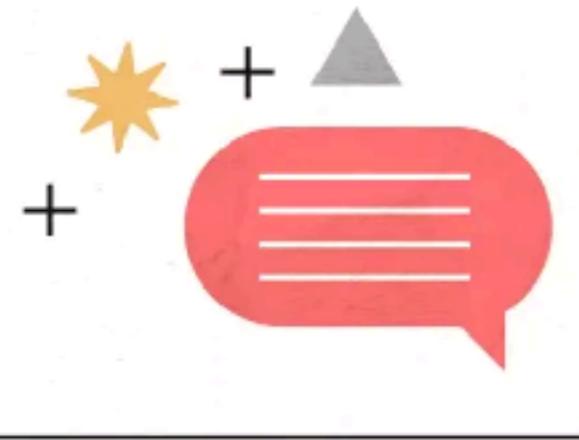
<https://asana.com/resources/data-driven-decision-making>



Sharing

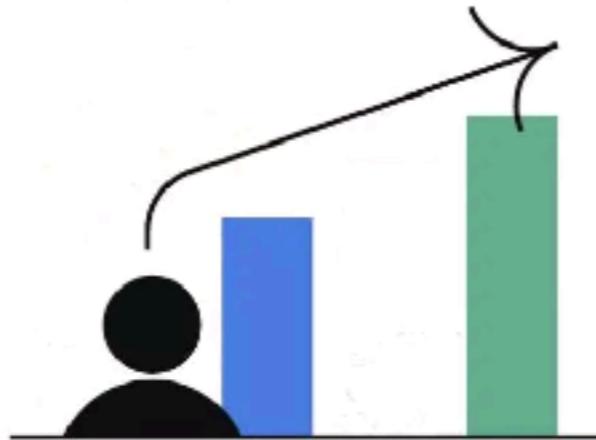
© 2020 - 2024 Siam Chamnankit Company Limited. All rights reserved.

# Tips for Data-driven Decision



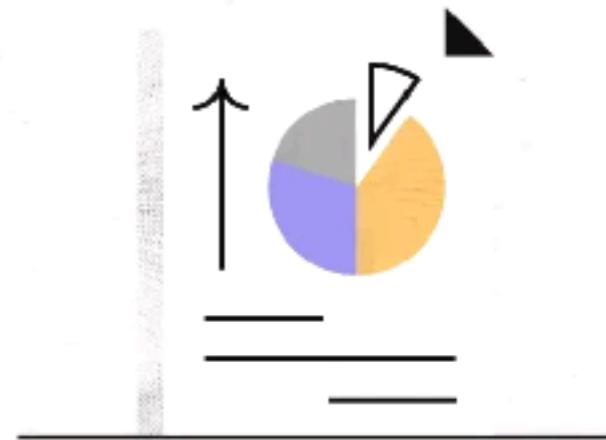
## Find the story

When analyzing data, look for patterns and meanings behind the numbers and figures.



## Consult the data

Before making any gut decisions, see if the facts align with your feelings.



## Learn data visualization

Make sure your data is visually appealing and easy to understand.

<https://asana.com/resources/data-driven-decision-making>



# Key Techniques

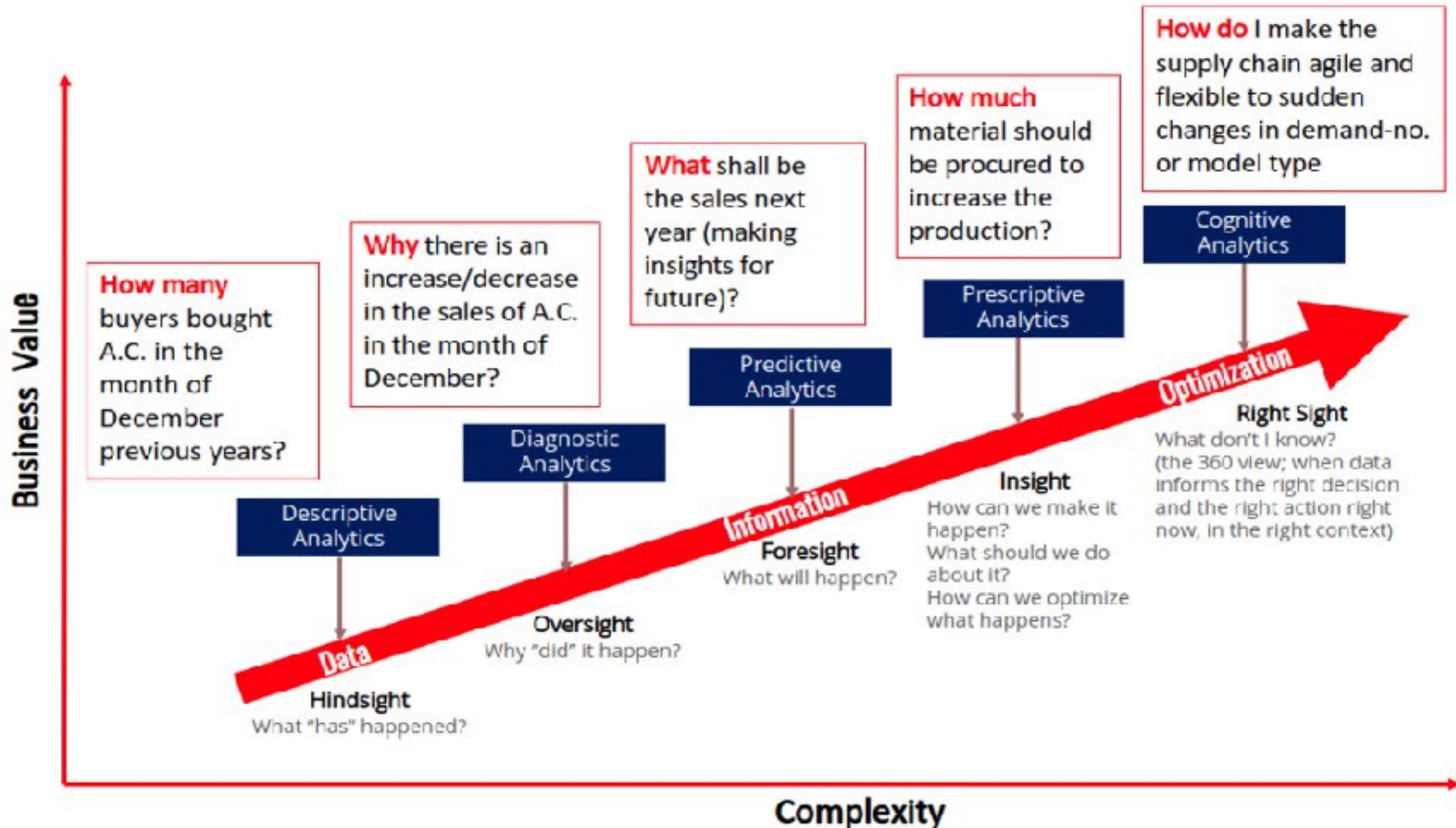
Predictive  
Analytics

Prescriptive  
Analytics

Diagnostic  
Analytics



# Key Techniques



# **Descriptive Analytics**

Descriptive analytics answers the question,

**“What happened?”**

by summarizing historical data to identify patterns  
and trends.



# Descriptive Analytics

It helps businesses understand **past performance** and gain insights from historical data to determine what has occurred over a given period.

Basic data aggregation

Reporting tool

Data visualisation tool

Statistic and trend



# Diagnostic Analytics

Diagnostic analytics answers the question,

**“Why did it happen?”**

by analyzing data to understand  
the root causes of past events.



# Diagnostic Analytics

It digs deeper into descriptive analytics to find the reasons behind specific outcomes or patterns.

Drill-down  
analytics

Correlation  
analysis

Root cause  
analysis

Statistic with  
regression analysis



# Predictive Analytics

Predictive analytics answers the question,

**“What will likely happen in the future?”**

by using historical data, statistical models, and machine learning techniques to forecast future trends.



# Predictive Analytics

It helps businesses anticipate **future outcomes** based on patterns and trends found in historical data.

Statistic  
modeling

Machine learning  
algorithm

Time series  
forecast



# Predictive Analytics



# Prescriptive Analytics

Prescriptive analytics answers the question,

**“What should be done?”**

by recommending specific actions  
or strategies to achieve desired outcomes or  
optimize processes.



# Prescriptive Analytics

It goes beyond predicting future outcomes by suggesting the best course of action based on predictions.

Optimization  
model

Simulation

Decision/scenario  
analysis



# Summary

Type	Main Question	Purpose	Techniques	Example
<b>Descriptive Analytics</b>	What happened?	Understanding past performance.	Data aggregation, reporting, visualization	Sales reports showing performance by region.
<b>Diagnostic Analytics</b>	Why did it happen?	Identifying causes of past outcomes.	Drill-down analysis, correlation, root cause	Analyzing why sales declined in a certain area.
<b>Predictive Analytics</b>	What will happen?	Forecasting future trends and outcomes.	Machine learning, statistical models	Predicting customer churn for an online store.
<b>Prescriptive Analytics</b>	What should we do?	Recommending actions to optimize future outcomes.	Optimization, simulation, decision models	Suggesting optimal pricing strategies for sales.



# Summary

## **Descriptive Analytics**

focuses on summarizing what has already happened

## **Diagnostic Analytics**

helps businesses understand why those things happened

## **Predictive Analytics**

forecasts future trends based on historical data

## **Prescriptive Analytics**

provides recommendations to optimize decision-making and future outcomes



# **Introduction to Business Intelligence (BI)**



# Business Intelligence (BI)

Refers to technologies, applications, and practices used for the collection, integration, analysis, and presentation of business data.

Collect

Integrate

Analysis

Visualize



# Purpose of BI

**BI is about reporting and dashboards  
that help businesses monitor their current state  
and past performance  
by organizing data into actionable insights.**



# Start with Data Sources

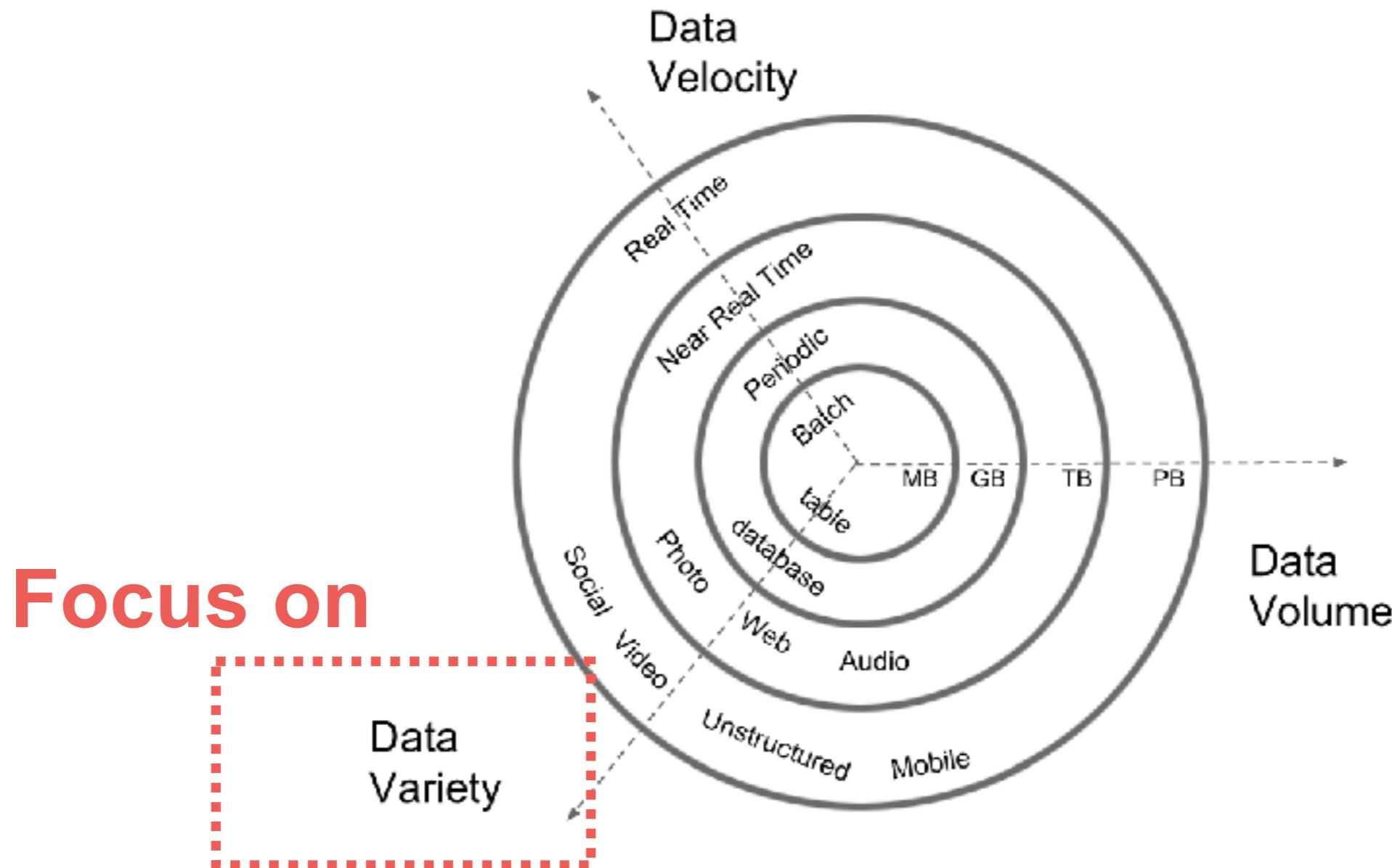


# Data Sources !!



# Characteristics for Big Data

## 3 V's !!!

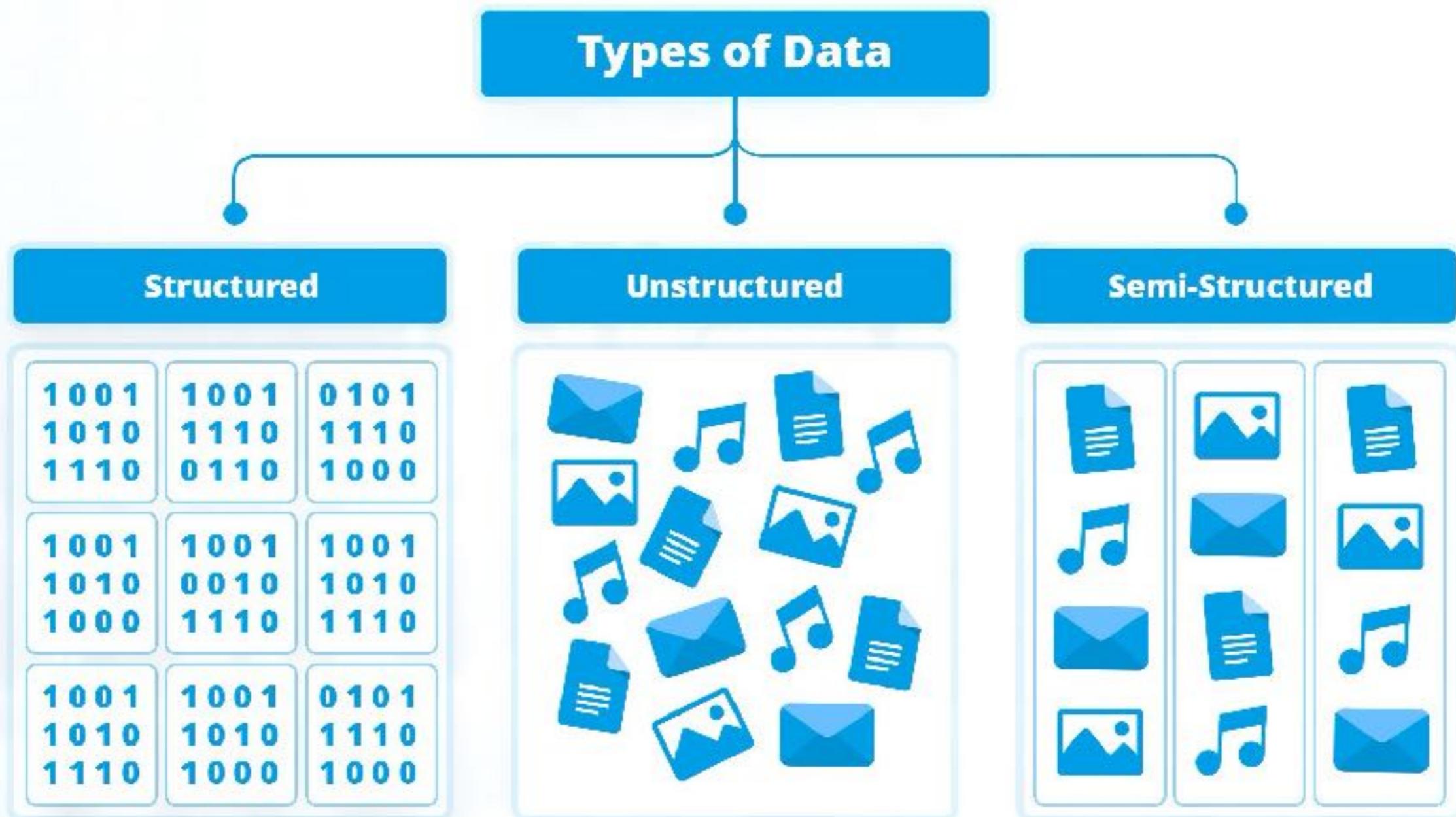


# Types of data sources

Structured  
Unstructured  
Semi-structured



# Types of data



Astera  
Enabling Data-Driven Innovation

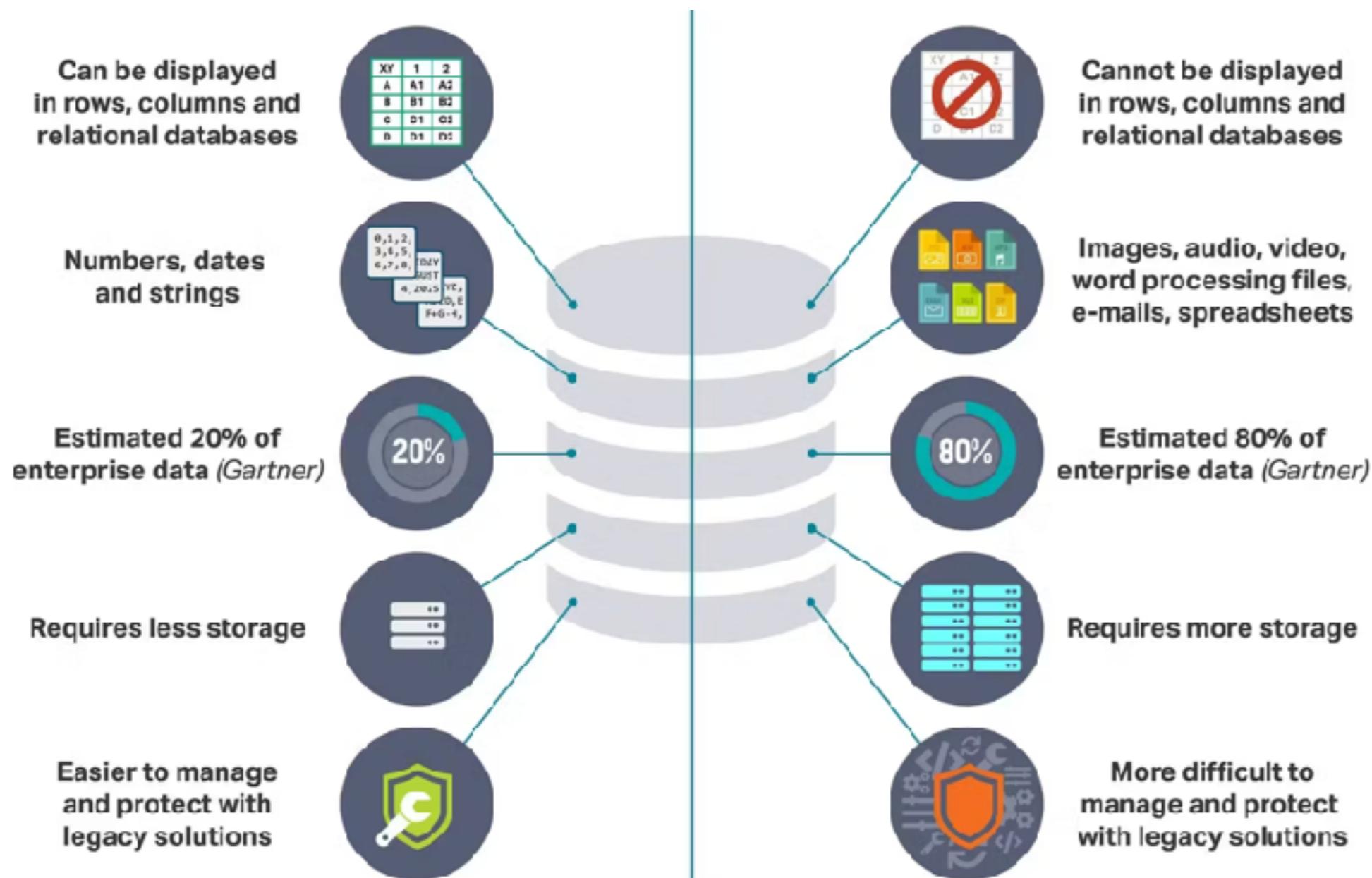


# Types of data sources

Feature	Structured Data	Unstructured Data	Semi-structured Data
Definition	Data organized in a predefined format with a clear structure	Data lacking a predefined format and organization	Data with some internal structure but lacking a rigid schema
Examples	Relational databases, spreadsheets, CSVs, APIs	Text documents, emails, images, videos, social media posts	JSON files, XML files, log files
Characteristics	Standardized format, easily searchable and analyzable, consistent and reliable	Diverse formats, challenging to process and analyze, rich and diverse information	Flexible format, adaptable to evolving data needs, requires specialized tools
Advantages	Easy to process and analyze, supports efficient data retrieval, suitable for statistical analysis	Rich and diverse information, captures real-world context, valuable for sentiment analysis and trend identification	Adaptable to evolving data needs, flexible and scalable, suitable for real-time applications
Disadvantages	Limited flexibility, unable to capture complex relationships, not suitable for all types of data	Difficult to process and analyze, requires specialized tools, data quality concerns	Limited data integration potential, lack of standardized formats, evolving data structures
Use Cases	Business intelligence, financial transactions, scientific research, data warehousing	Customer feedback analysis, social media monitoring, content analysis, multimedia processing	Real-time analytics, sensor data analysis, web scraping, scientific experiments
Tools and Techniques	SQL databases, spreadsheets, data warehouses, data analytics tools	Natural language processing (NLP), machine learning, sentiment analysis, image recognition	JSON parsers, XML parsers, stream processing tools, data pipelines



# Structured vs Unstructured data

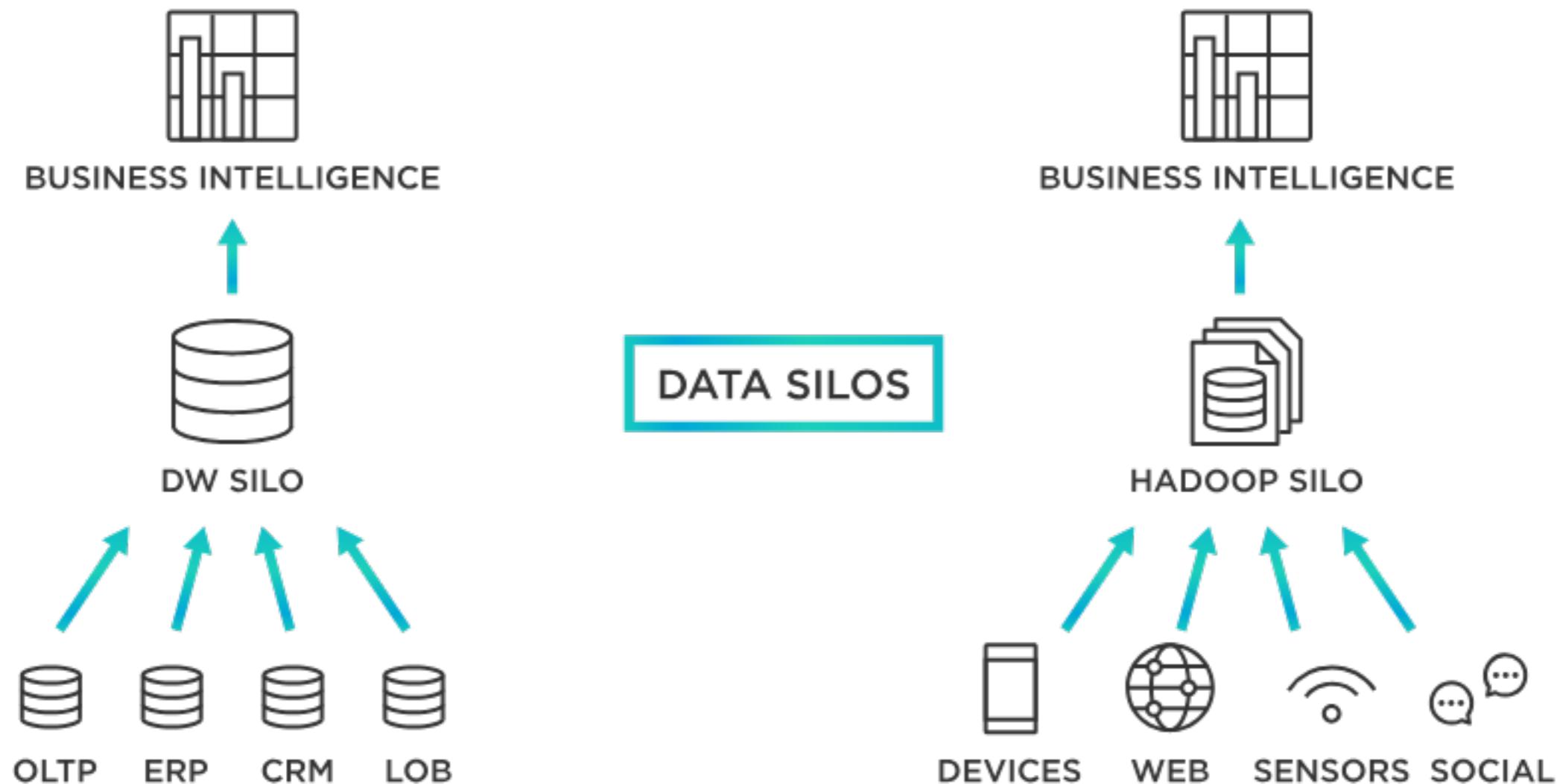


# Unstructured Data Challenges

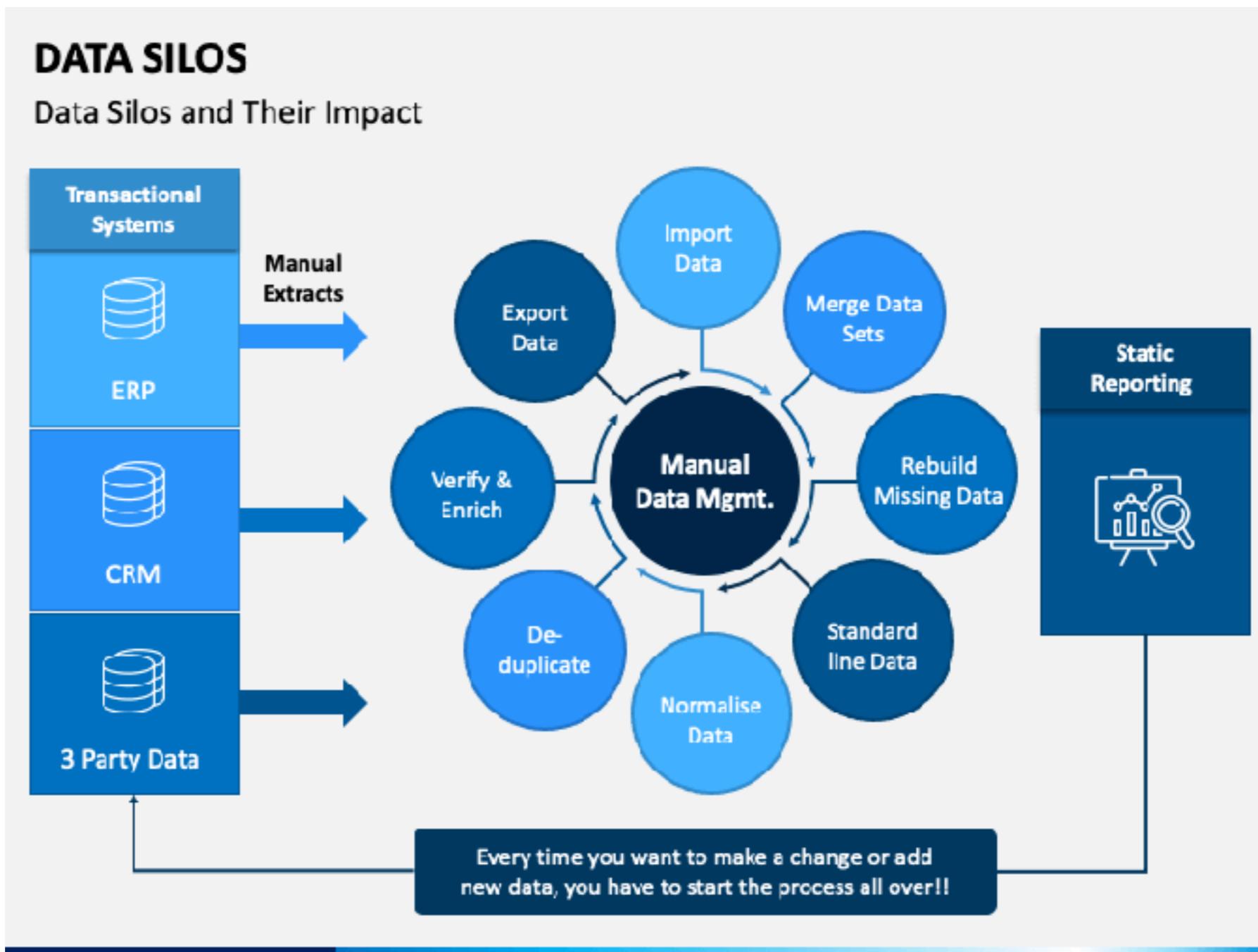
**Inability** to process growing data volumes  
Accessing **siloed** data  
Regulatory non-compliance  
**Reduced data usability**  
Increased vulnerability to cyber attacks



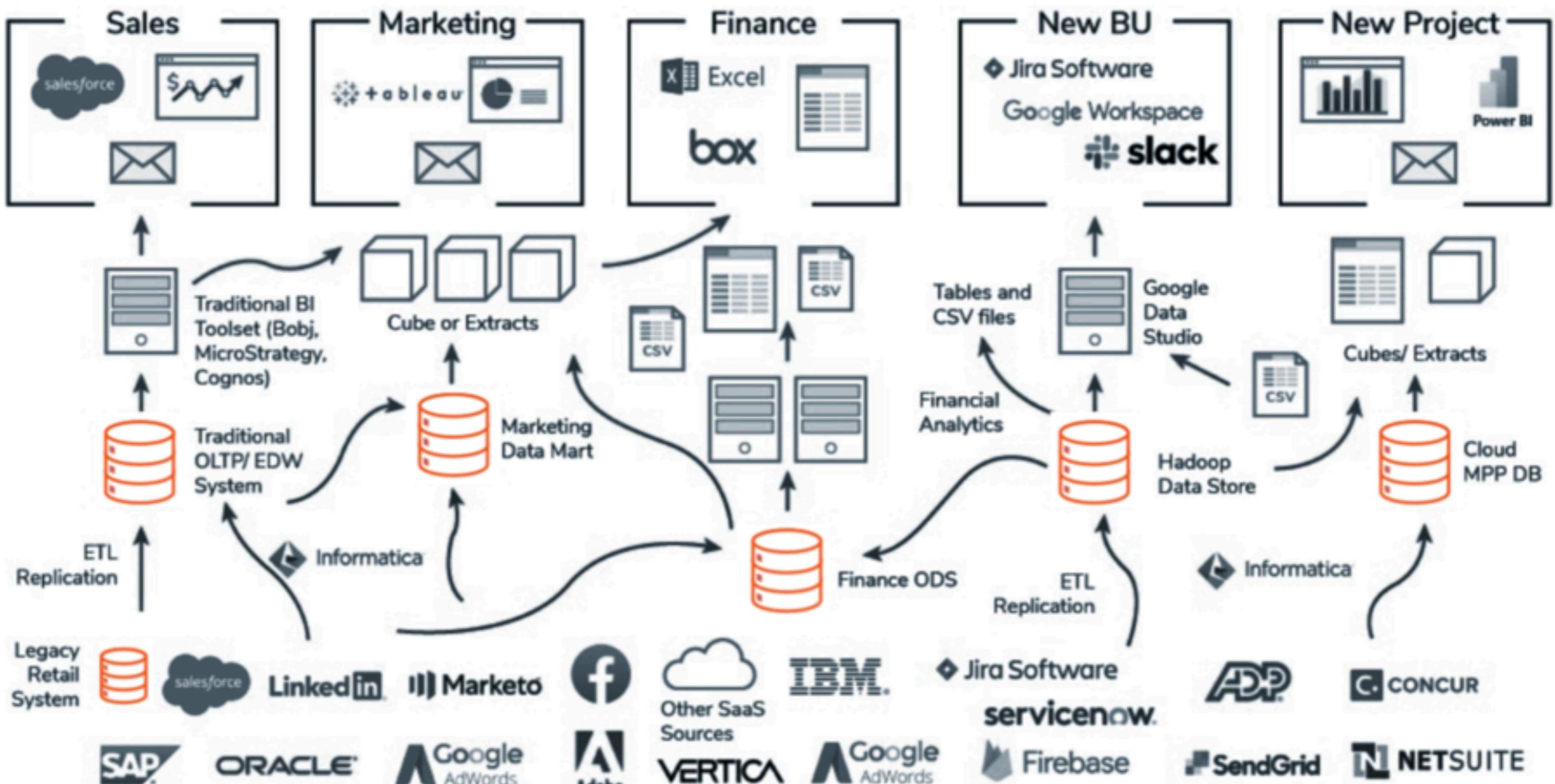
# Data siloed ?



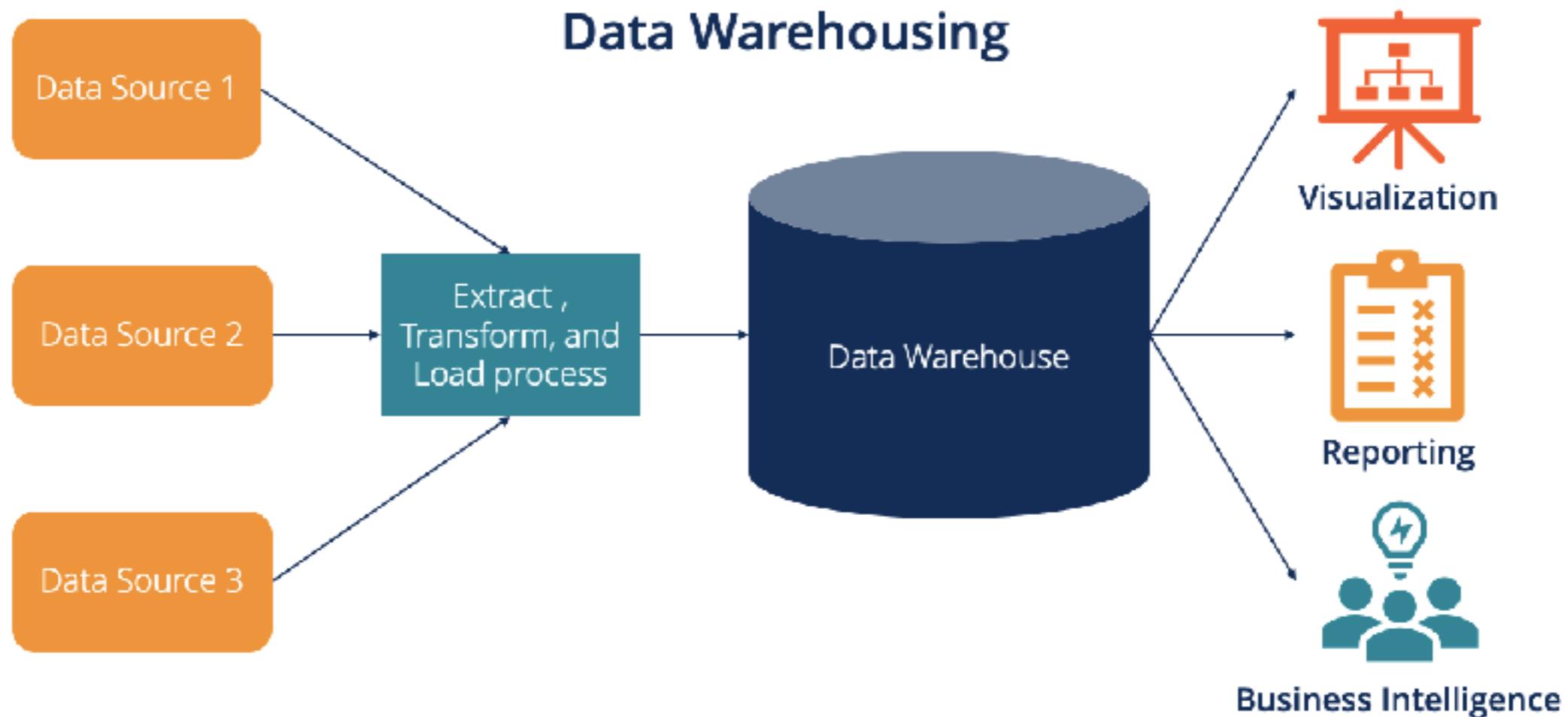
# Data siloed ?



# Data siloed ?

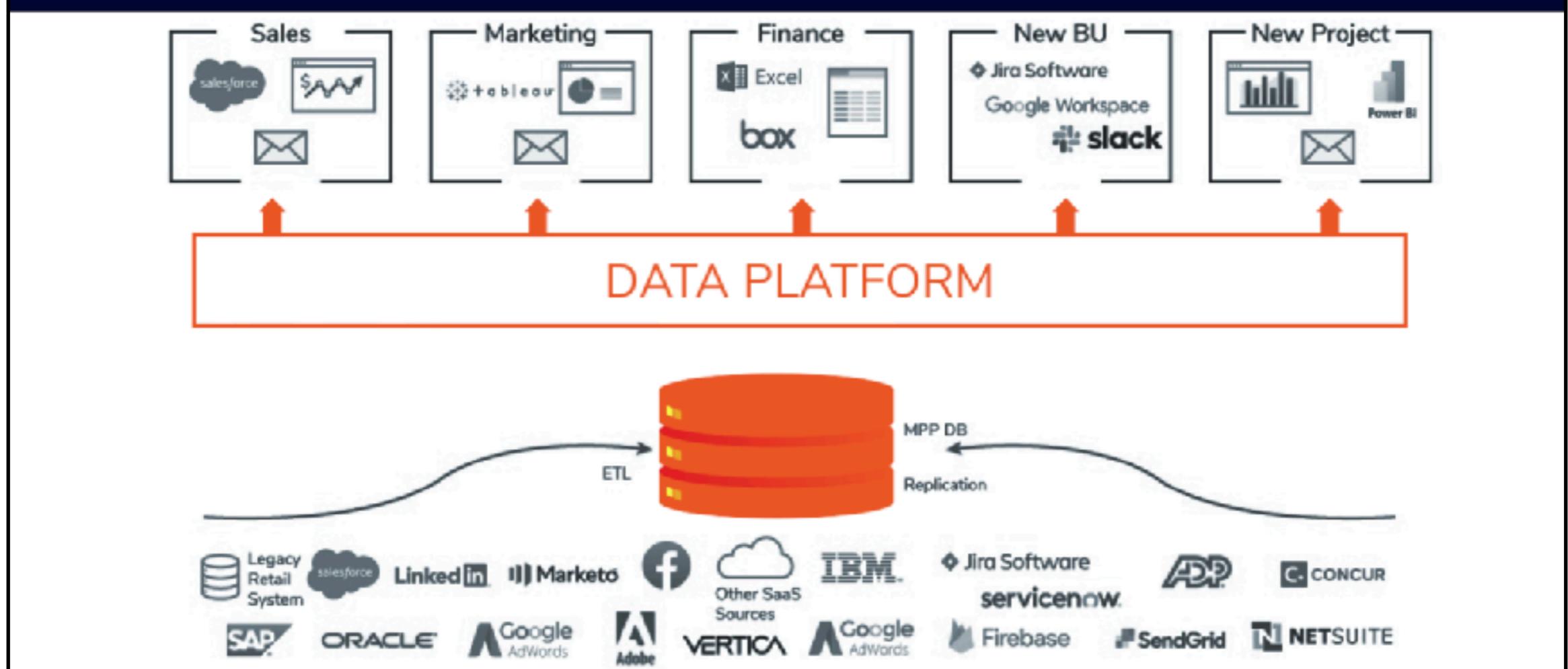


# Try to solve with data platform



# Data Platform

We clean it up



# **Basic of Data Management**



# Data management

Data management refers to the **process** of organizing, storing, maintaining, and retrieving data efficiently and securely.

It plays a crucial role in **ensuring** that data is accurate, available, and usable for decision-making, analysis, and operations



# Keys of data management

Data collection

Data storage

Data Security

Data quality

Data governance

Data Integration

Data backup/  
restore

Data  
accessibility

Data analysis



# Data storage ?

Database

Data  
Warehouse

Data Mart

Data Lake



# Data Warehouse

A data warehouse is a centralized repository that stores **structured** and processed data from **multiple sources**.

It is optimized for querying and analysis, primarily used for reporting, business intelligence, and decision-making.



# Data Mart

A data mart is a smaller, **more focused subset of a data warehouse**, designed to serve the needs of a **specific business unit or department** (e.g., sales, marketing).

It contains only relevant data for that particular group, **enabling quicker access and more tailored analysis**.



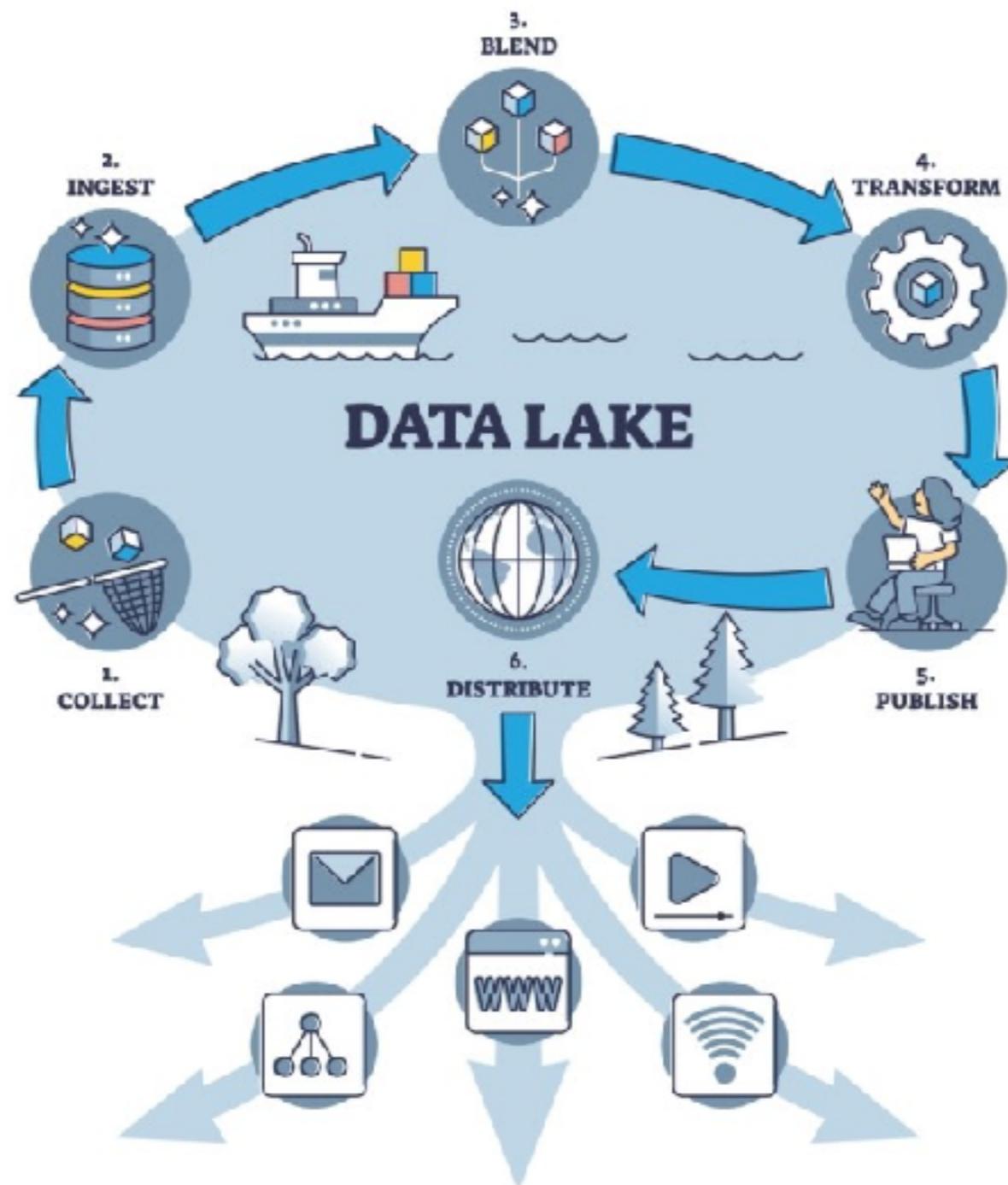
# Data Lake

A data lake is a **vast storage** repository that holds a large amount of raw data in its **native format** (structured, semi-structured, and unstructured) until it is needed for analysis.

It is designed for **storing data** before it is processed and transformed for analysis, offering flexibility and scalability.



# Data Lake

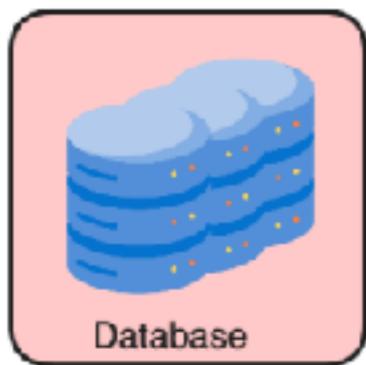


# Key differences

Feature	Data Warehouse	Data Mart	Data Lake
<b>Data Type</b>	Structured	Structured	Structured, semi-structured, and unstructured
<b>Scope</b>	Enterprise-wide (broad)	Department-specific (narrow)	Enterprise-wide (but raw)
<b>Storage Type</b>	Structured (rows/columns)	Structured	Raw format
<b>Data Processing</b>	Pre-processed (ETL)	Pre-processed (ETL)	Raw, processed when queried
<b>Primary Use</b>	Reporting, historical analysis	Department-specific insights	Big data analysis, machine learning
<b>Speed of Query</b>	Fast for complex queries	Faster for localized queries	Slower for complex queries
<b>Cost</b>	Generally more expensive	Less expensive than data warehouse	Cheaper for storage, but processing may incur costs
<b>Schema</b>	Schema-on-write	Schema-on-write	Schema-on-read



# Different between data types



VS



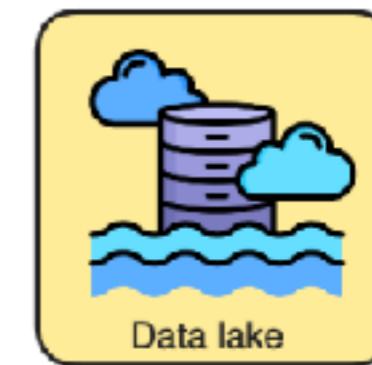
Data Warehouse

VS



Data mart

VS



Data lake

## Scope

Application-specific

Organization-wide,  
structured data.

Department-specific,  
structured data.

Organization-wide,  
any type of data

## Data Type

Structured

Structured

Structured

Structured,  
semi-structured,  
unstructured.

## Structure

Predefined schema

Schema on write

Schema on write (inherited  
from data warehouse)

Schema on read

## Use Case

Operational  
applications(OLTP)

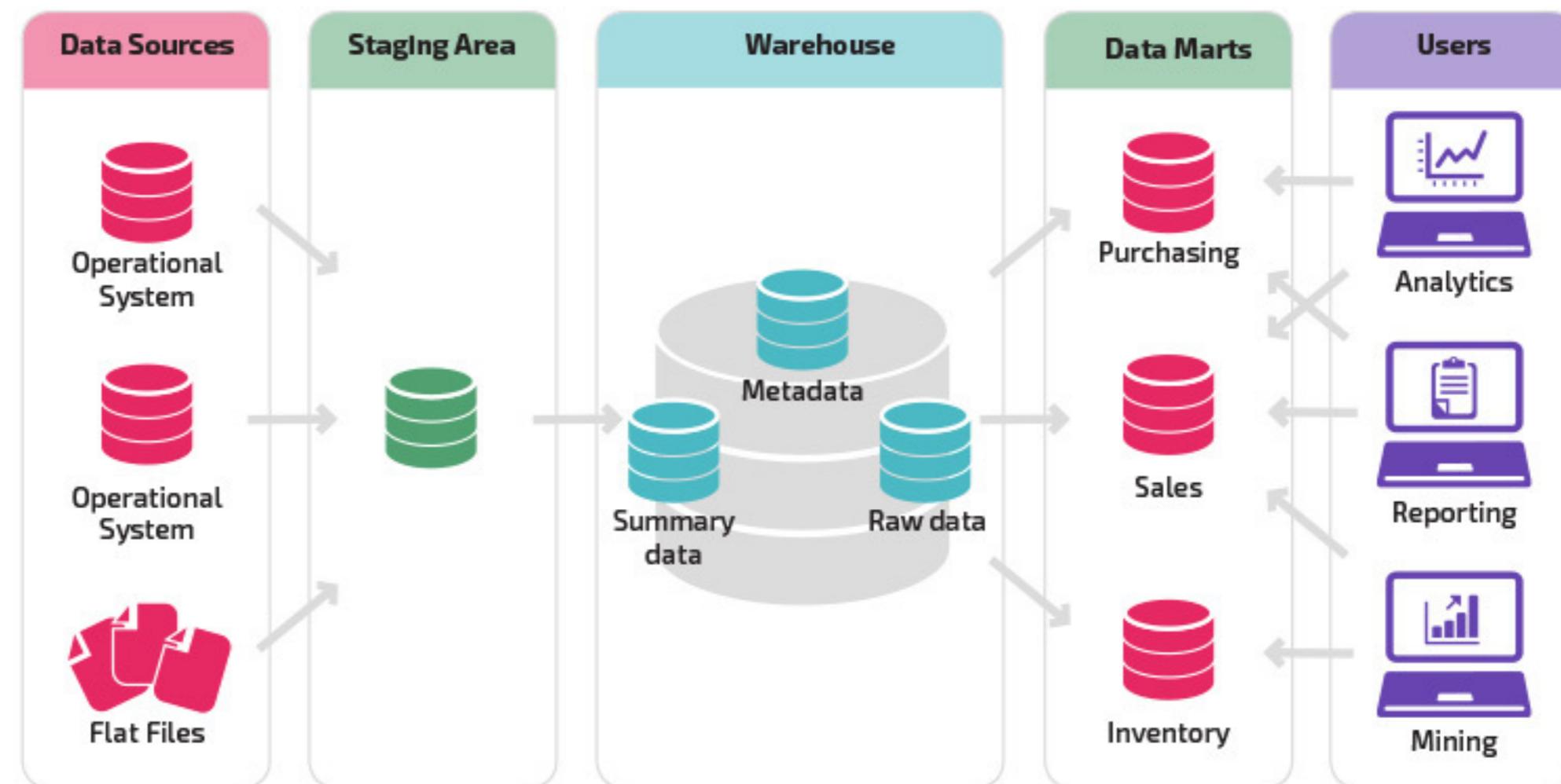
Business intelligence,  
historical  
analysis(OLAP).

Specific business function  
analysis

Big data analytics,  
data exploration.



# Data Warehouse vs Data Mart



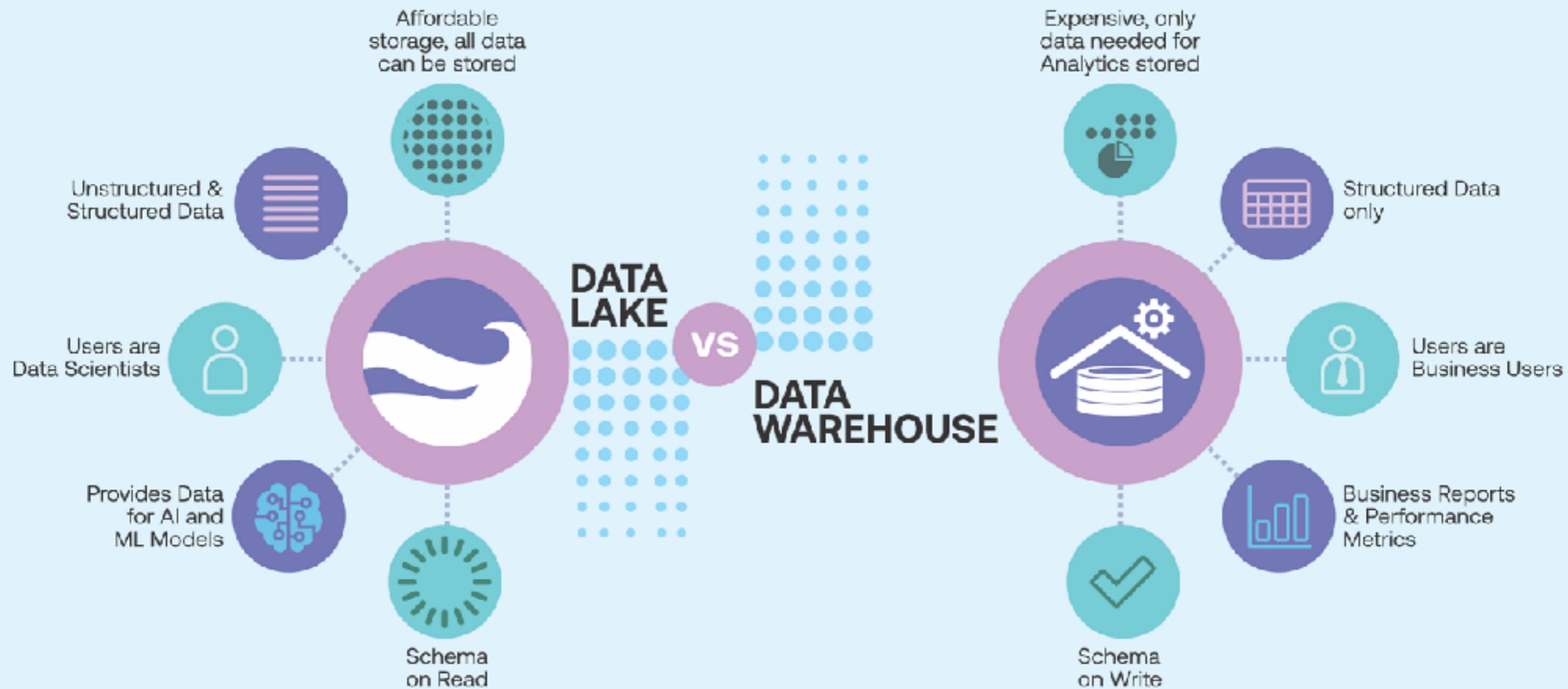
<https://panoply.io/data-warehouse-guide/data-mart-vs-data-warehouse/>



Sharing

© 2020 - 2024 Siam Chamnankit Company Limited. All rights reserved.

# Data Warehouse vs Data Lake



<https://panoply.io/data-warehouse-guide/data-mart-vs-data-warehouse/>



Sharing

© 2020 - 2024 Siam Chamnkit Company Limited. All rights reserved.

# Data pipeline

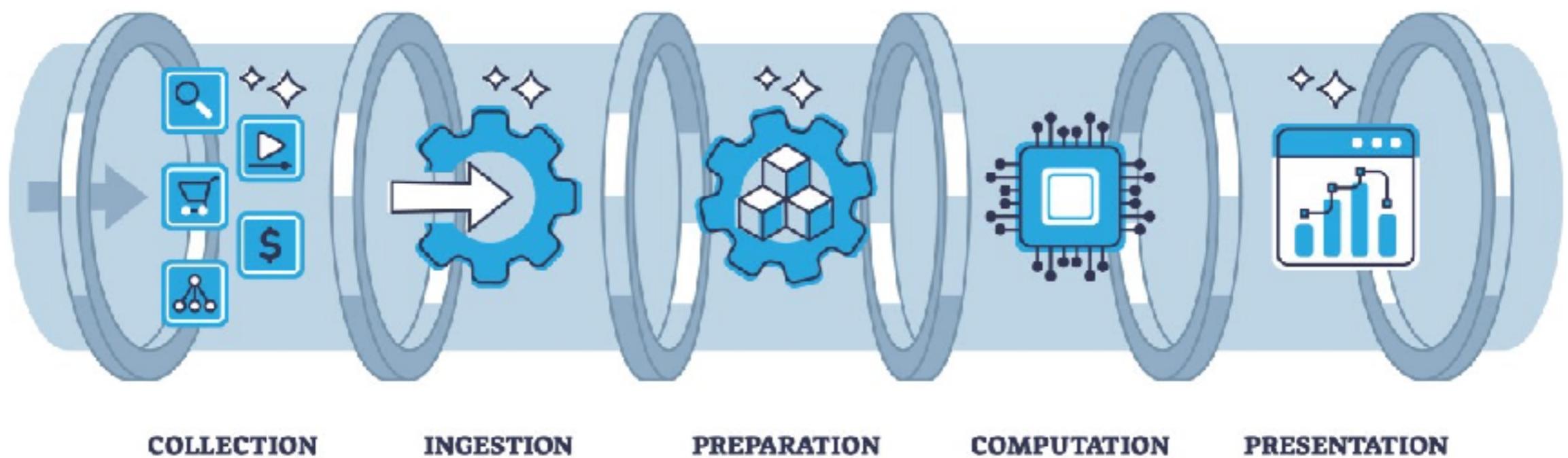


# Data pipeline

A **data pipeline** refers to a series of **processes** and tools used to automate the flow of data from source systems to destination systems, such as databases, data lakes, or data warehouses.



# Data pipeline



# Data Ingestion

The process of collecting data from various sources and bringing it into the pipeline.

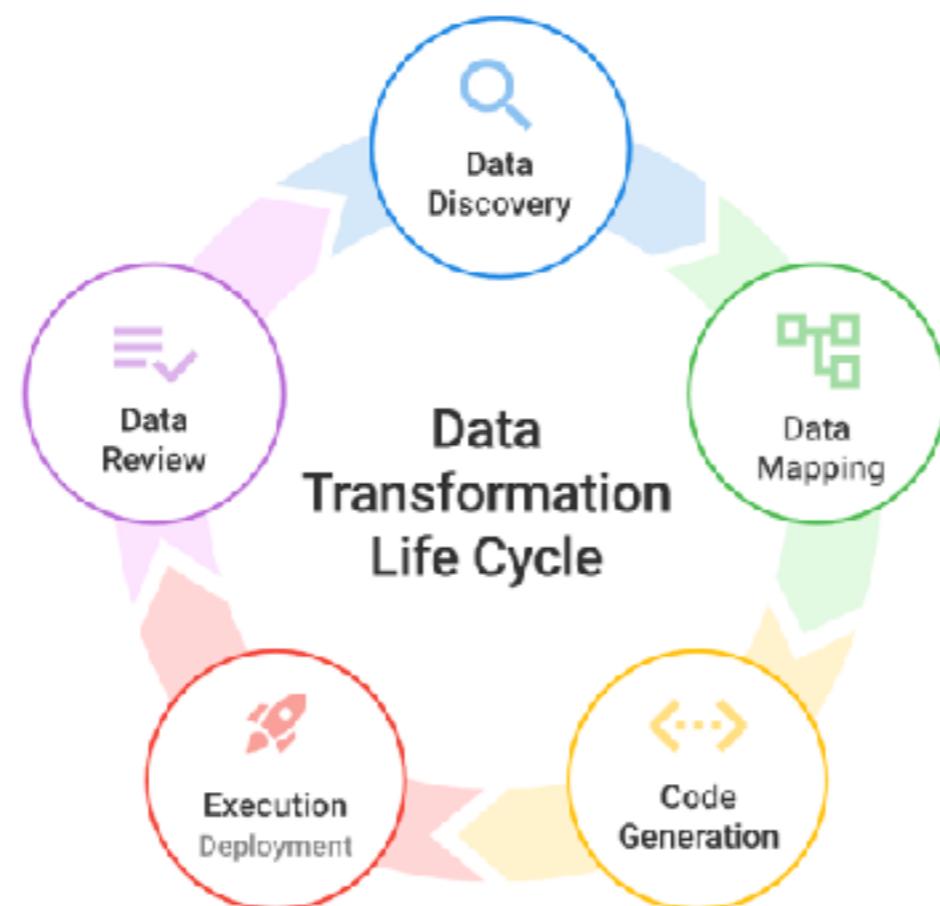
Batch  
processing

Real-time  
processing



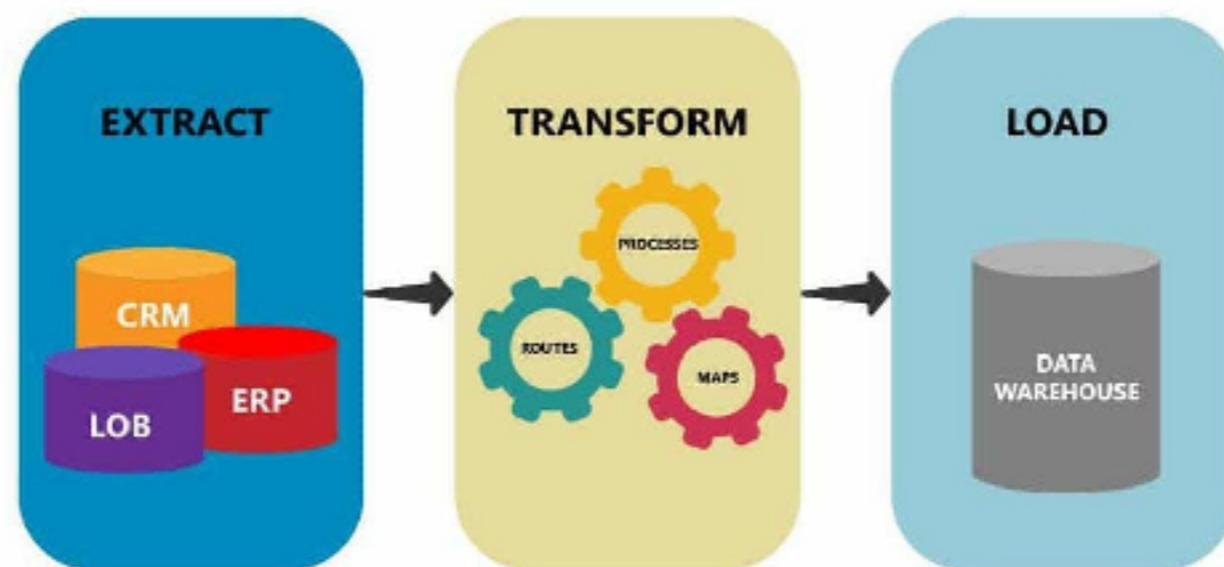
# Data Transformation

This step involves cleaning, validating, and transforming the data into a standardized format or structure suitable for analysis or storage.



# Data Loading

The transformed data is loaded into a target destination, which could be a data warehouse, database, or cloud storage platform.



# Orchestration

This component **manages the overall workflow** and scheduling of the pipeline, ensuring the proper execution and coordination of data processing tasks.



# How to cleaning data ?



# How to cleaning data ?



# Key Techniques with MS Excel

Remove  
duplication

Handle missing  
data

Standardize  
Data Formatting

Remove  
Unwanted  
Characters

Convert Text to  
Columns

Data Validation

Find and  
Replace  
Inconsistent

Fix Outliers

Data Cleansing



# Workshop



# **Big Data Analysis Techniques**

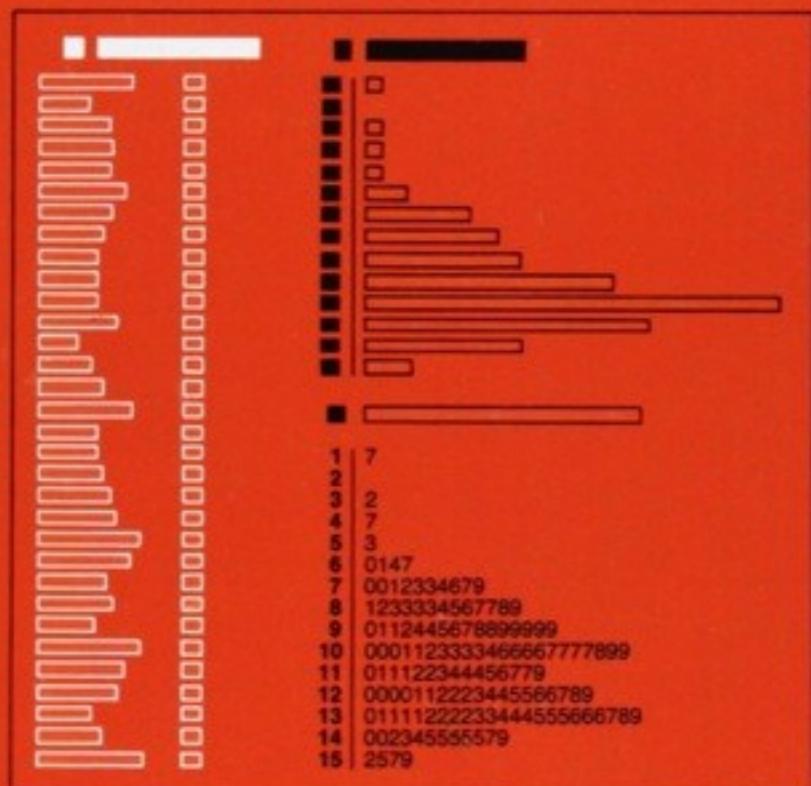


# **Exploratory Data Analysis (EDA)**



John W. Tukey

## EXPLORATORY DATA ANALYSIS



### Objectives of EDA :

- Suggest hypotheses about the causes of observed phenomena
- Assess assumptions on which statistical inference will be based
- Support the selection of appropriate tools and techniques
- Provide a basis for further data collection through surveys or experiments

John Tukey,  
*Exploratory Data Analysis*  
(New York: Pearson, 1977)



# Exploratory Data Analysis

Critical step in the **data science process** that involves summarizing the main characteristics of a dataset, often using visual methods, before applying more formal **statistical techniques**.

Summarize

Visualize

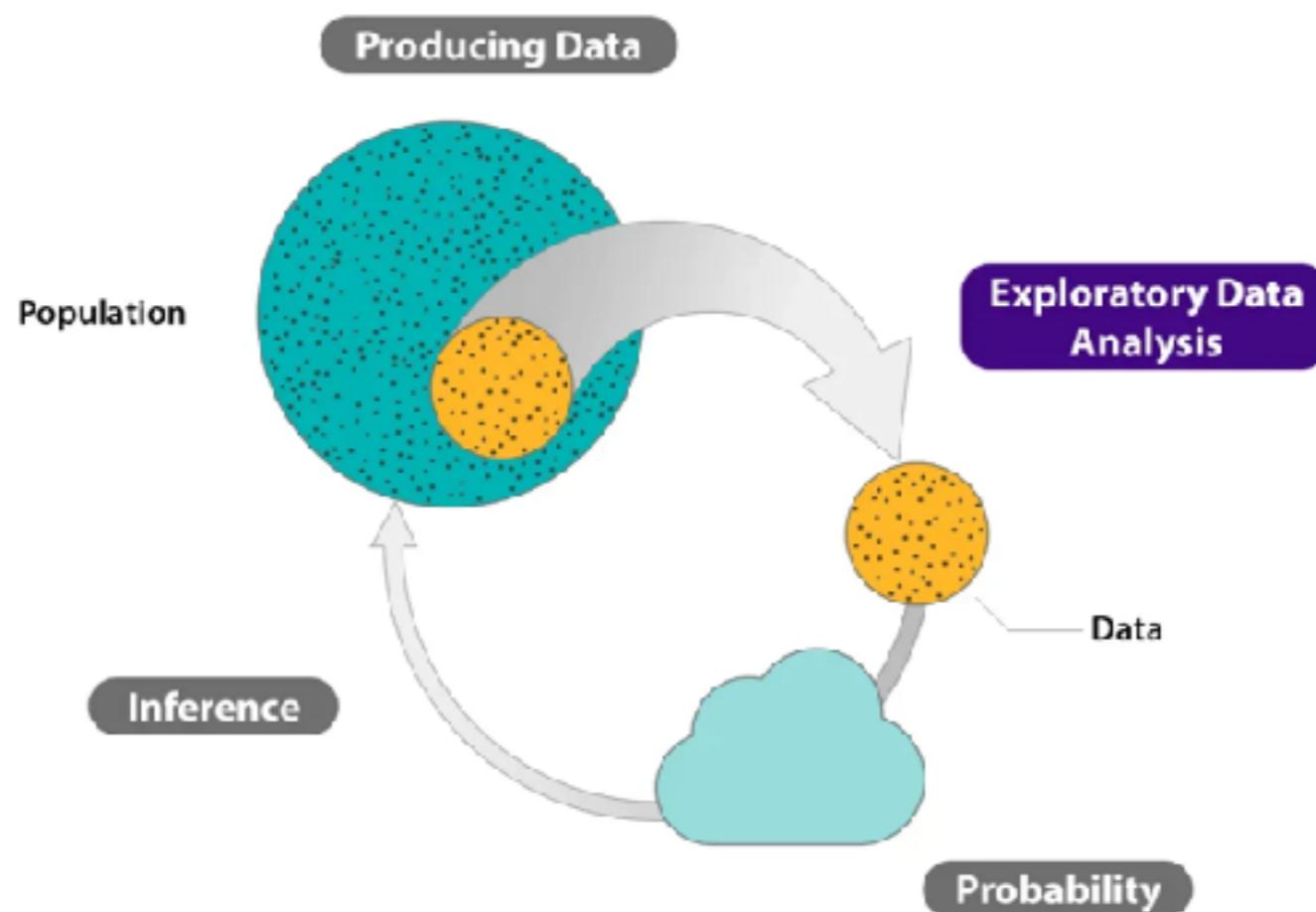
Identify

Missing data/values



# Exploratory Data Analysis

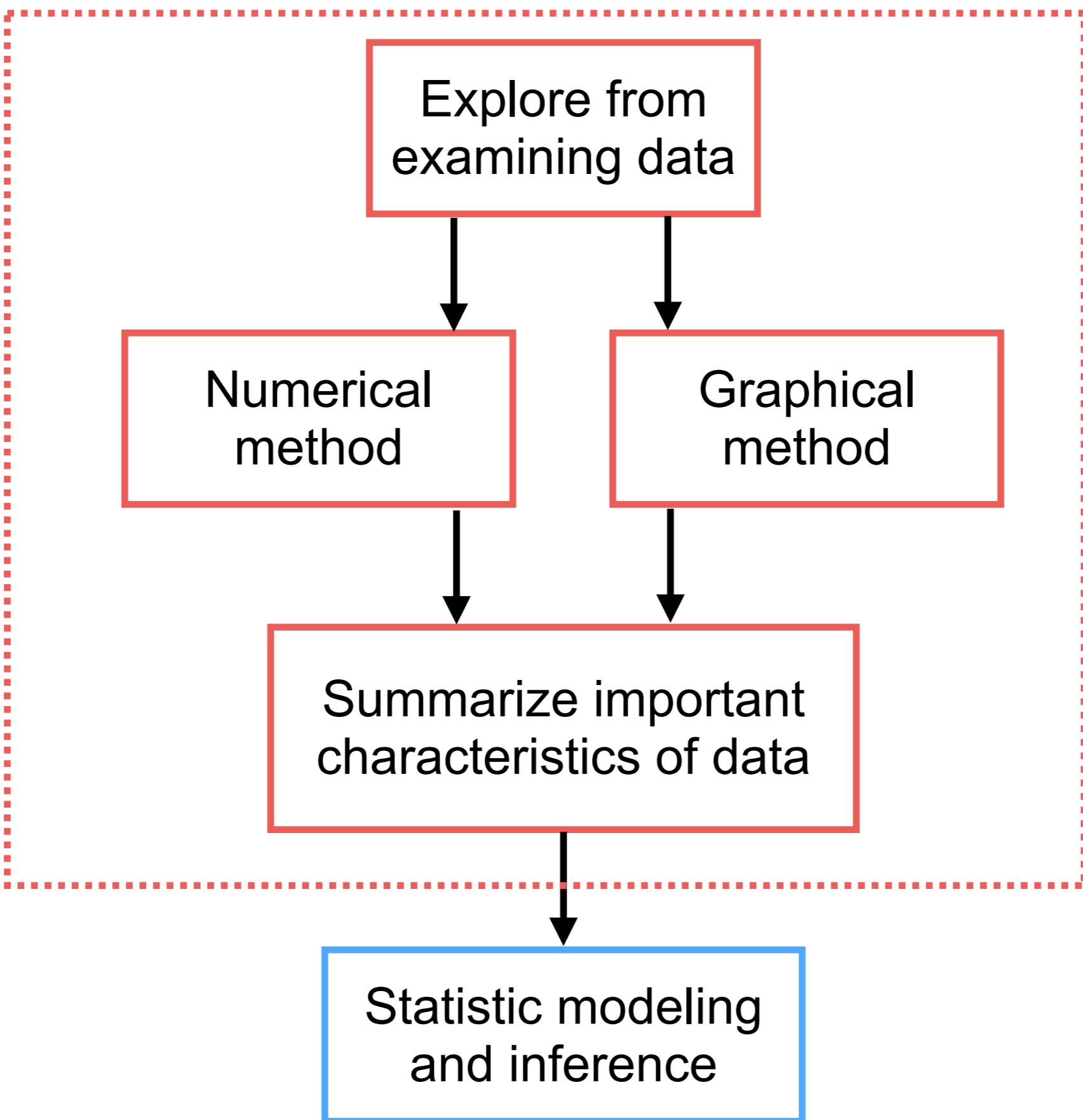
EDA helps in understanding the data, identifying patterns, detecting anomalies, and testing hypotheses.



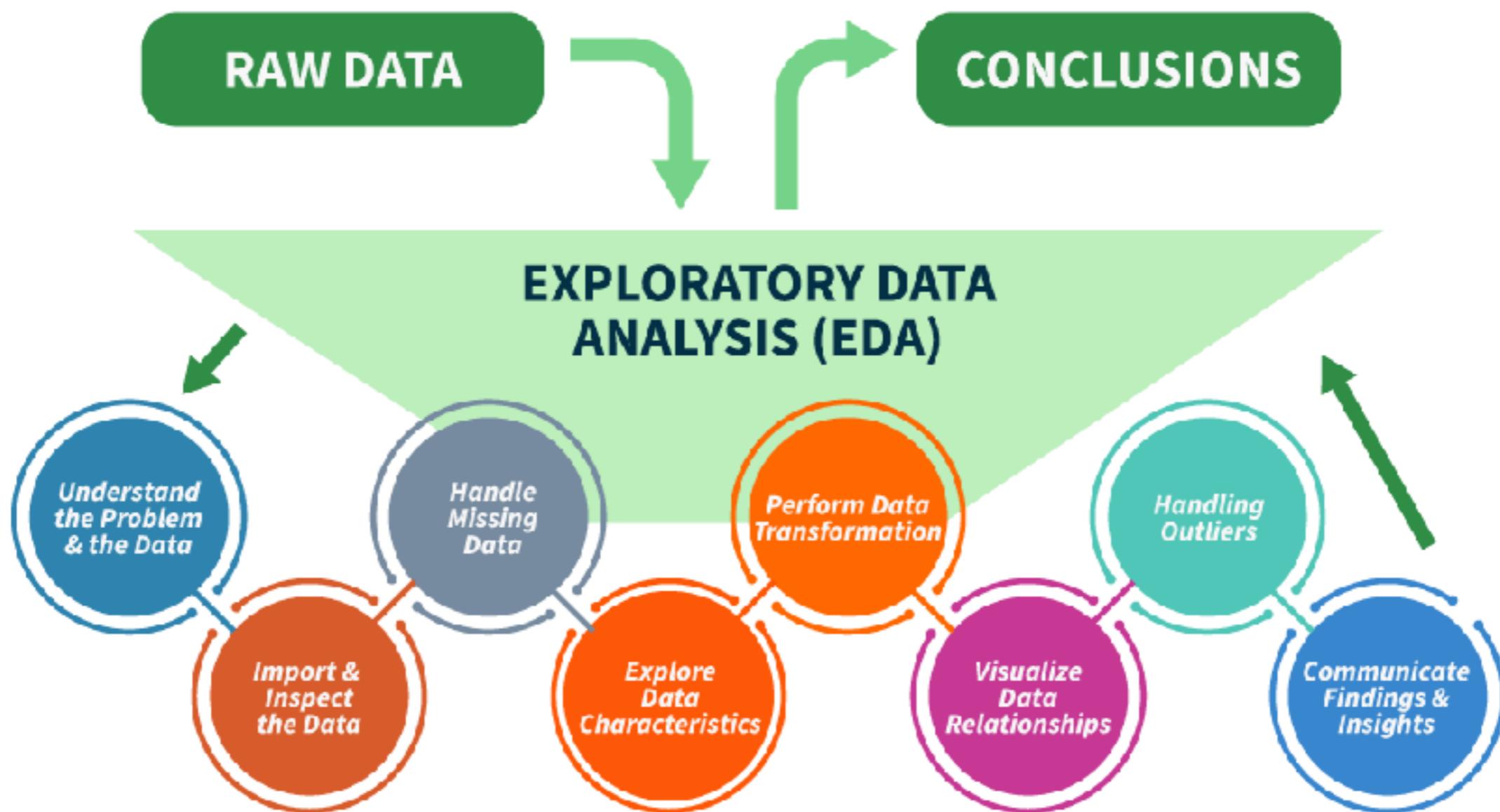
# Why EDA ?

Better understand of data  
Validate assumption about data  
Identify errors and outliers in data  
Interesting trends, patterns and relationships  
Insight and new ideas/hypothesis  
Input for statistic modeling





# Steps for Performing Exploratory Data Analysis



# Steps

Problem and data understanding

    Data cleaning

    Pattern discovery

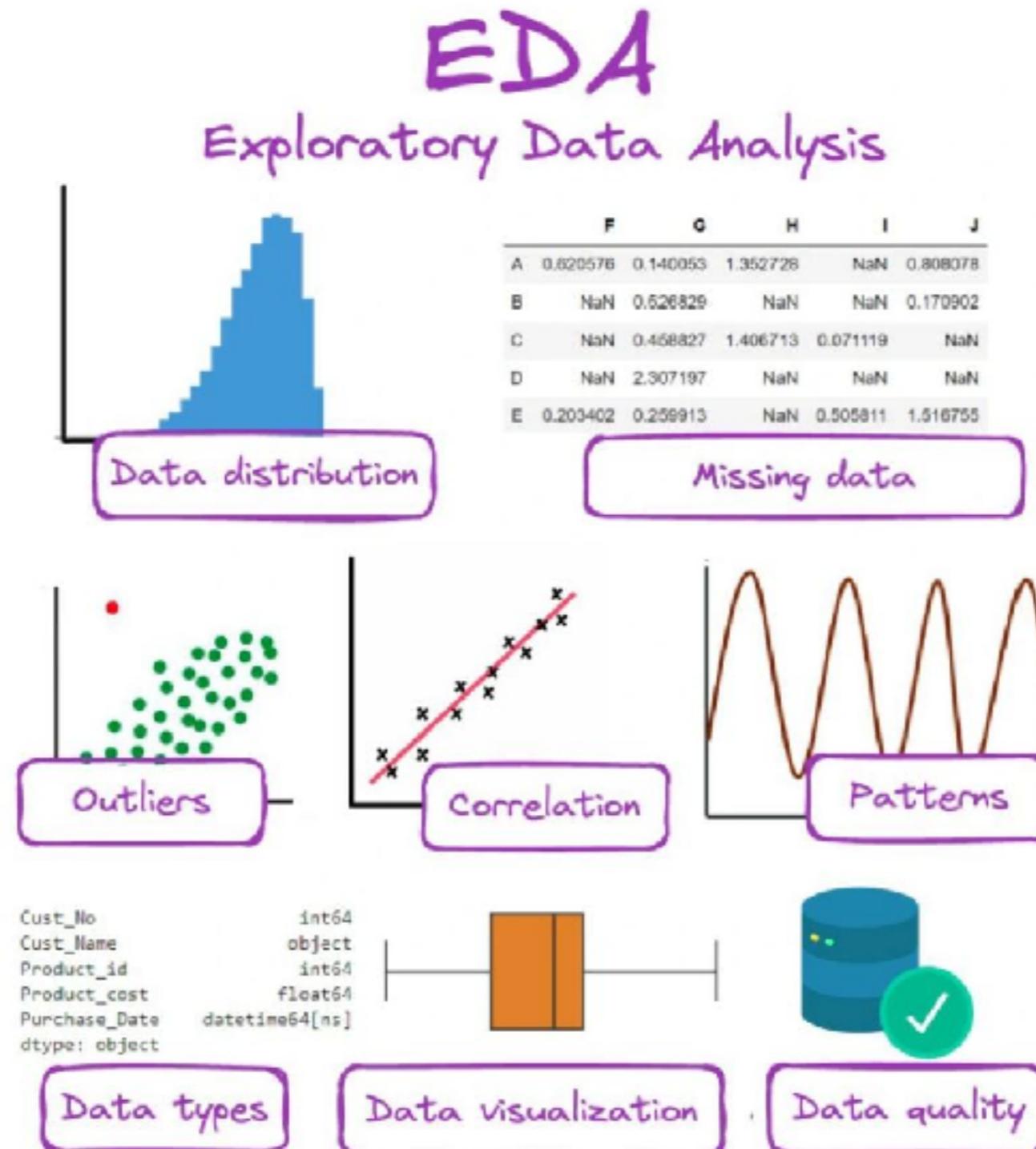
    Data visualization

    Model selection

    Quality control



# Workshop with EDA



# **Exploratory Data Analysis (EDA)**



# Analysis Process

Iterative process

Data  
Transformation

Data  
Analysis

Data  
Visualization

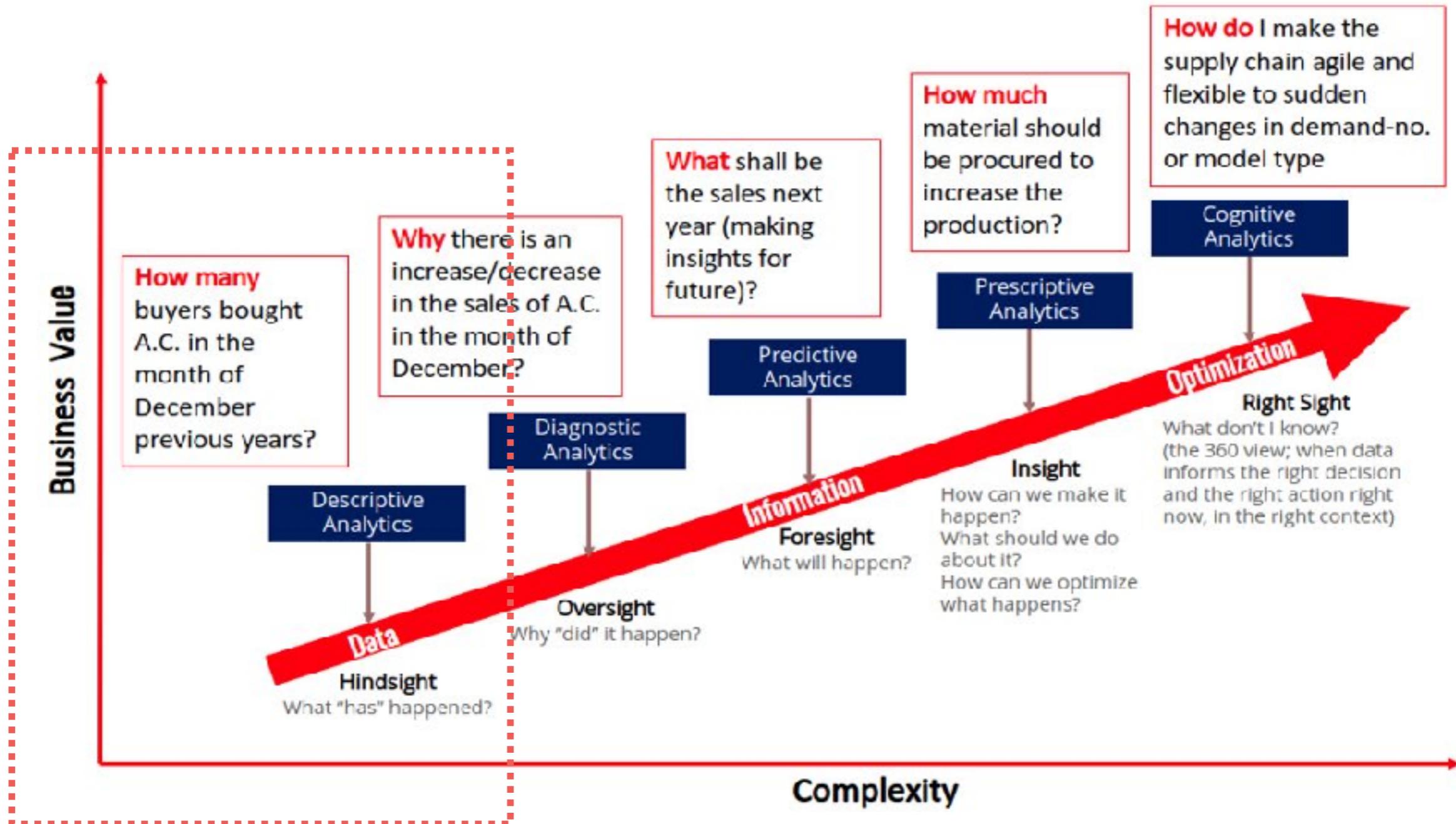
Data manipulation  
Data wrangling



# **Statistical Analysis for data analysis**



# Key Techniques



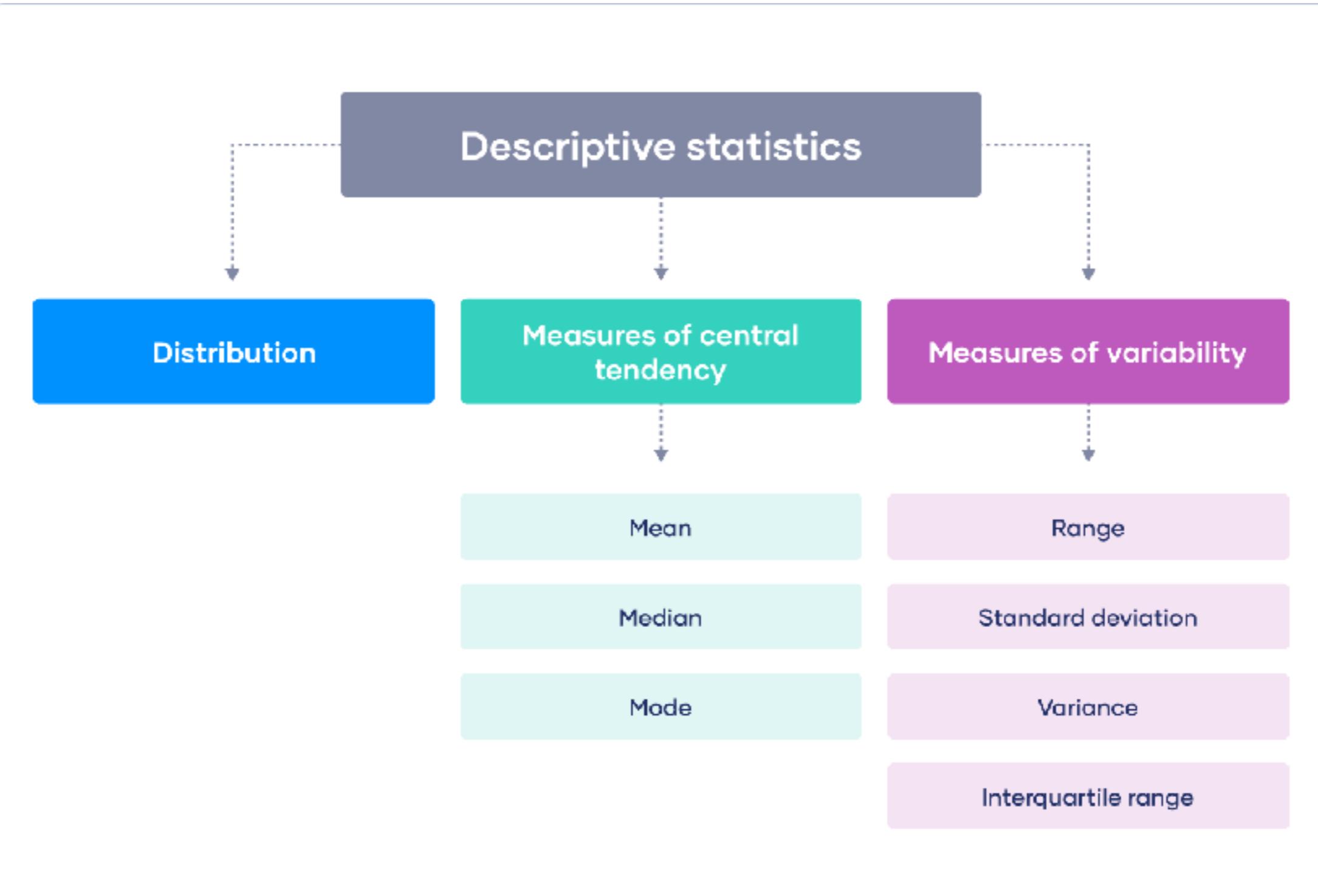
# **Descriptive statistic**

Statistical methods used to **summarize** and describe the **main features** of a dataset.

Help in **understanding the data** by providing a clear overview through numerical measures and visual representations.



# Types of descriptive statistic



# Measures of Central Tendency

## Mean

The average of all data points

## Median

The middle value when data points are arranged in order

## Mode

The value that appears most frequently in the dataset



# Mean

2, 3, 4

1      2      ③

$$2 + 3 + 4 = 9$$

$$\text{MEAN} = 9 \div 3 = 3$$

wikiHow to Find Mean, Median, and Mode

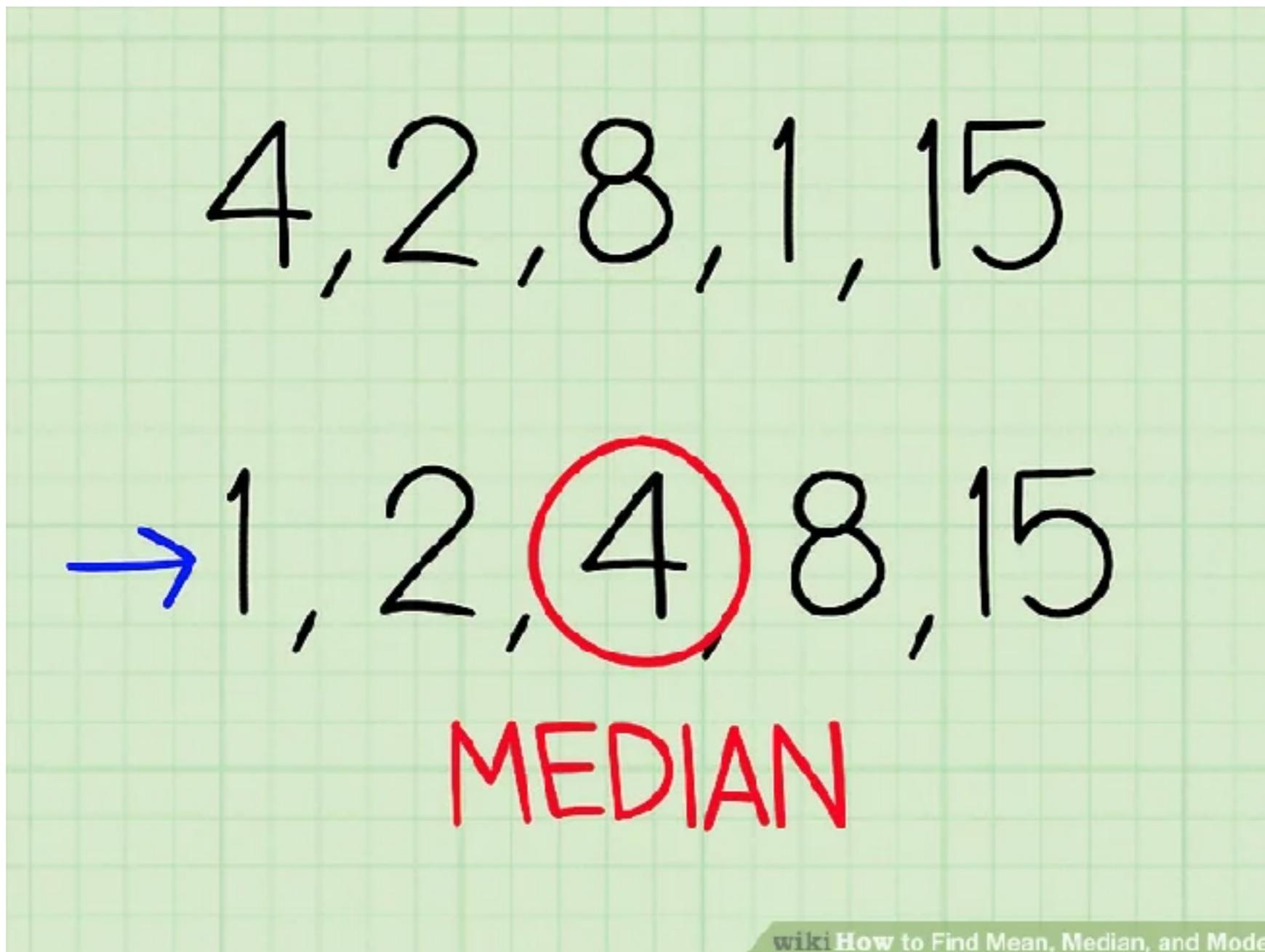
<https://www.wikihow.com/Find-Mean,-Median,-and-Mode>



Sharing

© 2020 - 2024 Siam Chamnankit Company Limited. All rights reserved.

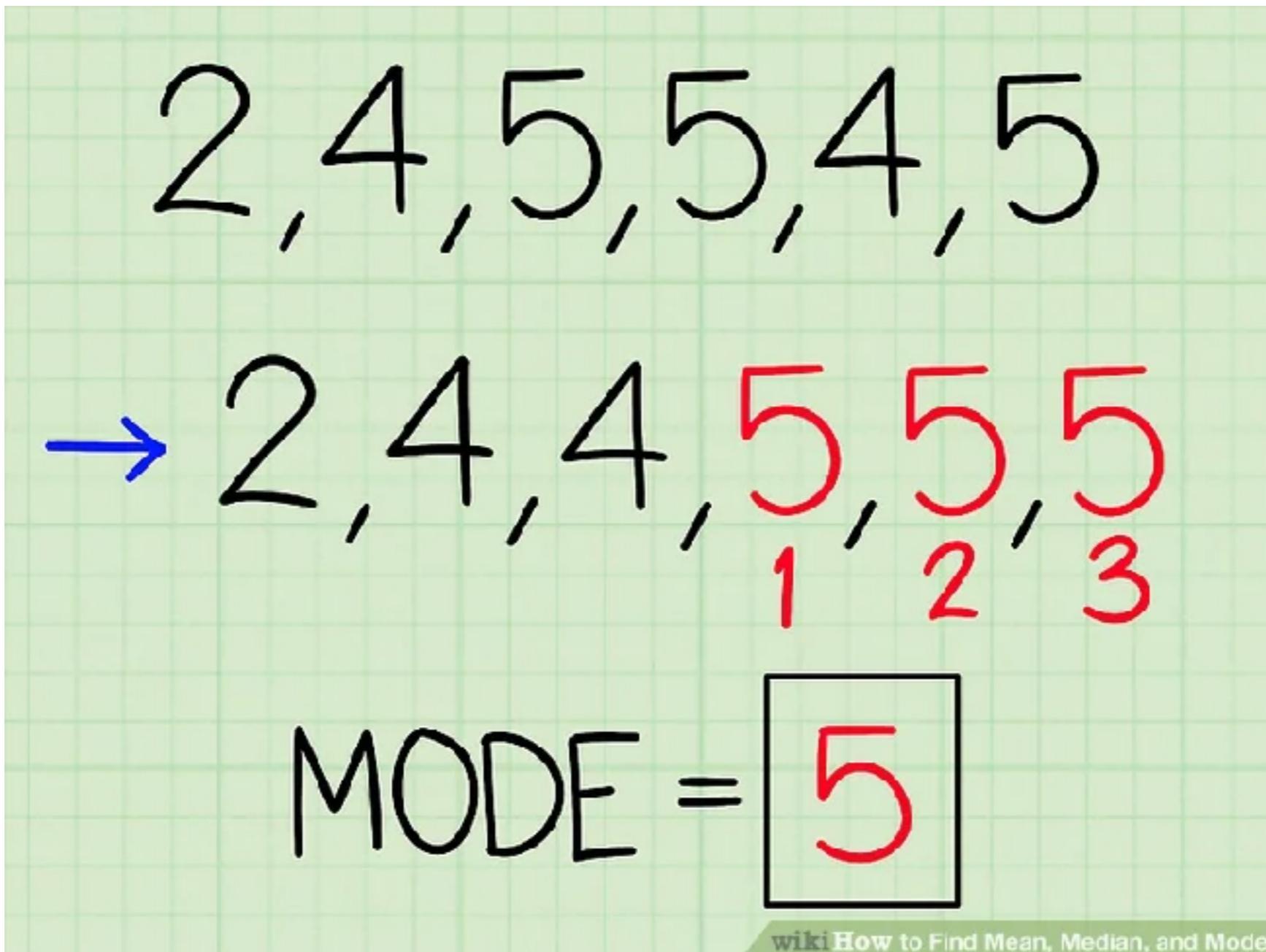
# Median



<https://www.wikihow.com/Find-Mean,-Median,-and-Mode>



# Mode



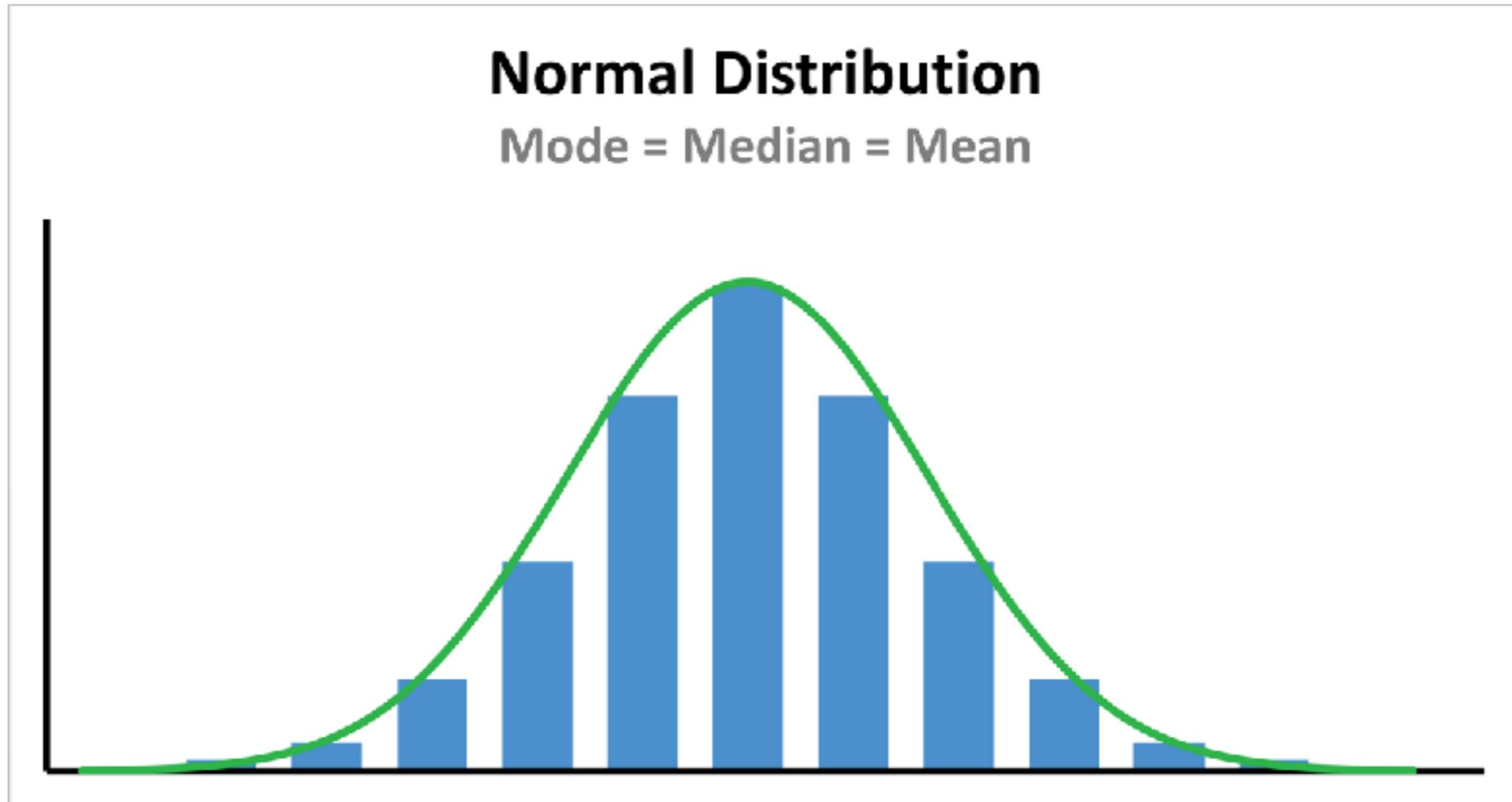
<https://www.wikihow.com/Find-Mean,-Median,-and-Mode>



Sharing

© 2020 - 2024 Siam Chamnankit Company Limited. All rights reserved.

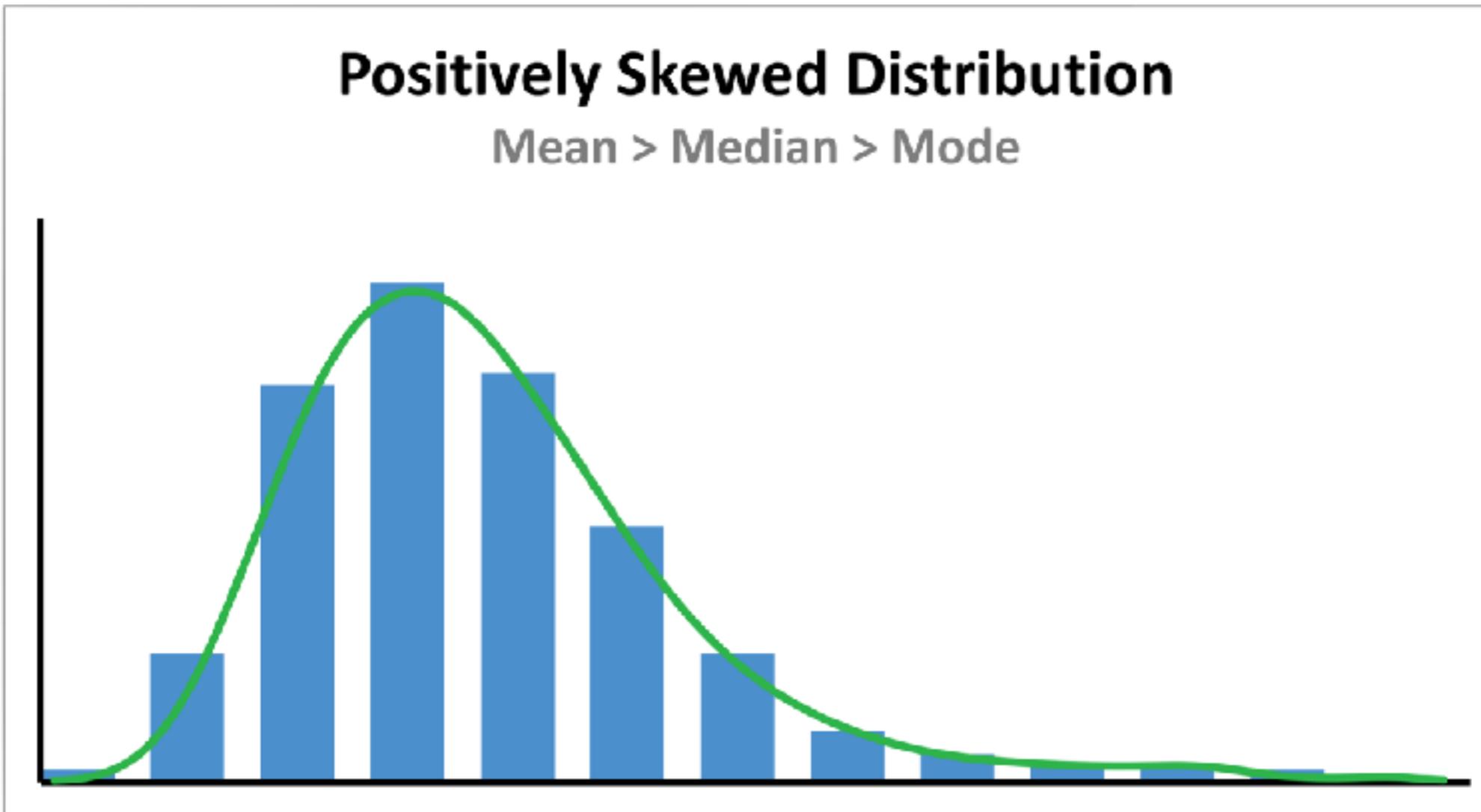
# Normal Distribution (Gaussian)



# Positive Skewed Distribution

## Positively Skewed Distribution

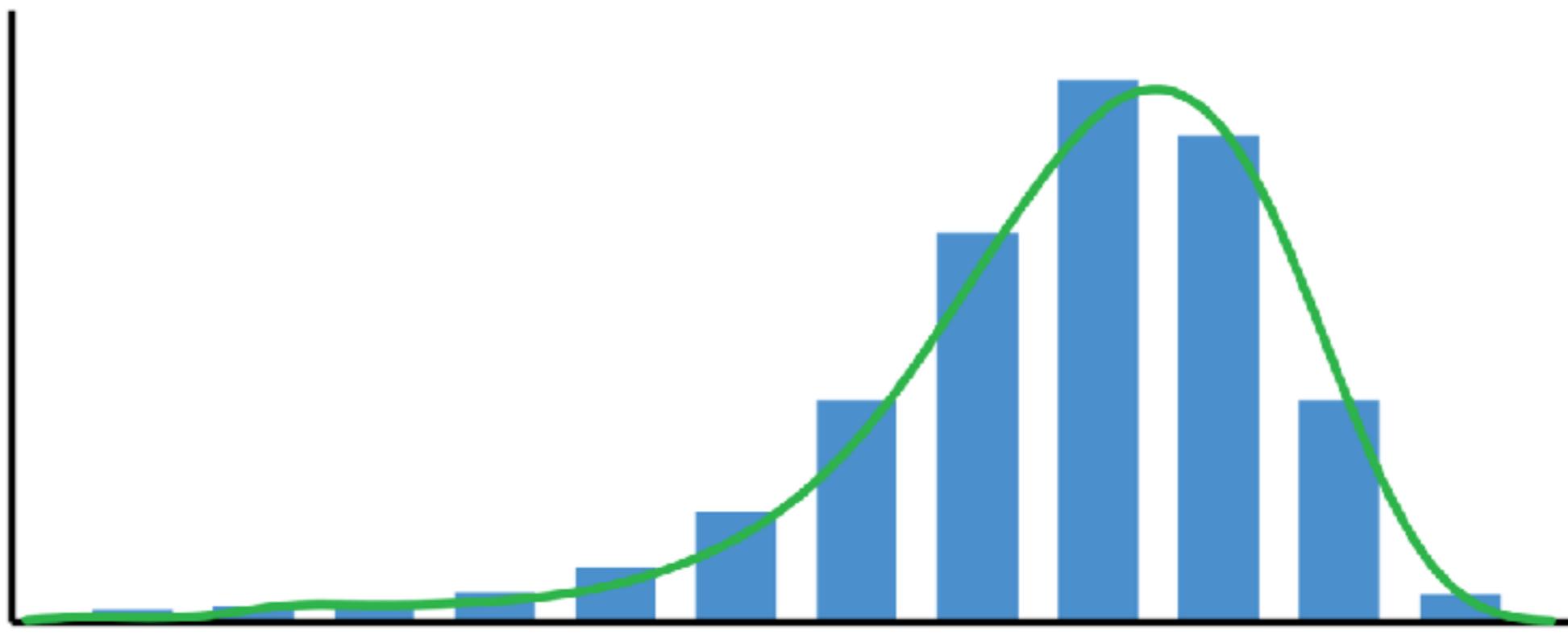
Mean > Median > Mode



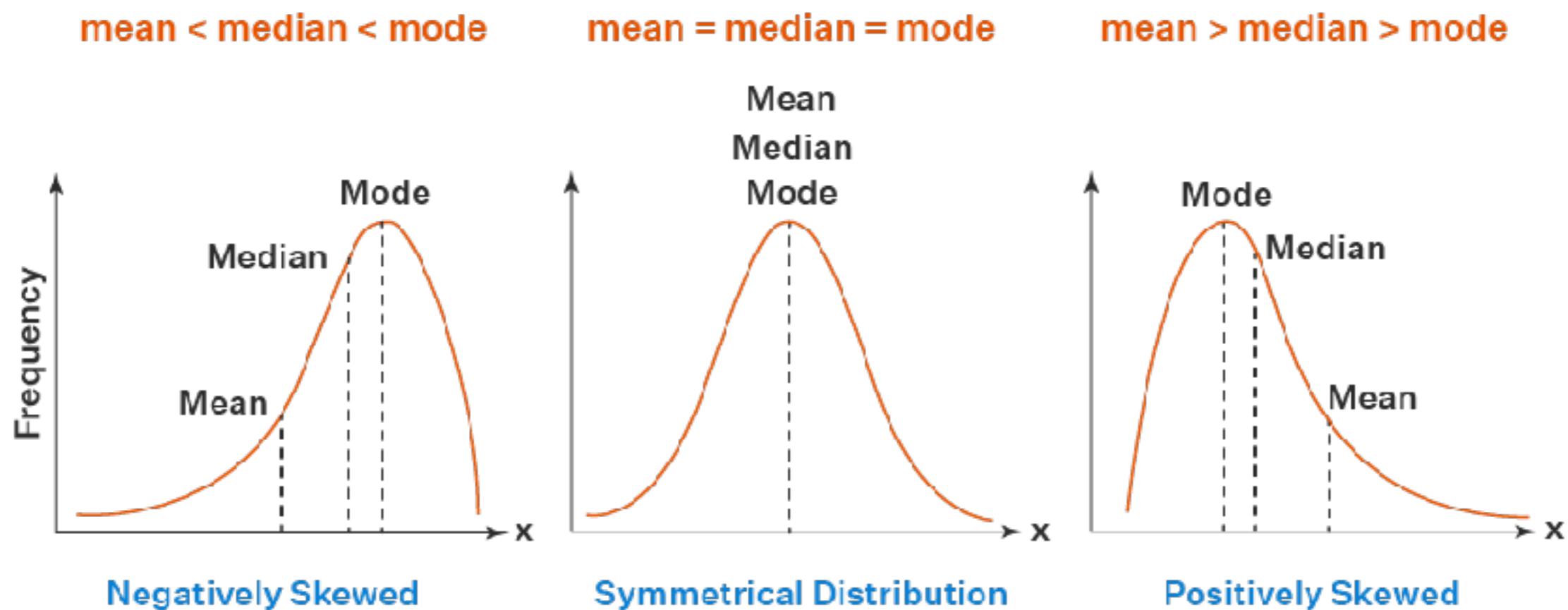
# Negative Skewed Distribution

## Negatively Skewed Distribution

Mean < Median < Mode



# Distribution



# Measures of Variability

## Range

The difference between the maximum and minimum values

## Variance

The average of the squared differences from the mean

## Standard Deviation

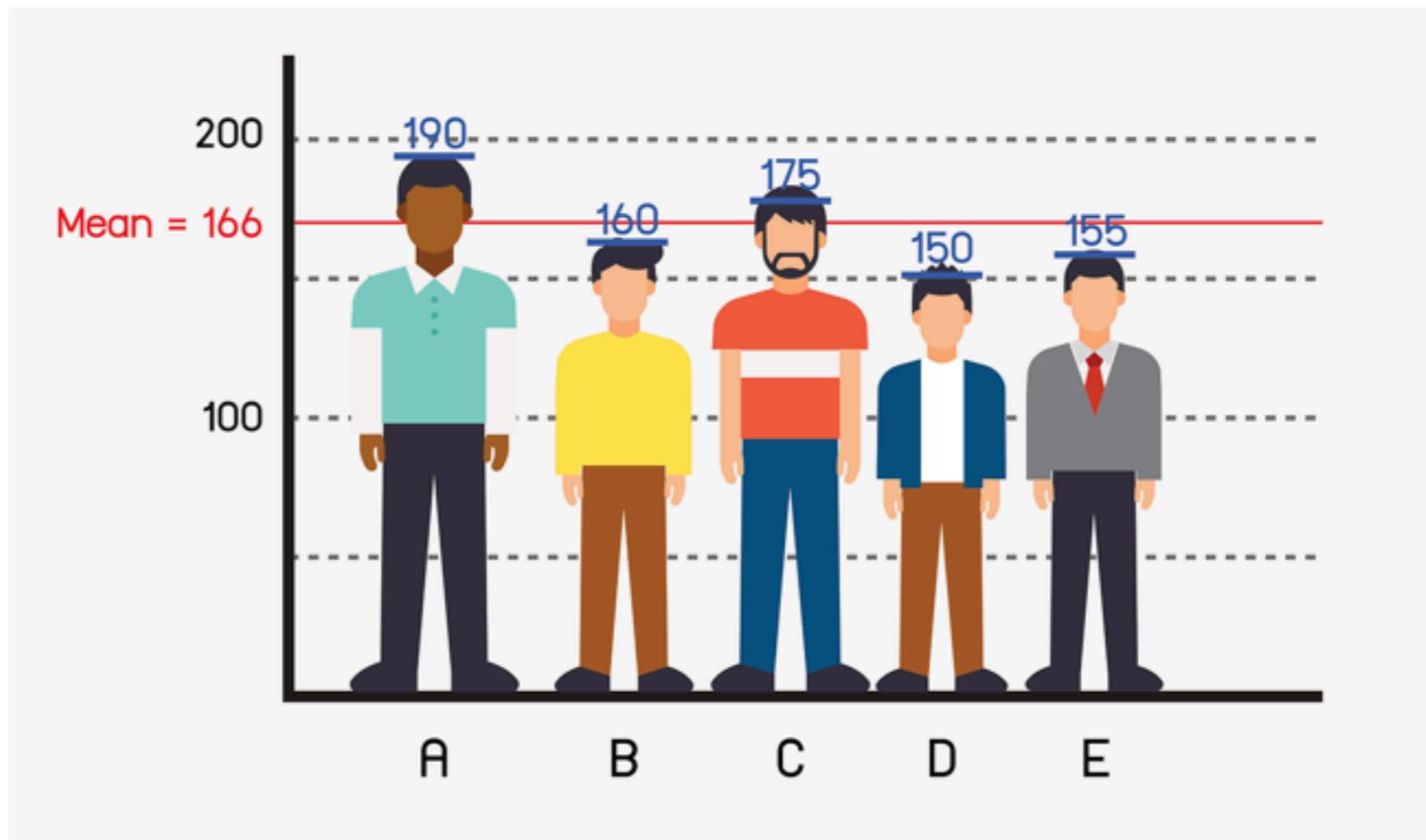
The square root of the variance, showing how much the values deviate from the mean

## Interquartile Range (IQR)

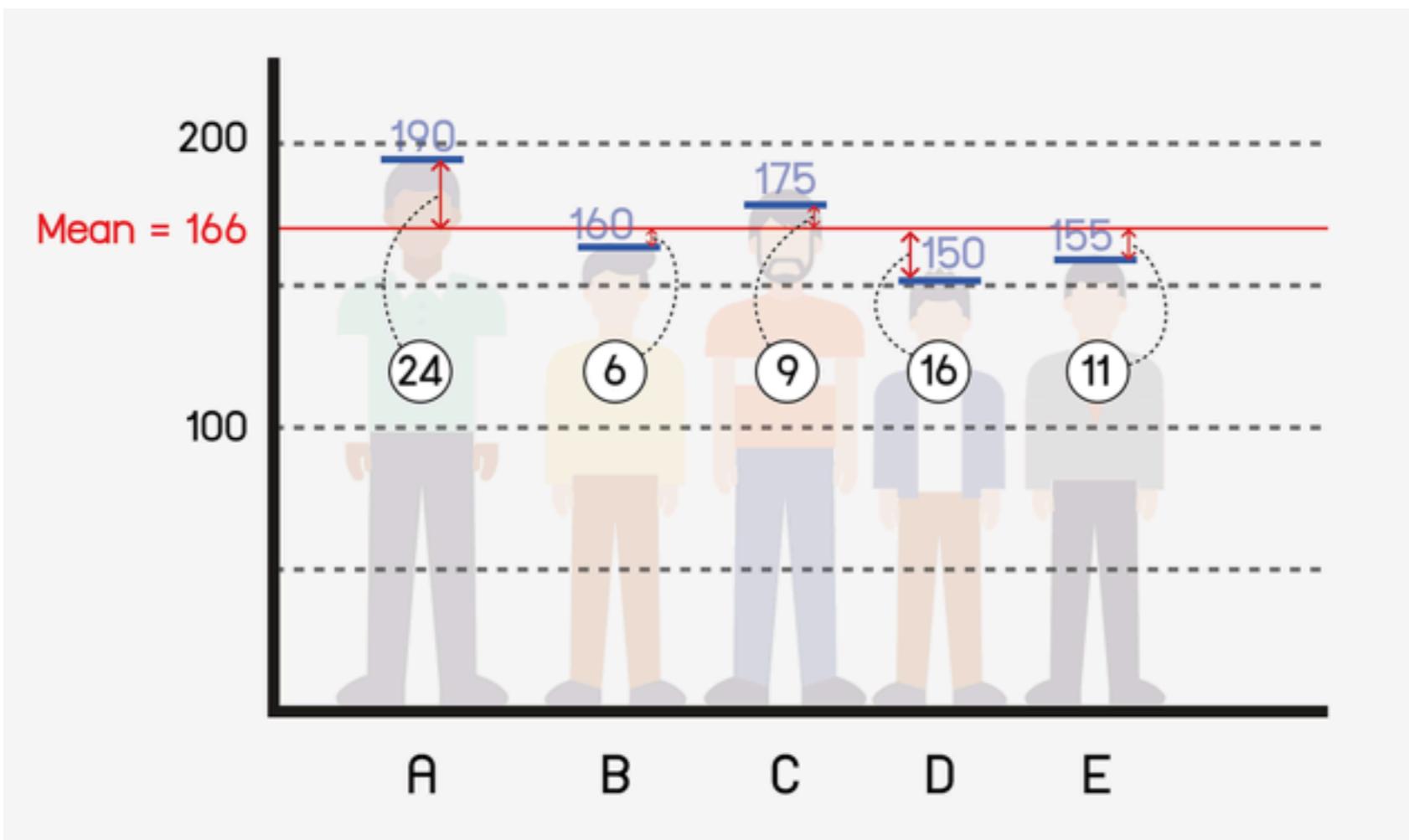
The difference between the first (25th percentile) and third (75th percentile) quartiles



# Variance ? (1)



# Variance ? (2)



# Variance ? (3)

ค่าเฉลี่ย คือ  
Mean

$$\frac{190 + 160 + 175 + 150 + 155}{5} = 166$$

ค่าแปรปวน คือ  
Variance

$$\frac{24^2 + (-6)^2 + 9^2 + (-16)^2 + (-11)^2}{5} = 214$$

ค่าเบี่ยงเบนมาตรฐาน คือ  
Standard Deviation

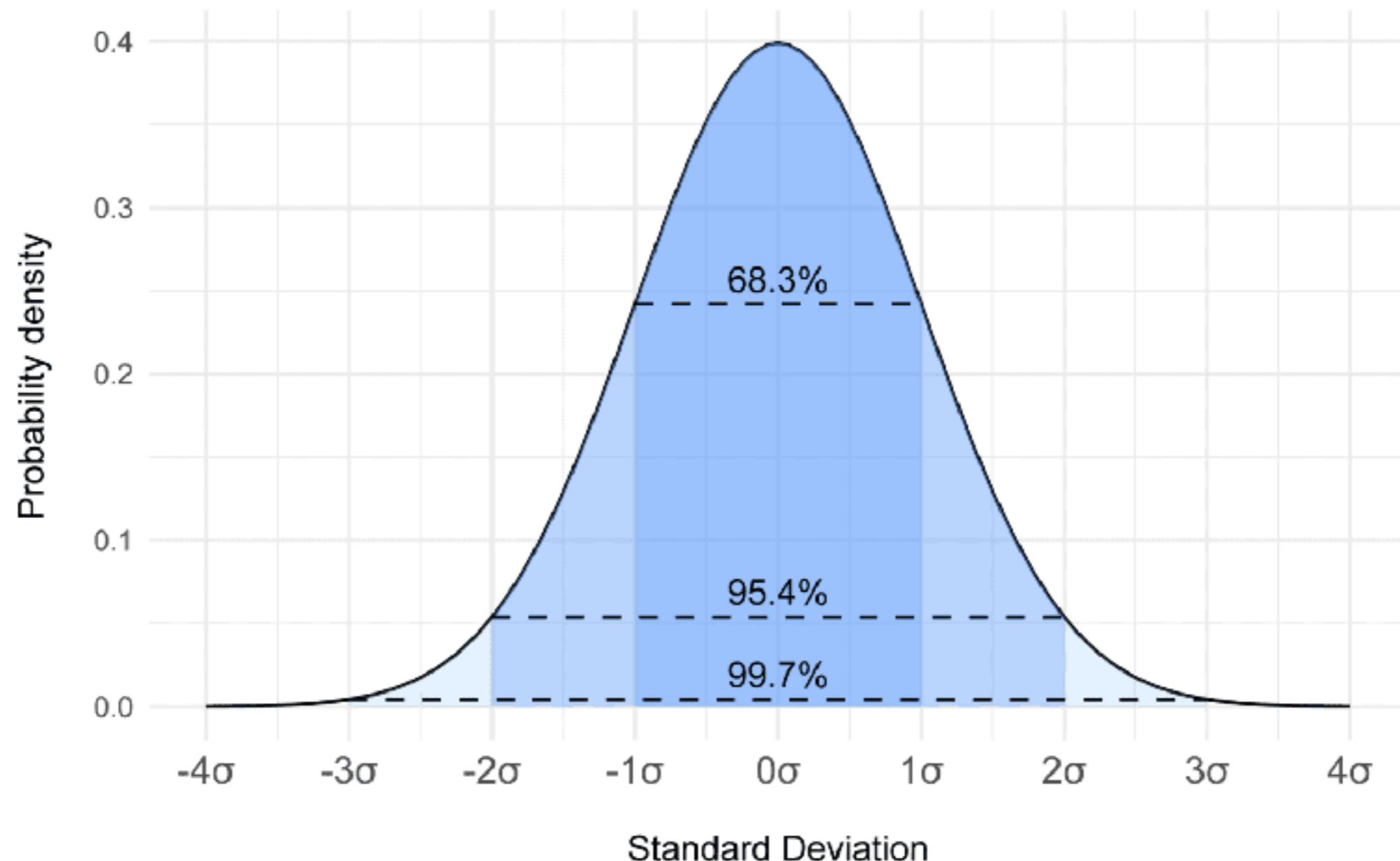
$$\sqrt{214} = 14.63$$



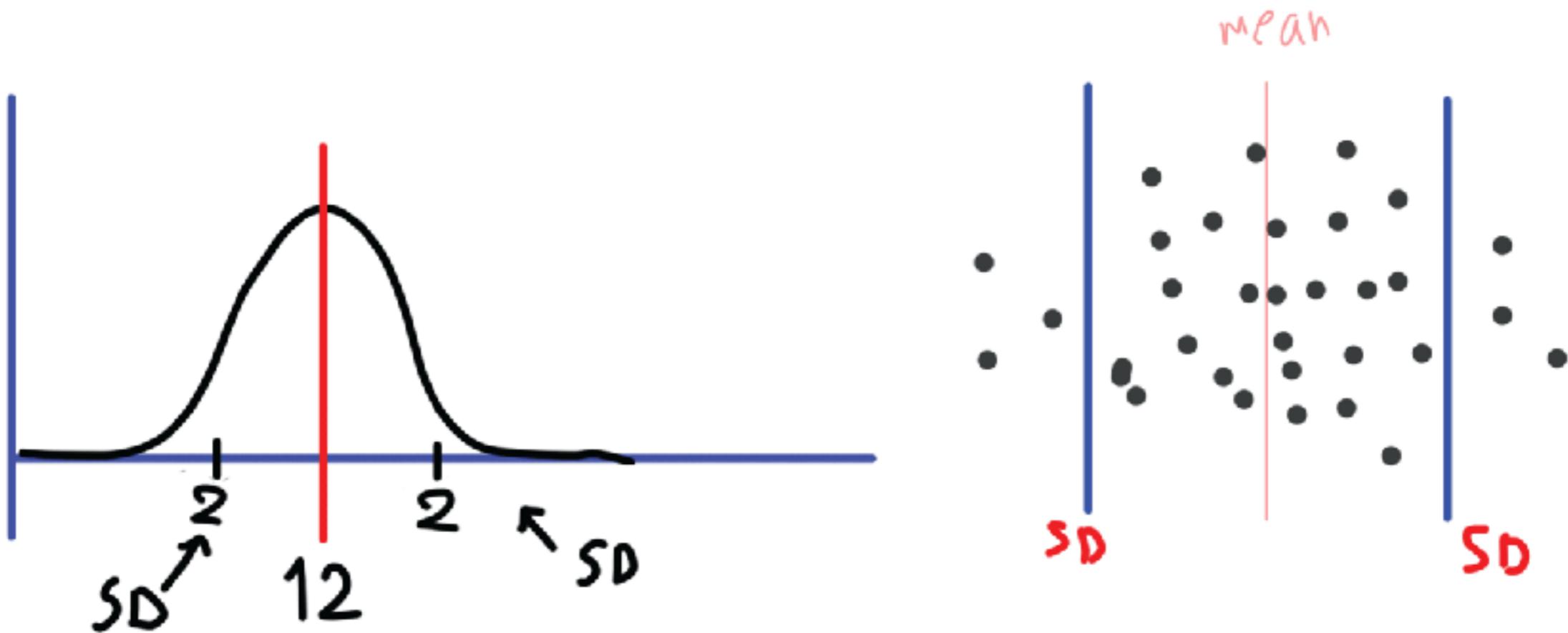
# Workshop



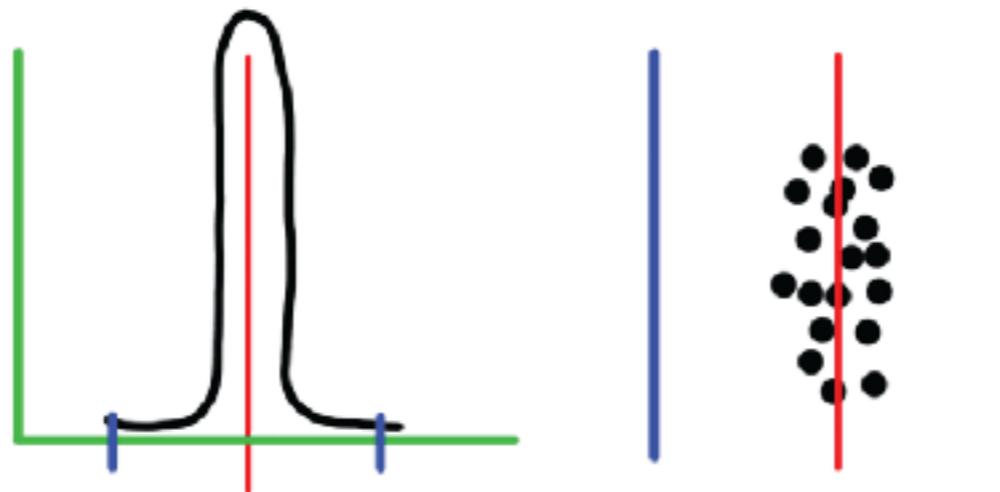
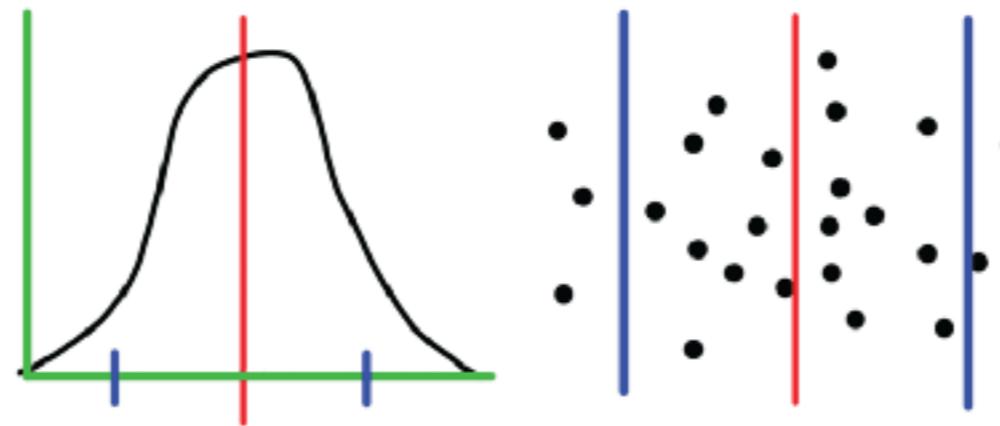
# Standard Deviation (1)



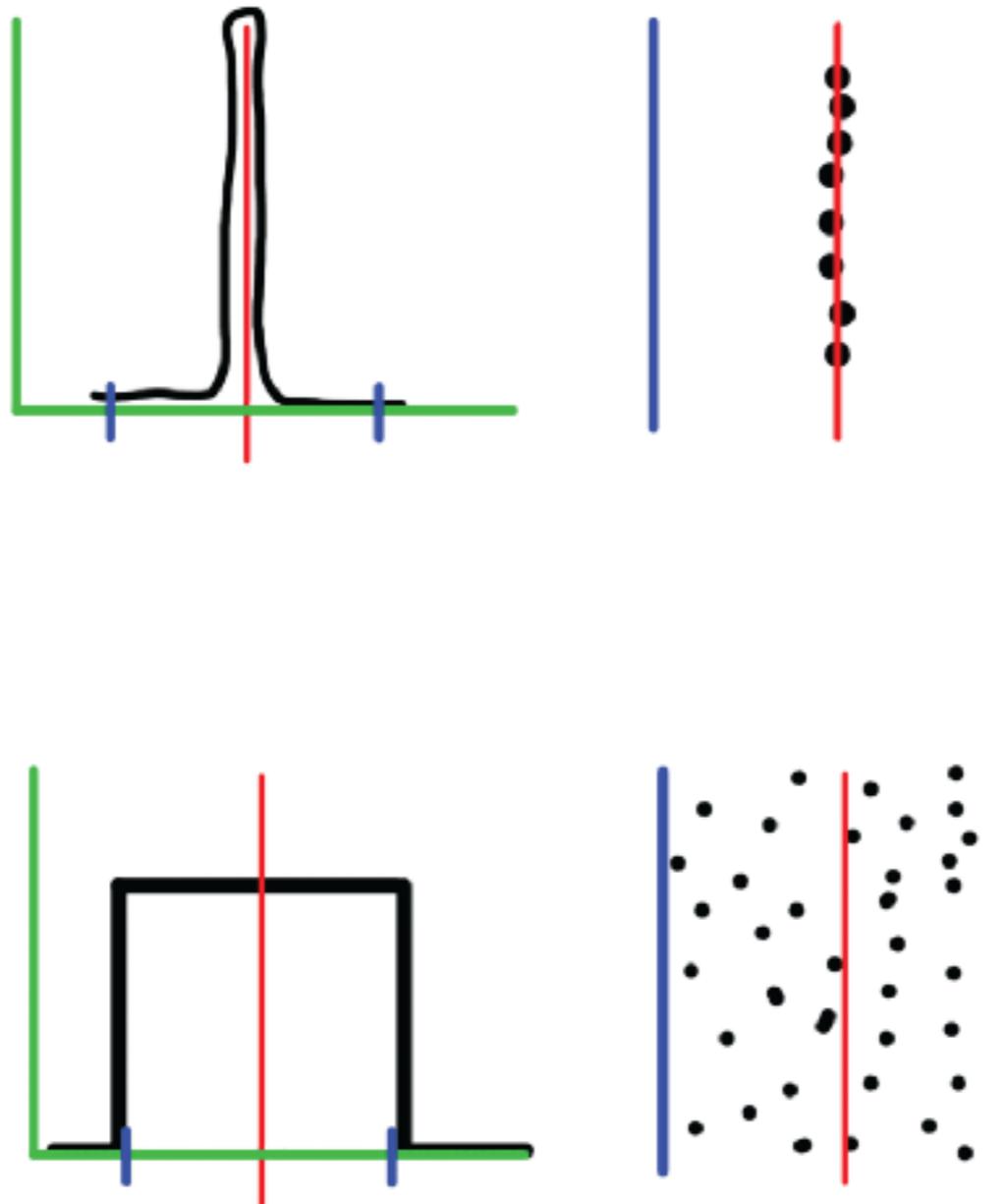
# Standard Deviation (2)



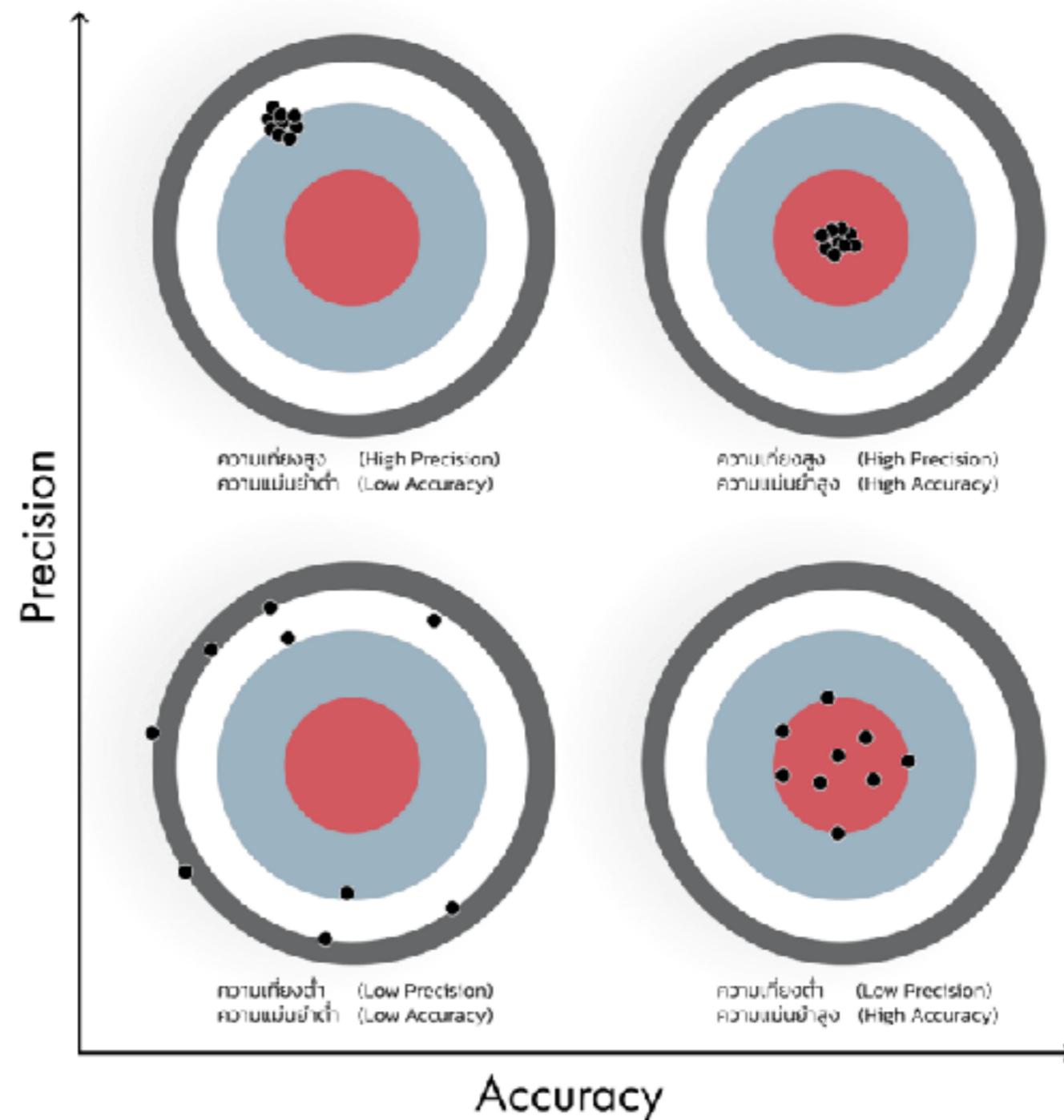
# Standard Deviation (3)



# Standard Deviation (4)



# Standard Deviation and Precision



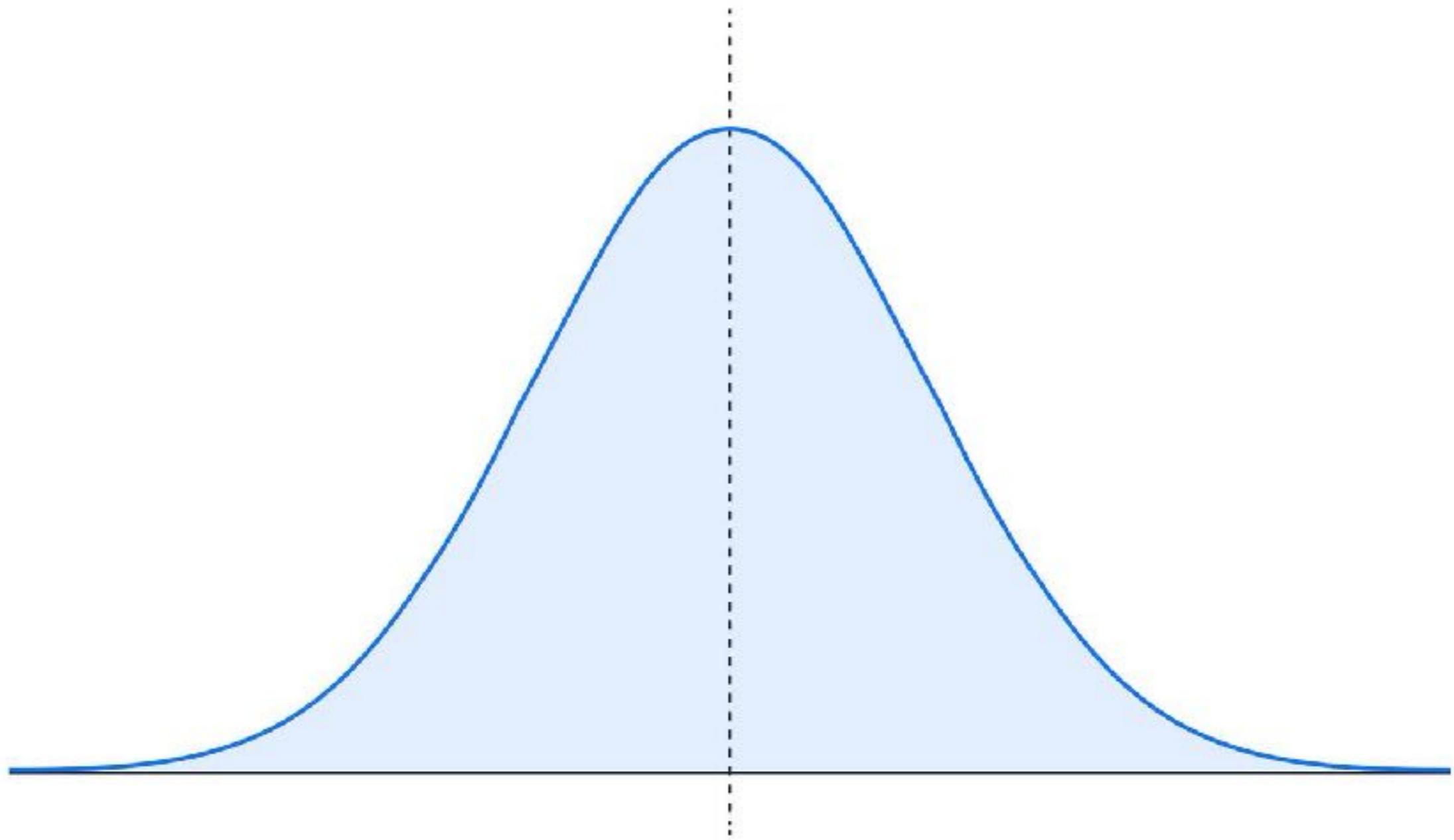
# Percentile

Measure used in statistics to indicate the relative standing of a value within a dataset.

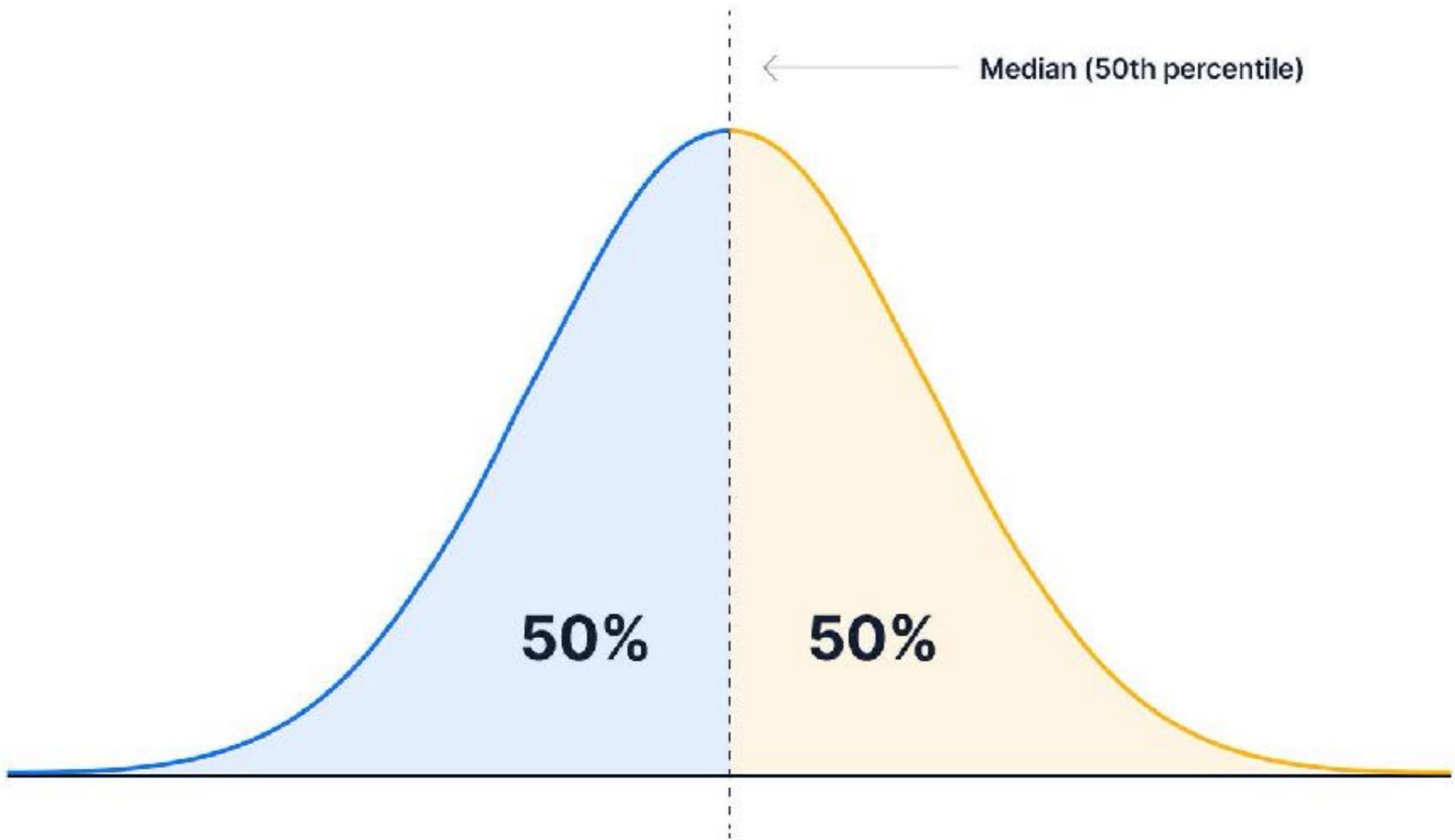
It represents the **percentage of values** in the data that are below a certain point.



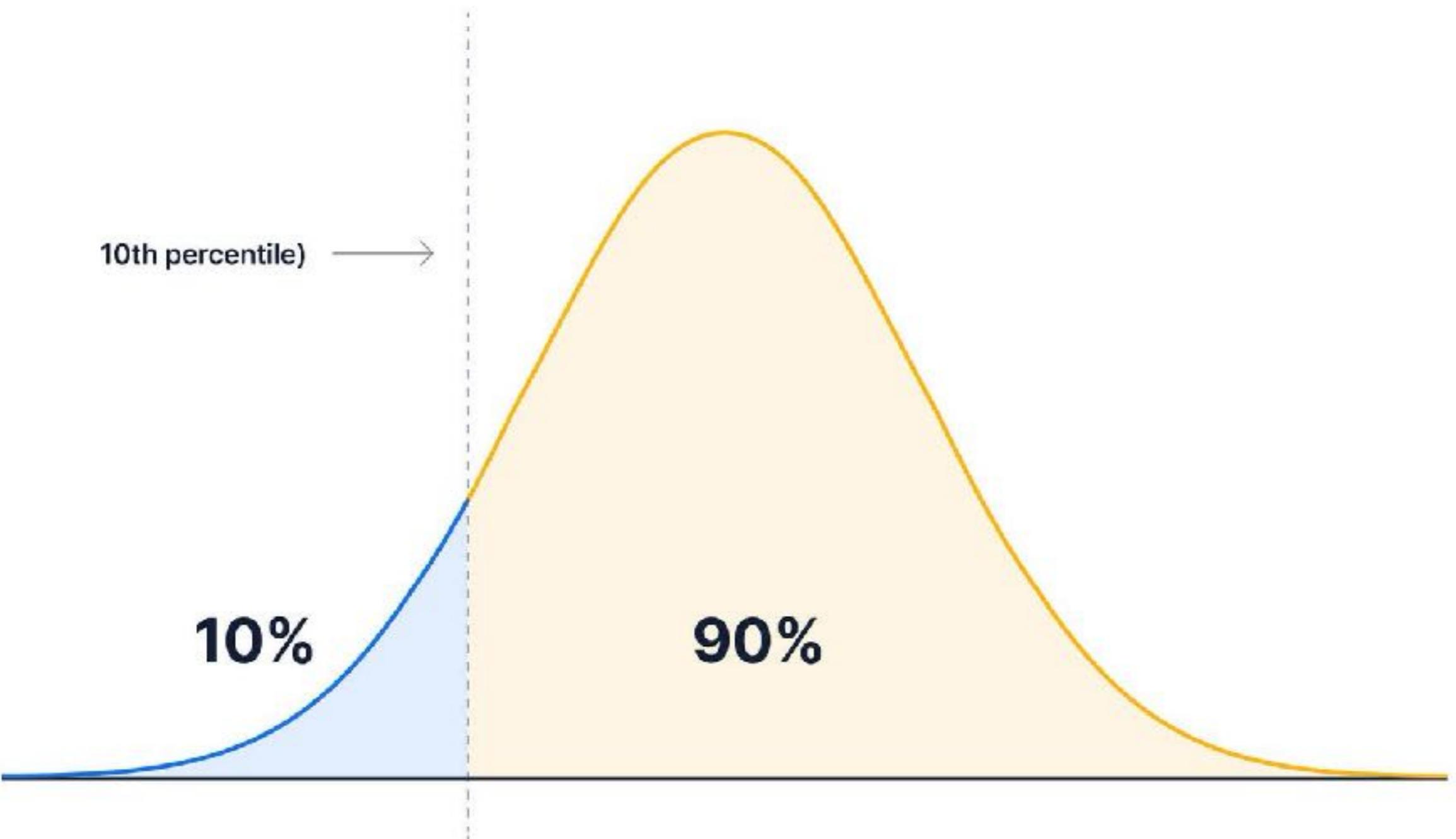
# Normal distribution



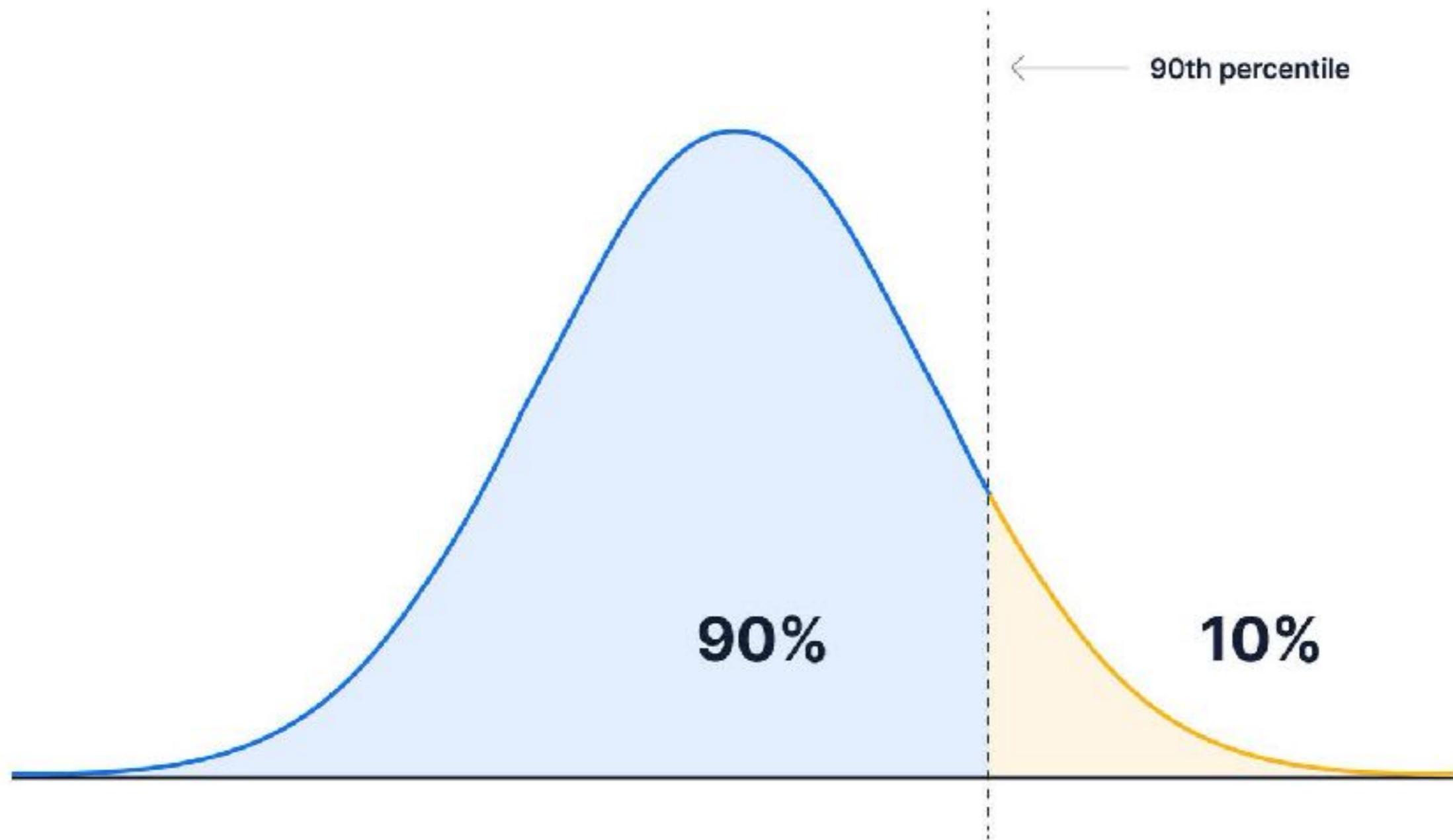
# 50th percentile



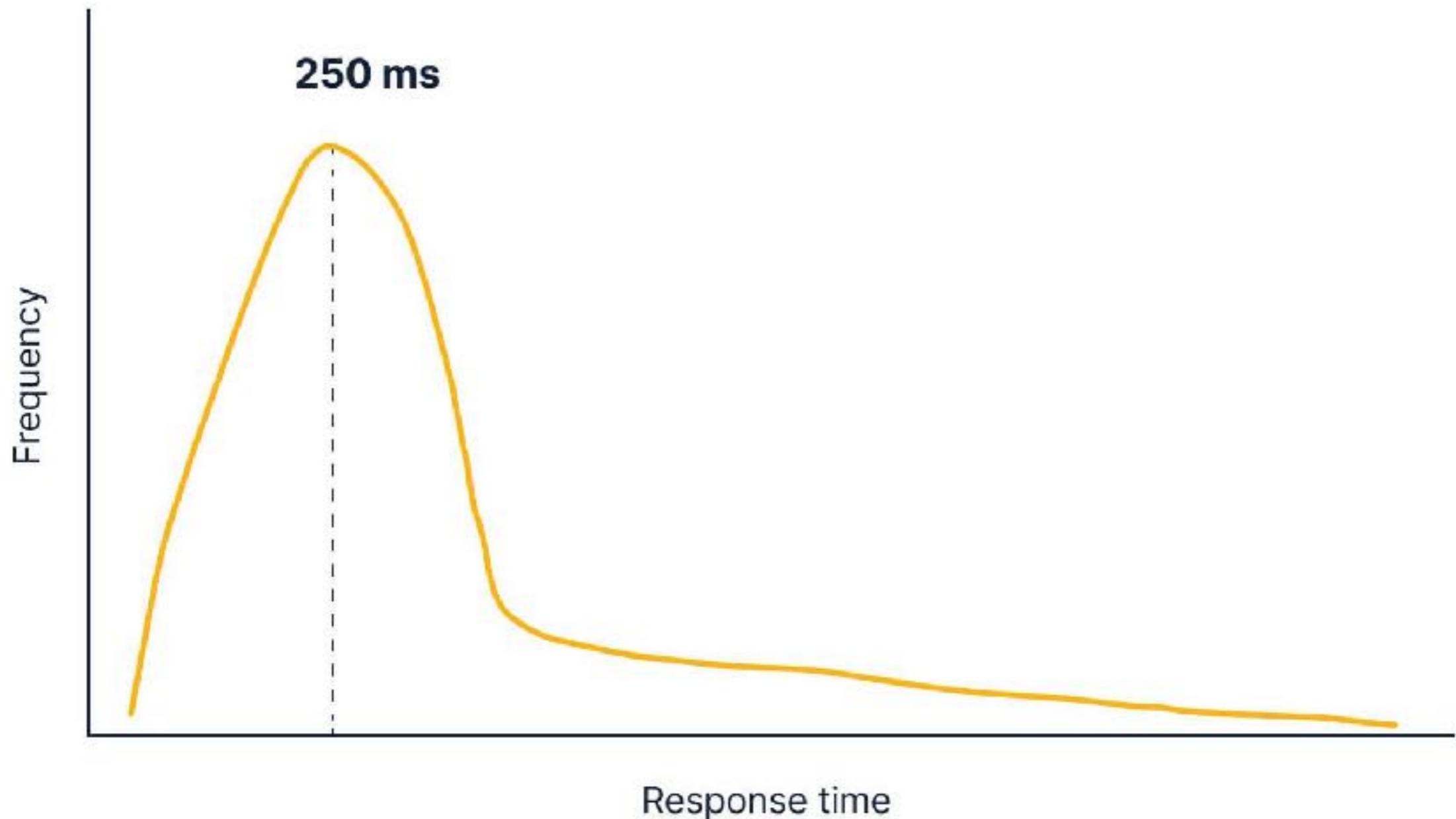
# 10th percentile



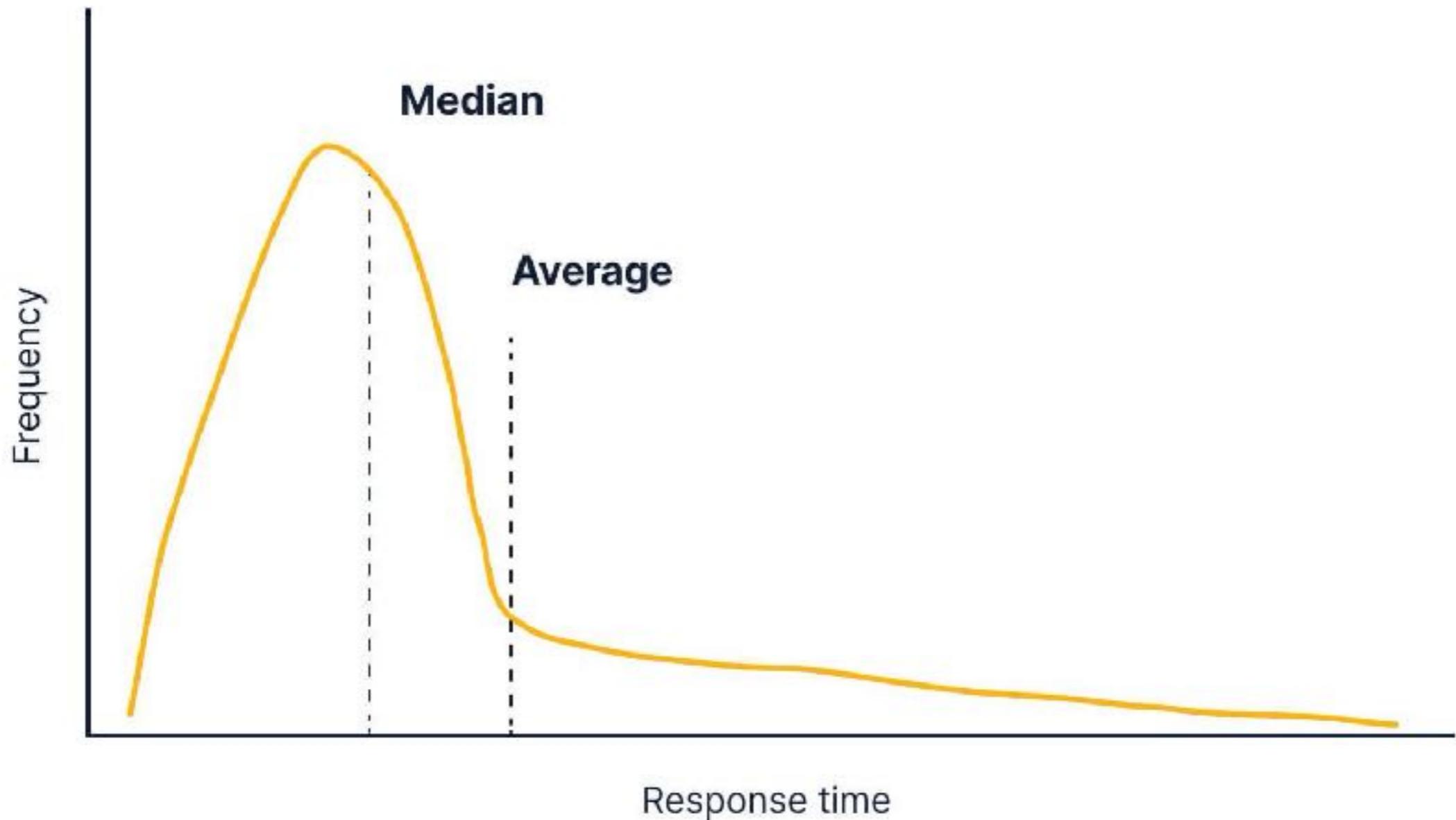
# 90th percentile



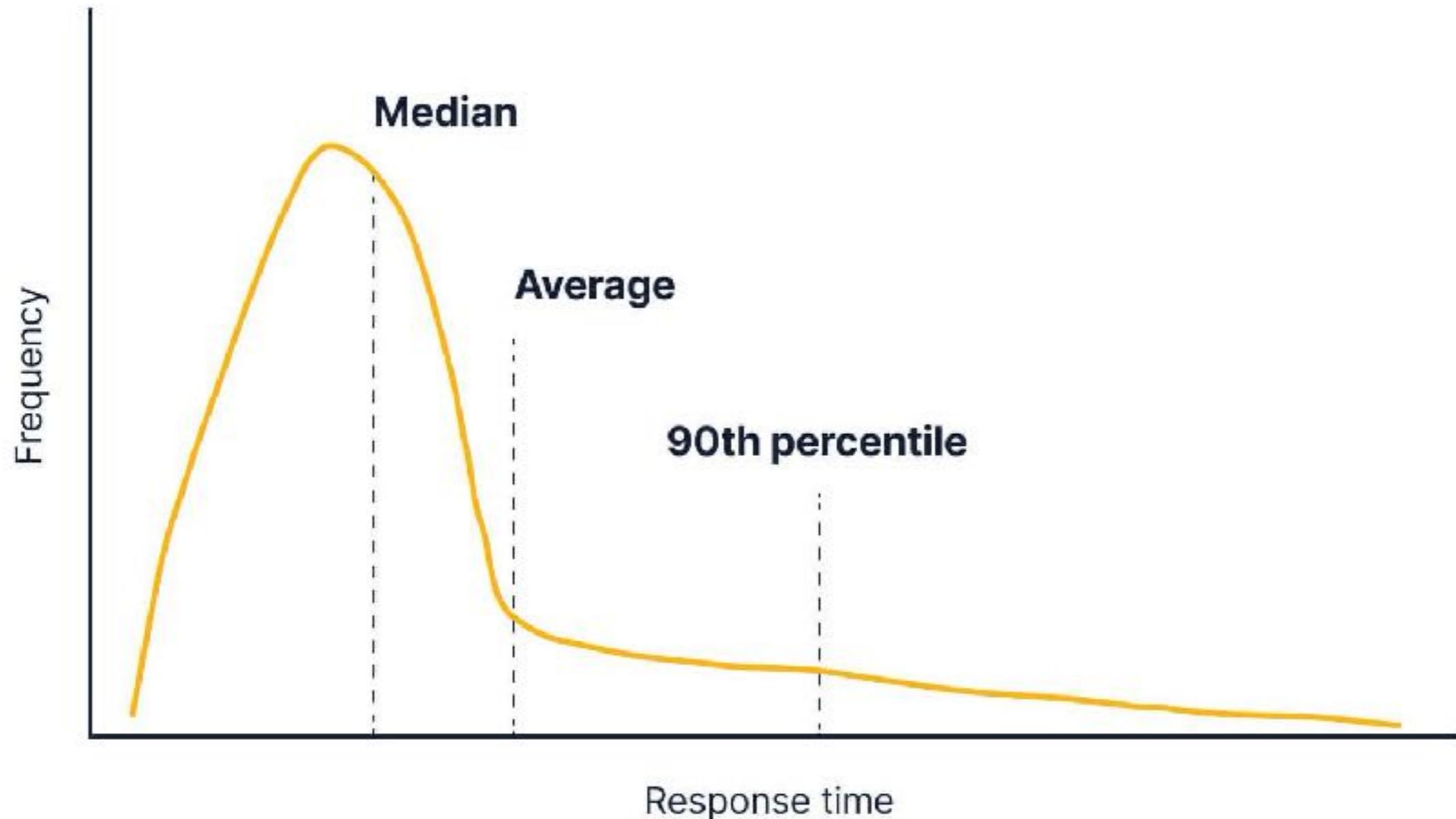
# Long tails (1)



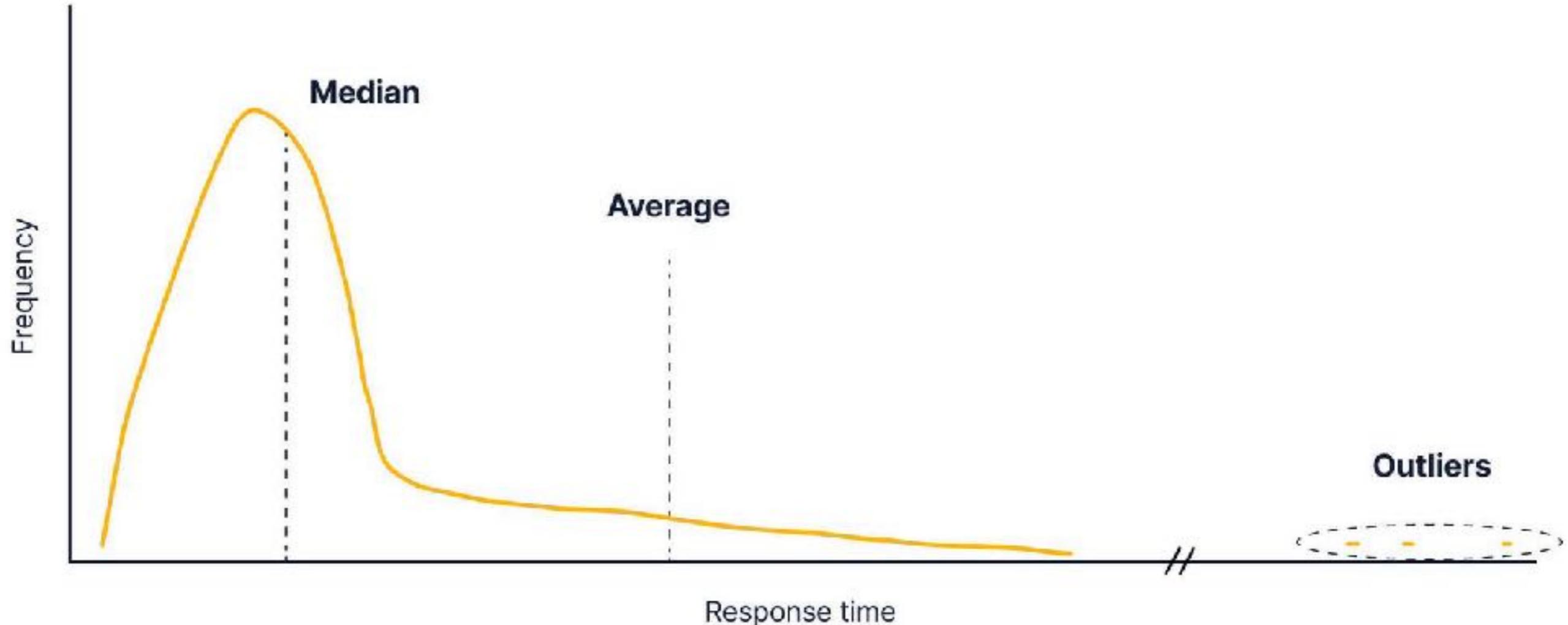
# Long tails (2)



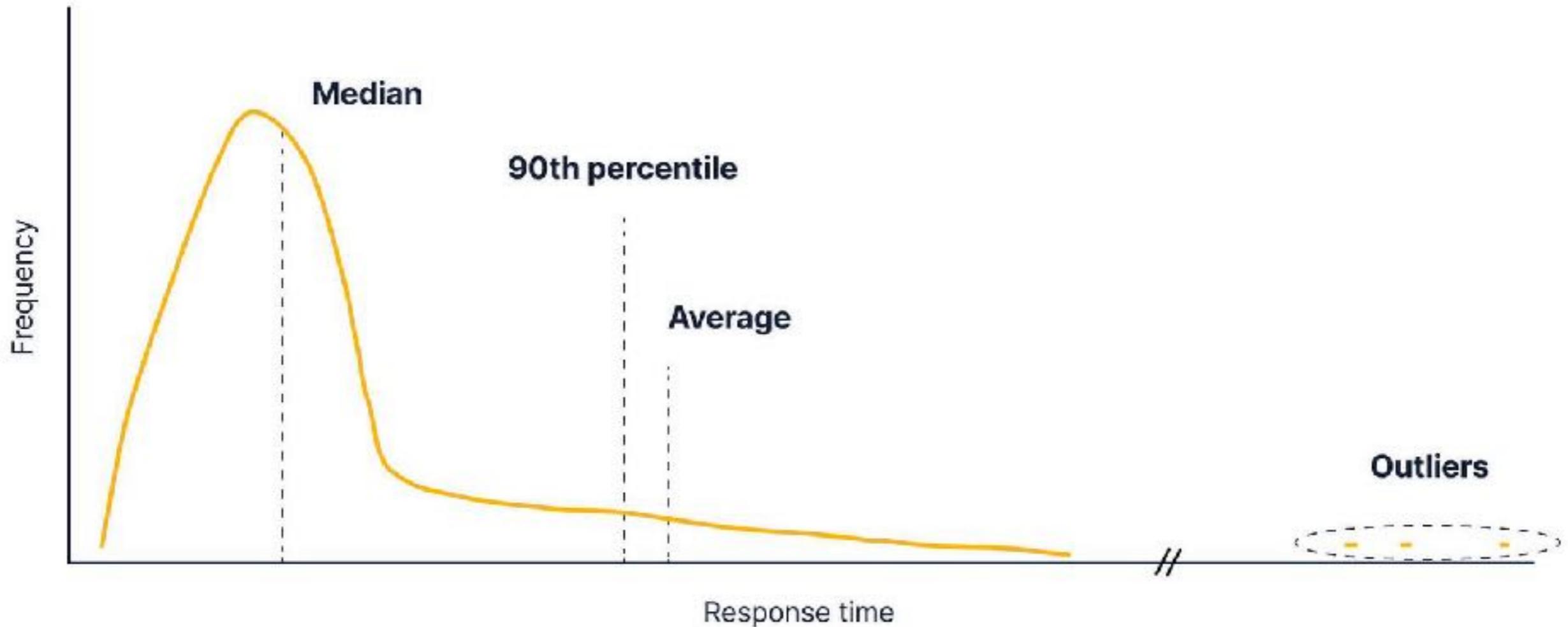
# Long tails (2)



# Long tails with outliers



# Long tails with outliers



# Workshop



# **Business Intelligence (BI) and Visualization**



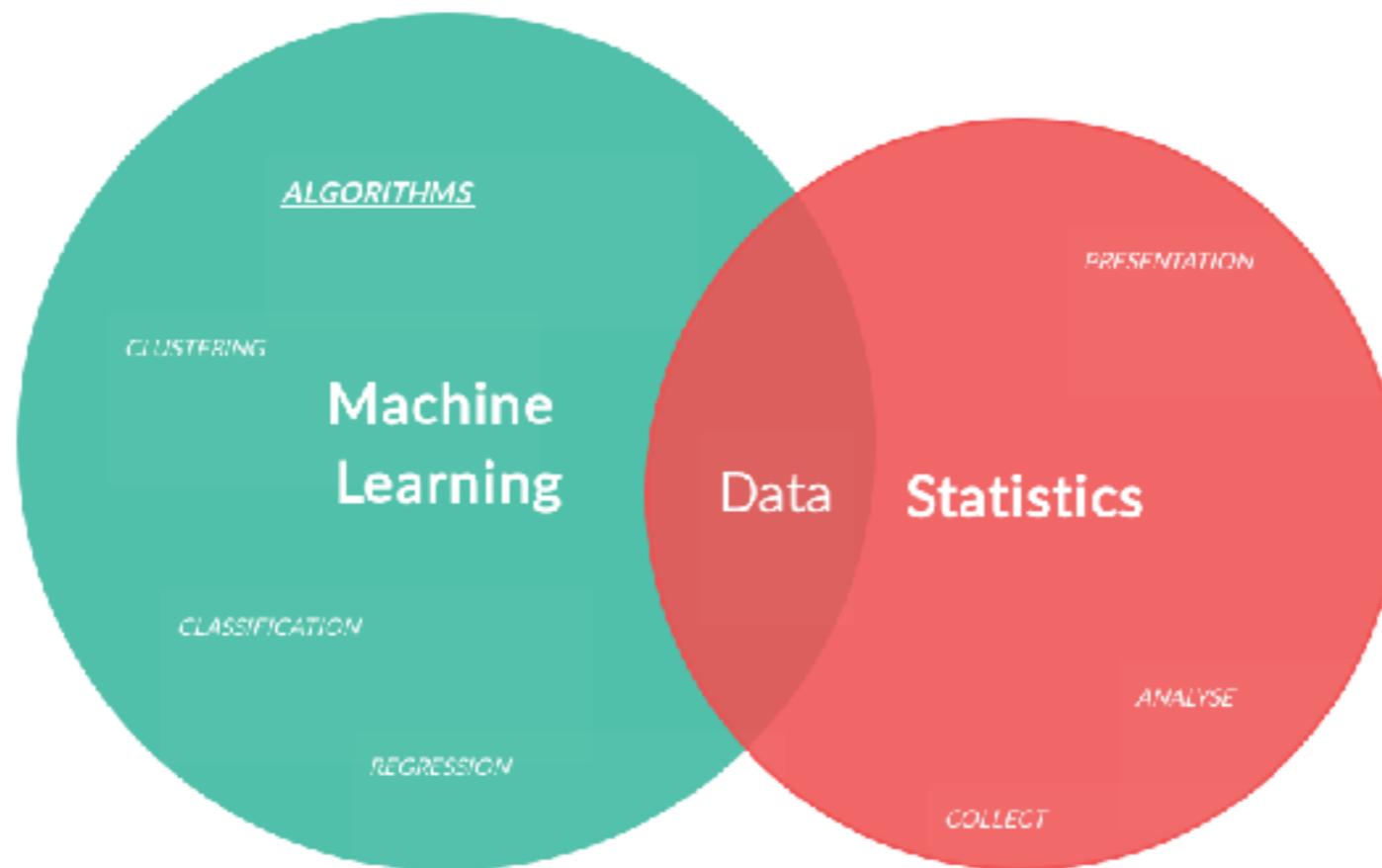
# Data visualization

Data visualization is a method that uses **visuals**,  
both static and interactive,  
to help people **understand**  
the large amount of data being collected.

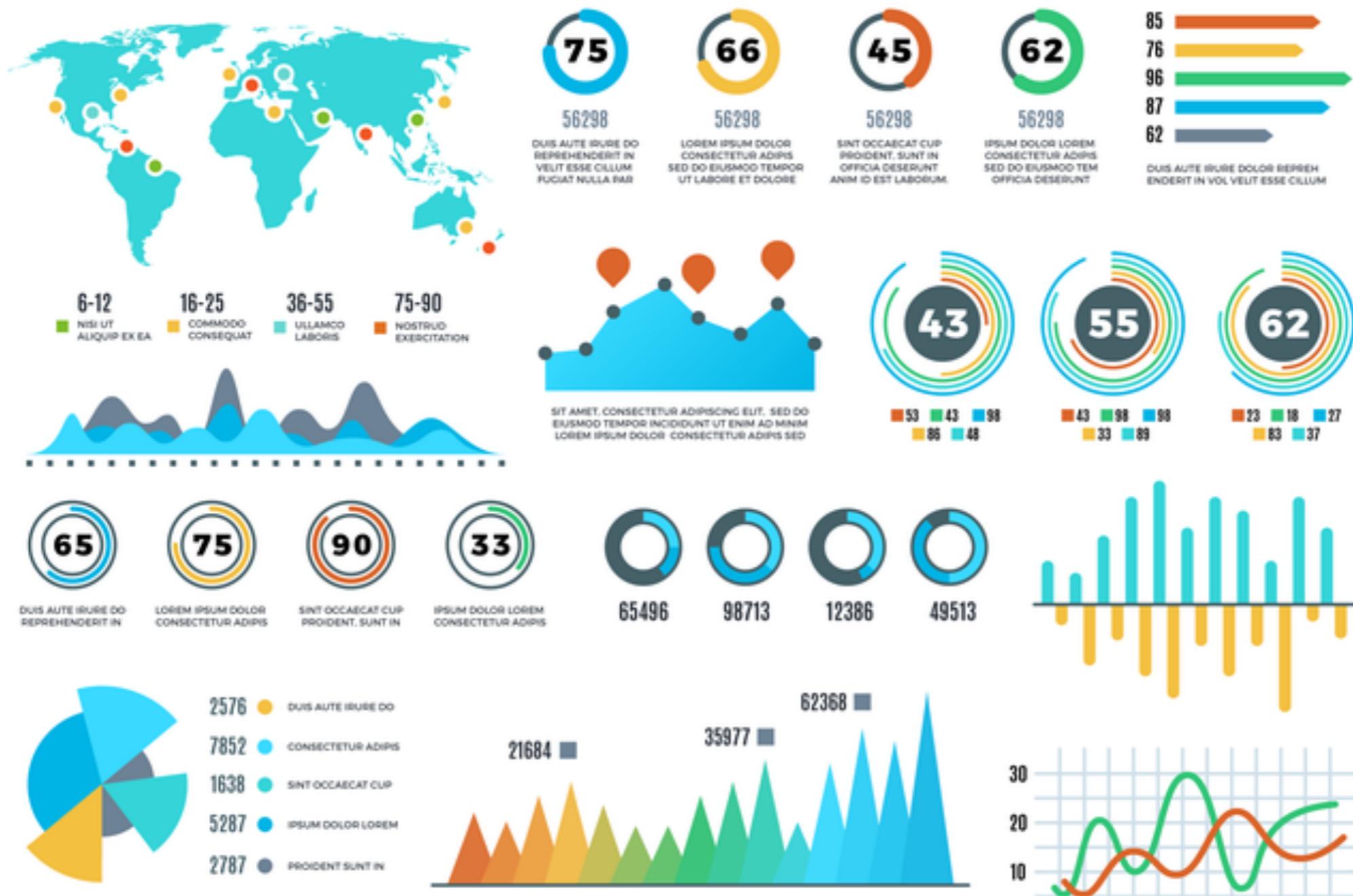


# Data visualization

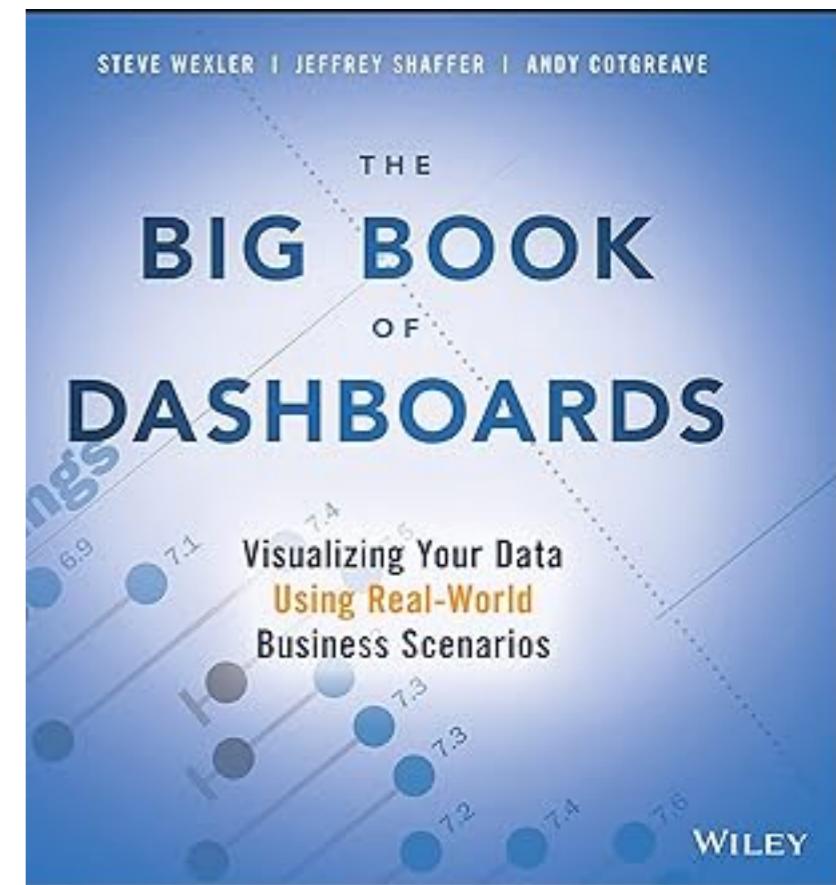
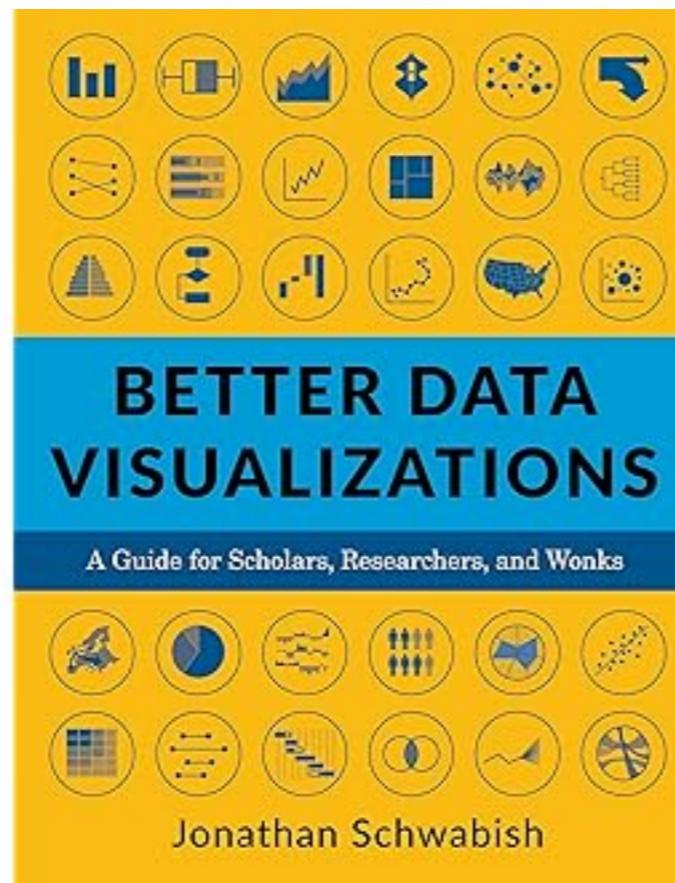
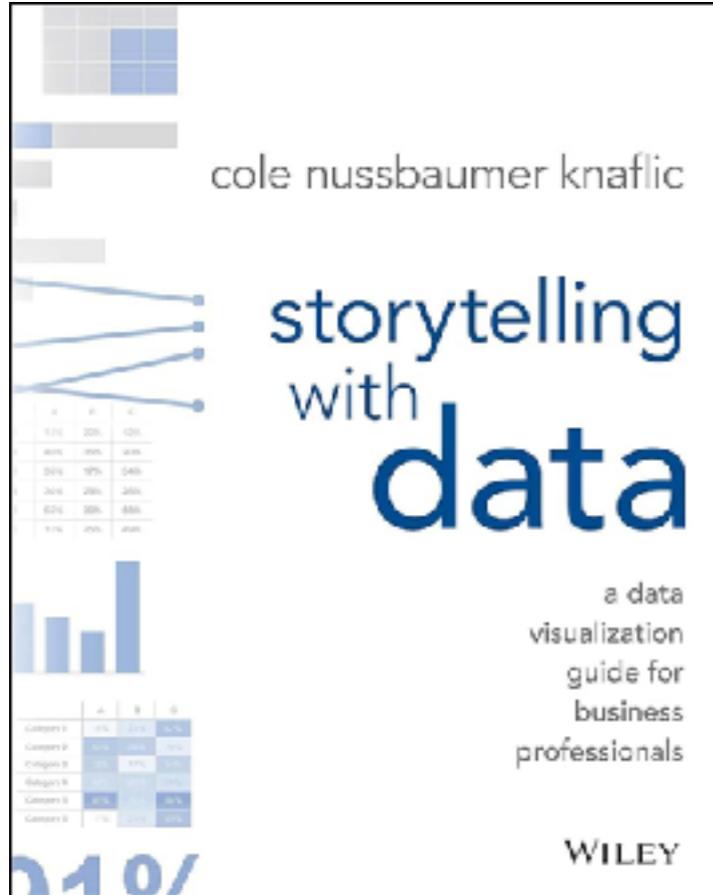
Data visualization is an **important skill** in applied statistics and machine learning.



# Data visualization



# Books



# Data visualization process

1

Collecting data

2

Clean your data

3

Choose a chart type

4

Prepare data

5

Visualize data

6

Presentation



# How to choose a chart type ?

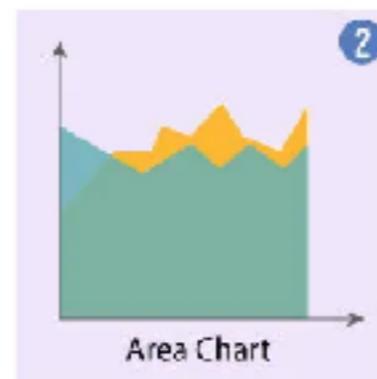
Before choosing a visual chart or graph,  
it is important to understand your **audience**



# TYPES OF DATA VISUALIZATION CHARTS



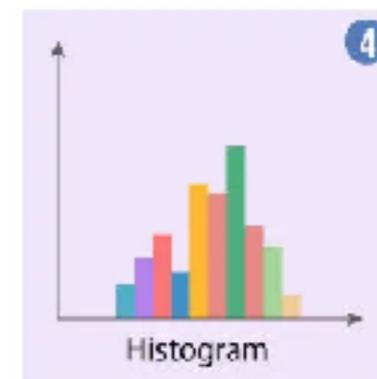
Display trends over time



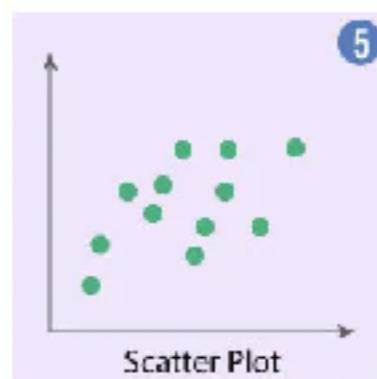
A line chart with areas below the lines filled with colors



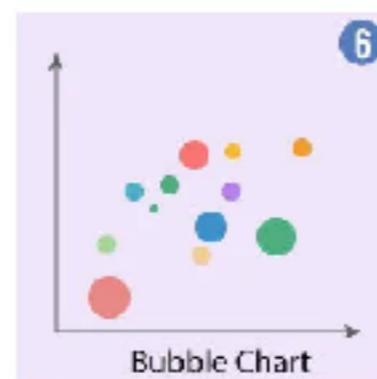
Display trends with multiple variables



Display the shape and spread of continuous dataset samples



Show correlation in a dataset



Show and compare the relationship between the labelled circles



Show the contribution of data point inside a whole dataset



Visualize the distance between intervals



Show data with location as a variable



Show magnitude of a phenomenon



# Key Concepts in Data Visualization

Data  
Representation

Clarity

Context

Accuracy

Color Theory

Storytelling

Interactivity



# Q/A



# Workshop



# Data is everywhere

Descriptive

Predictive

Prescriptive

Metrics  
Historical data

Insights  
Modeling

Data products

เกิดอะไรขึ้นในอดีต ?

จะเกิดอะไรขึ้นในอนาคต ?

จะทำให้สิ่งที่  
ผู้ใช้งานต้องการ  
เกิดขึ้นได้อย่างไร



# Understand your data

# of rows

# of columns

Column  
description

Alignment

Styling

Cleaning data

Outlining

Summarize

Visualize

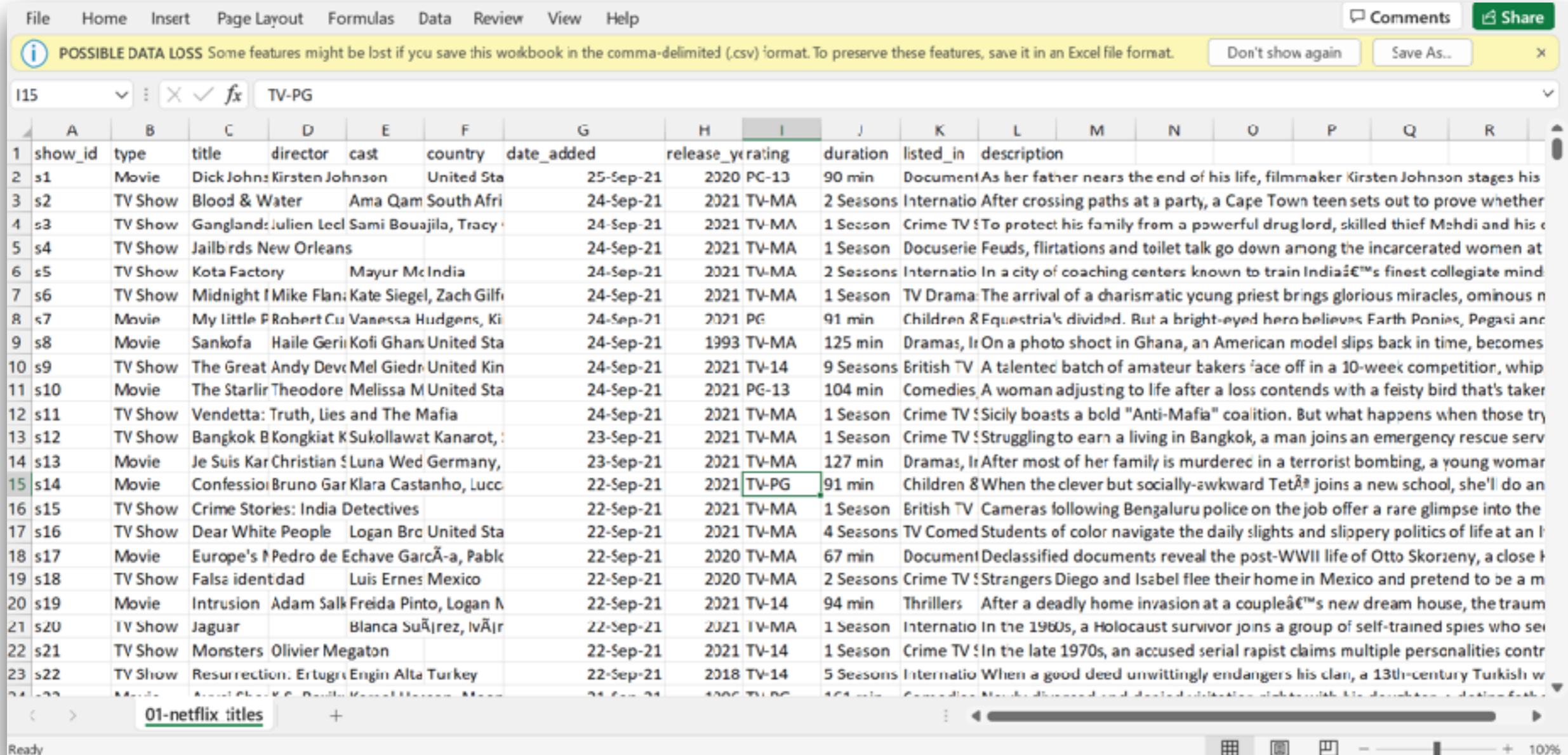


# **Workshop #1**

## **01-netflix\_titles**



# Open Data in Excel



show_id	type	title	director	cast	country	date_added	release_year	rating	duration	listed_in	description
s1	Movie	Dick Johnson	Kirsten Johnson		United States	25-Sep-21	2020	PG-13	90 min	Documentary	As her father nears the end of his life, filmmaker Kirsten Johnson stages his final film.
s2	TV Show	Blood & Water	Ama Qamir	South Africa		24-Sep-21	2021	TV-MA	2 Seasons	International	After crossing paths at a party, a Cape Town teen sets out to prove whether she can change her family's violent past.
s3	TV Show	Gangland: Julien Leclercq	Sami Bouajila, Tracy Kessinger		France	24-Sep-21	2021	TV-MA	1 Season	Crime TV	To protect his family from a powerful drug lord, skilled thief Mehdi and his crew must infiltrate a gangland prison.
s4	TV Show	Jailbirds	New Orleans		United States	24-Sep-21	2021	TV-MA	1 Season	Documentary	Feuds, flirtations and toilet talk go down among the incarcerated women at a maximum-security prison.
s5	TV Show	Kota Factory		Mayur McIndia	India	24-Sep-21	2021	TV-MA	2 Seasons	International	In a city of coaching centers known to train India's finest collegiate mind, a group of students prepare for their exams.
s6	TV Show	Midnight	I	Mike Flanagan, Kate Siegel, Zach Gilford	United States	24-Sep-21	2021	TV-MA	1 Season	TV Drama	The arrival of a charismatic young priest brings glorious miracles, ominous new powers and a secret past to a small town.
s7	Movie	My Little	P	Robert Cuccioli, Vanessa Hudgens, Kit Harington	United States	24-Sep-21	2021	PG	91 min	Children & Family	Equestria's divided. But a bright-eyed hero believes Earth Ponies, Pegasus and Unicorns can work together to save the day.
s8	Movie	Sankofa	Haile Gerima	Kofi Ghan	United States	24-Sep-21	1993	TV-MA	125 min	Dramas, International	On a photo shoot in Ghana, an American model slips back in time, becomes a child and must navigate a dangerous world.
s9	TV Show	The Great	Andy Devine, Mel Giedroyc		United Kingdom	24-Sep-21	2021	TV-14	9 Seasons	British TV	A talented batch of amateur bakers face off in a 10-week competition, whipping up some drama along the way.
s10	Movie	The Star	Levi	Theodore Tahu Rhodes, Melissa Mays	United States	24-Sep-21	2021	PG-13	104 min	Comedies	A woman adjusting to life after a loss contends with a feisty bird that's taken over her body.
s11	TV Show	Vendetta: Truth, Lies and The Mafia				24-Sep-21	2021	TV-MA	1 Season	Crime TV	Sicily boasts a bold "Anti-Mafia" coalition. But what happens when those trying to bring them down team up?
s12	TV Show	Bangkok B	Kongkiat K	Sukollawat Kanarot, Nopparat	Thailand	23-Sep-21	2021	TV-MA	1 Season	Crime TV	Struggling to earn a living in Bangkok, a man joins an emergency rescue service.
s13	Movie	Je Suis Kar	Christian S	Luna Wedde	Germany, France	23-Sep-21	2021	TV-MA	127 min	Dramas, International	After most of her family is murdered in a terrorist bombing, a young woman must confront her past and find a way to move forward.
s14	Movie	Confession	Bruno Ganz	Klara Castanho, Lucca	Portugal	22-Sep-21	2021	TV-PG	91 min	Children & Family	When the clever but socially-awkward Teté joins a new school, she'll do anything to fit in.
s15	TV Show	Crime Stories: India Detectives				22-Sep-21	2021	TV-MA	1 Season	British TV	Cameras following Bengaluru police on the job offer a rare glimpse into the daily lives of these officers.
s16	TV Show	Dear White People	Logan Browning	United States		22-Sep-21	2021	TV-MA	4 Seasons	TV Comedies	Students of color navigate the daily slights and slippery politics of life at an Ivy League school.
s17	Movie	Europe's	Pedro de Echave	García, Pablo	Spain	22-Sep-21	2020	TV-MA	67 min	Documentary	Declassified documents reveal the post-WWII life of Otto Skorzeny, a close friend of Adolf Hitler.
s18	TV Show	Falsa Identidad	Luis Ernesto	Mexico	Mexico	22-Sep-21	2020	TV-MA	2 Seasons	Crime TV	Strangers Diego and Isabel flee their home in Mexico and pretend to be a married couple.
s19	Movie	Intrusion	Adamalko	Freida Pinto, Logan Marshall-Grimes	United States	22-Sep-21	2021	TV-14	94 min	Thrillers	After a deadly home invasion at a couple's new dream house, the trauma continues.
s20	TV Show	Jaguar		Blanca Suárez, Iván Arbelaez	Spain	22-Sep-21	2021	TV-MA	1 Season	International	In the 1960s, a Holocaust survivor joins a group of self-trained spies who seek justice.
s21	TV Show	Monsters	Olivier Megaton			22-Sep-21	2021	TV-14	1 Season	Crime TV	In the late 1970s, an accused serial rapist claims multiple personalities controlled by different parts of his brain.
s22	TV Show	Resurrection: Ertugrul	Engin Altay	Alta	Turkey	22-Sep-21	2018	TV-14	5 Seasons	International	When a good deed unwittingly endangers his clan, a 13th-century Turkish warrior must make a difficult choice.

<https://www.kaggle.com/datasets/shivamb/netflix-shows>



Sharing

© 2020 - 2024 Siam Chamnkit Company Limited. All rights reserved.

# Step 1

Observe your data !!

# of rows

# of columns

Column  
description



# **Exploratory Data Analysis (EDA)**

Data  
transformation

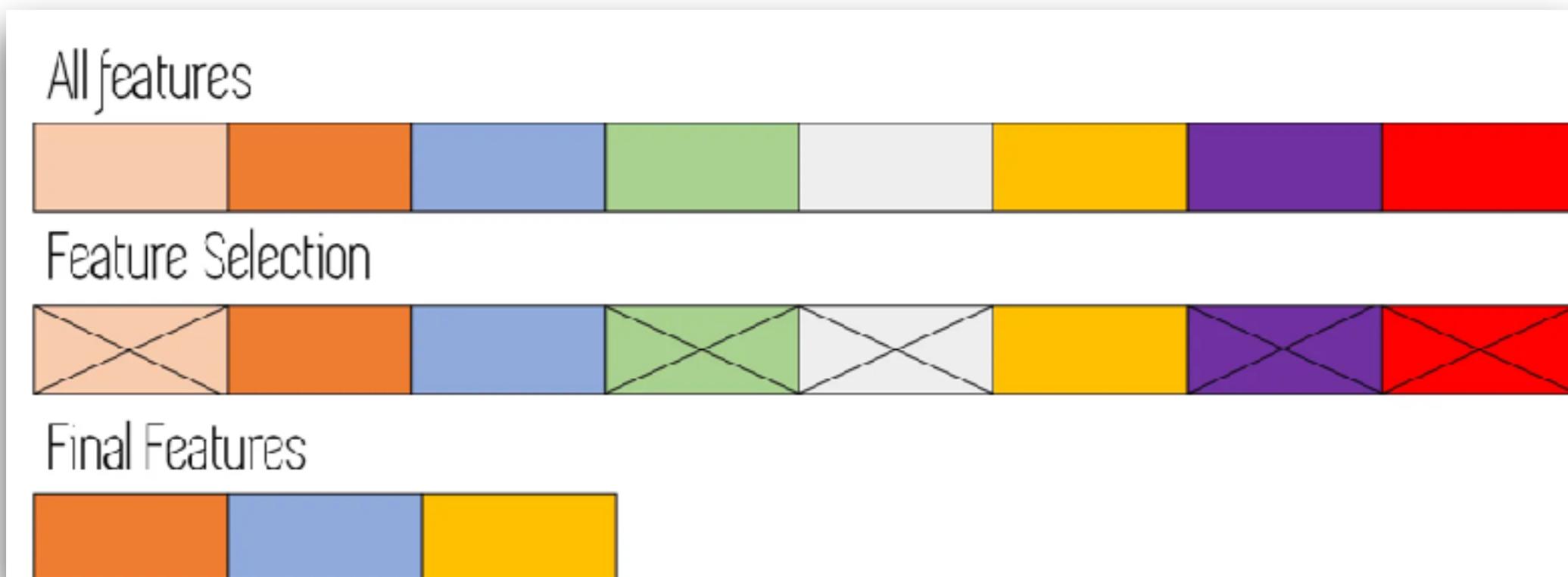
Data  
analysis

Data  
visualization



# Step 2 :: Feature selection

Consists on selecting the best features for our models and algorithms, by taking these insights from the data



# Feature selection ?

Title and description columns are low priority



# Step 3 :: Start with Question ?



# Sample questions ?

Top 10 director has cast more movie ?

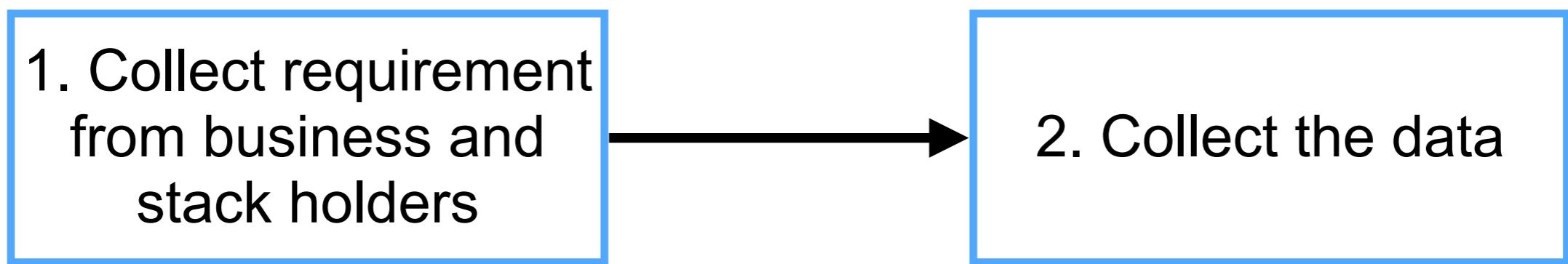
Top 5 country has produced more movies and TV shows ?

Which Genre has more in Movie Type ?

Count the movies & TV Shows before & after the year 2010 ?



# Steps in data analytic

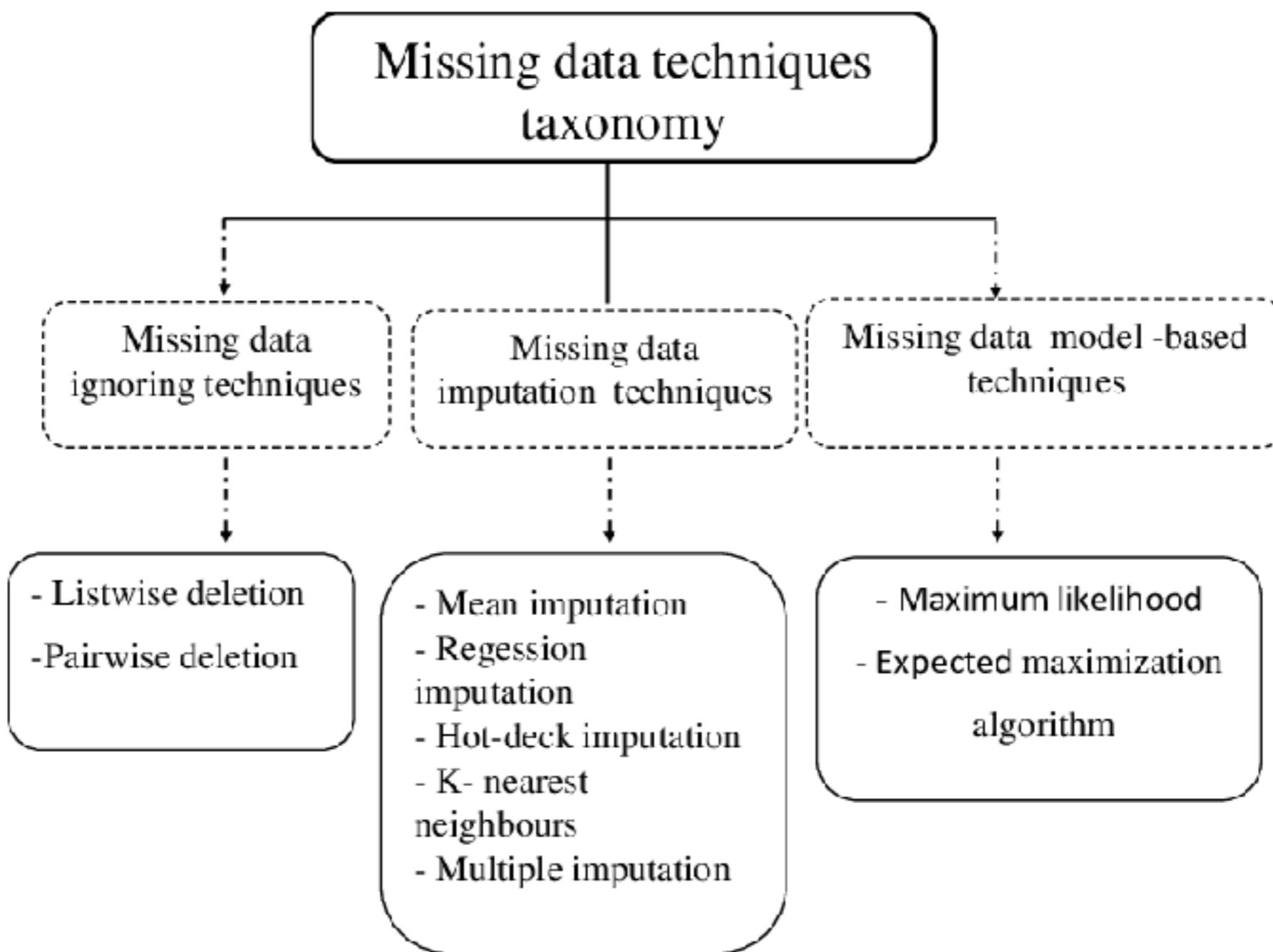


# Step 4 :: Clean your data

Empty/missing value  
Remove duplication and unwanted data  
Outliner data



# Missing data techniques



# Working with Excel #1

**Press CTRL + T**

To make it as a table for readability purpose



# Working with Excel #2

Look at each column  
and check it have **empty values**

COUNTBLANK()

COUNTA()



# Working with Excel #2

Check your data with = COUNTA(table[column])

1	Counting Value
2	show_id
3	type
4	title
5	director
6	cast
7	country
8	date_added
9	release_year
10	rating
11	duration
12	listed_in



# Working with Excel #3

Duplicated data in each columns

Listed\_in

Director

Cast

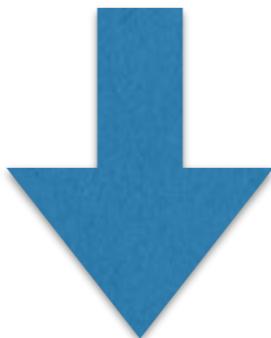
Country



# Working with Excel #3

Sample with listed\_in column

Crime TV Shows, International TV Shows, TV Action & Adventure



Crime TV  
Shows

International TV  
Shows

TV Action &  
Adventure



# Working with Excel #3

## Select Data -> Text to columns

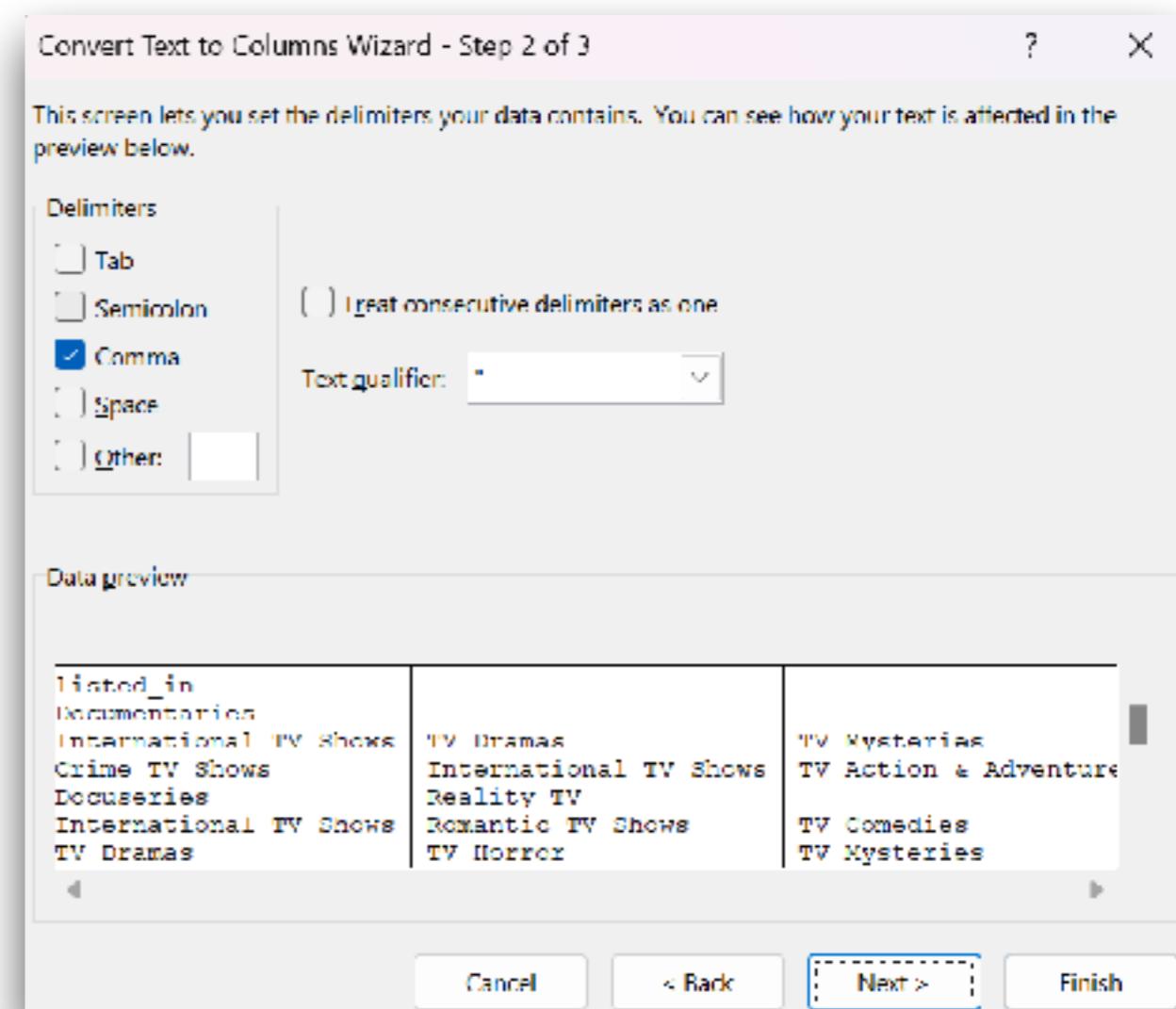
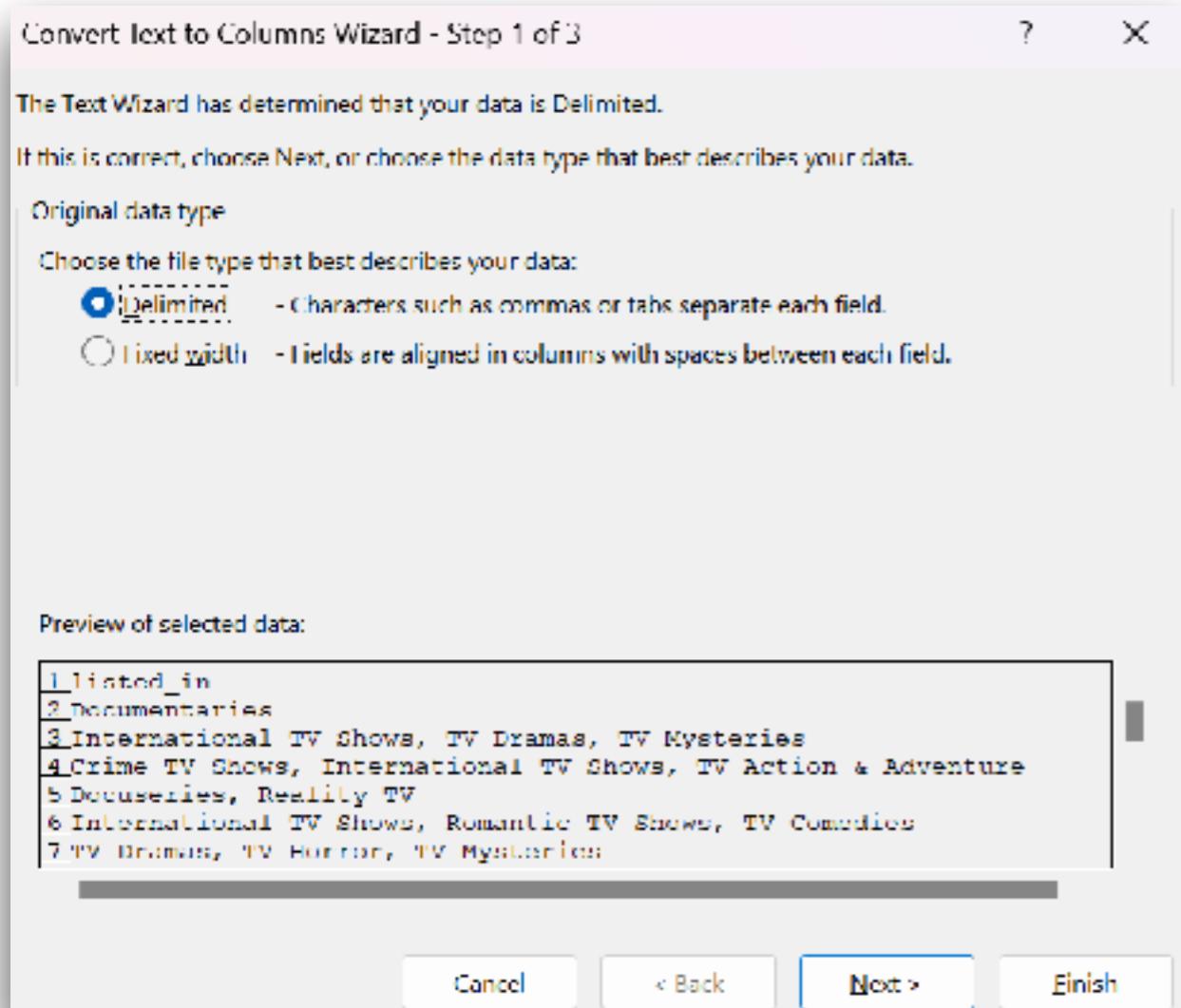
The screenshot shows a Microsoft Excel spreadsheet with the 'Data' tab selected. In the 'Data Tools' group, the 'Text to Columns' button is highlighted with a red dashed box. A tooltip for 'Text to Columns' is displayed, stating: 'Split a single column of text into multiple columns.' It also mentions: 'For example, you can separate a column of full names into separate first and last name columns.' Below that, it says: 'You can choose how to split it up: fixed width or split at each comma, period, or other character.' At the bottom right of the tooltip is a 'Tell me more' link.

3	I later	Ama Qam South Afric	24-Sep-21	2021 TV-MA	2 Seasons	International TV Shows, TV Crime TV Shows, International Docuseries, Reality TV
4	Julien Lelio Sami Bouajila, Tracy C		24 Sep 21	2021 TV-MA	1 Season	TV Dramas, TV Horror, TV Children & Family Movies
5	ew Orleans		24-Sep-21	2021 TV-MA	1 Season	Dramas, Independent Movie British TV Shows, Reality TV Comedies, Dramas
6	TV	Mayur Mc India	24-Sep-21	2021 TV-MA	2 Seasons	Crime TV Shows, Docuseries, International TV Dramas, International Movie British TV Shows, Reality TV Crime TV Shows, Docuseries, International TV Dramas, International Movie British TV Shows, Reality TV
7	Mike Flanagan, Kate Siegel, Zach Gilfoyle		24 Sep 21	2021 TV-MA	1 Season	Children & Family Movies, Thrillers
8	Robert Cullinan Vanessa Hudgens, Kira		24-Sep-21	2021 PG	91 min	International TV Shows, TV Action & Adventure
9	Haile Gerin Kofi Ghan	United States	24-Sep-21	1993 TV-MA	125 min	International TV Shows, TV Action & Adventure, TV Dramas
10	Andy Devine Mel Giedroyc United King		24-Sep-21	2021 TV-14	9 Seasons	International TV Shows, TV Action & Adventure, TV Dramas
11	Theodore T Melissa M United Stat		24-Sep-21	2021 PG-13	104 min	International TV Shows, TV Action & Adventure, TV Dramas
12	Truth, Lies and The Mafia		24-Sep-21	2021 TV-MA	1 Season	International TV Shows, TV Action & Adventure, TV Dramas
13	Kongkiat Kongkollawat Kanarot, S		23-Sep-21	2021 TV-MA	1 Season	International TV Shows, TV Action & Adventure, TV Dramas
14	Christian S Luna Wed Germany, C		23-Sep-21	2021 TV-MA	127 min	International TV Shows, TV Action & Adventure, TV Dramas
15	Bruno Garcia Klara Castanho, Lucca		22-Sep-21	2021 TV-PG	91 min	International TV Shows, TV Action & Adventure, TV Dramas
16	ies: India Detectives		22-Sep-21	2021 TV-MA	1 Season	International TV Shows, TV Action & Adventure, TV Dramas
17	the People Logan Browning United Stat		22-Sep-21	2021 TV-MA	4 Seasons	International TV Shows, TV Action & Adventure, TV Dramas
18	Pedro de Echave Garcia, Pablo		22-Sep-21	2020 TV-MA	67 min	International TV Shows, TV Action & Adventure, TV Dramas
19	tidad Luis Ernesto Mexico		22-Sep-21	2020 TV-MA	2 Seasons	International TV Shows, TV Action & Adventure, TV Dramas
20	Adam Salky Freida Pinto, Logan M		22-Sep-21	2021 TV-14	94 min	International TV Shows, TV Action & Adventure, TV Dramas
21	Blanca Suárez, Iván		22-Sep-21	2021 TV-MA	1 Season	International TV Shows, TV Action & Adventure, TV Dramas
22	Olivier Megaton		22-Sep-21	2021 TV-14	1 Season	International TV Shows, TV Action & Adventure, TV Dramas
23	on: Ertugrul Engin Alta Turkey		22-Sep-21	2018 TV-14	5 Seasons	International TV Shows, TV Action & Adventure, TV Dramas



# Working with Excel #3

## Choose delimiter with comma (,)



# Working with Excel #3

## Result in Excel file

listed_in	Column1	Column2
Documentaries		
International TV Shows	TV Dramas	TV Mysteries
Crime TV Shows	International TV Shows	TV Action & Adventure
Docuseries	Reality TV	
International TV Shows	Romantic TV Shows	TV Comedies
TV Dramas	TV Horror	TV Mysteries
Children & Family Movies		
Dramas	Independent Movies	International Movies
British TV Shows	Reality TV	
Comedies	Dramas	
Crime TV Shows	Docuseries	International TV Shows
Crime TV Shows	International TV Shows	TV Action & Adventure
Dramas	International Movies	
Children & Family Movies	Comedies	
British TV Shows	Crime TV Shows	Docuseries
TV Comedies	TV Dramas	
Documentaries	International Movies	
Crime TV Shows	Spanish-Language TV Shows	TV Dramas
Thrillers		
International TV Shows	Spanish-Language TV Shows	TV Action & Adventure
Crime TV Shows	Docuseries	International TV Shows



# Try by yourself

Listed\_in

Director

Cast

Country



# Step 5 :: Start analyze the data

Top 10 director has cast more movie ?

Top 5 country has produced more movies and TV shows ?

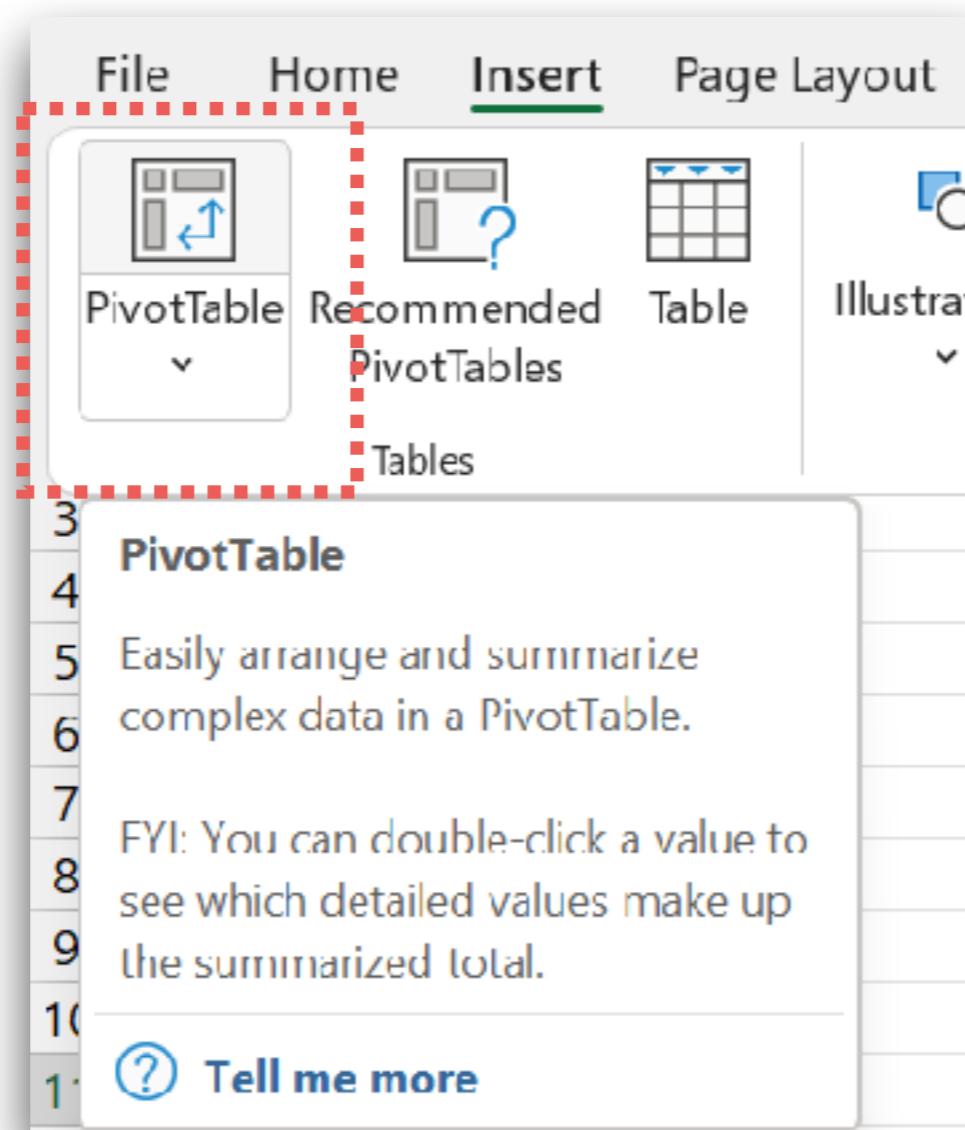
Which Genre has more in Movie Type ?

Count the movies & TV Shows before & after the year 2010 ?



# Top 10 director has cast more movie ?

## Working with Pivot Table



# Top 10 director has cast more movie ?

Choose Columns, Rows and Values

The screenshot shows the 'PivotTable Fields' dialog box with the following settings:

- Choose fields to add to report:** A search bar at the top.
- Fields list:**
  - type**
  - title**
  - director**
  - cast**
  - country**
- Drag fields between areas below:** A section for dragging fields between columns, rows, and values.
- Filters:** A section for applying filters.
- Columns:** A list under 'Columns' containing 'type'.
- Values:** A list under 'Values' containing 'Count of director'.
- Rows:** A list under 'Rows' containing 'director'.
- Defer Layout Update:** A checkbox at the bottom.

**Annotations:**

- 1. Column = type** (Red text)
- 2. Rows = director** (Red text)
- 3. Values = director** (Red text)



# Top 10 director has cast more movie ?

Count of director Row Labels	Column Labels		
	Movie	TV Show	Grand Total
A. L. Vijay		2	2
A. Raajdheep		1	1
A. Salaam		1	1
A.R. Murugadoss		2	2
Ã“skar ThÃ³r Axelsson		1	1
Ã€lex Pastor, David Pastor		2	2
Ã‡agan Irmak		1	1
Ãlex de la Iglesia		2	2
Ãlvaro Brechner		1	1
Ãlvaro Delgado-Aparicio L.		1	1
Ãlvaro Longoria, Gerardo Olivares		1	1
Ãngel GÃ³mez HernÃ¡ndez		1	1
Ãngeles ReinÃ©		1	1
Ãsold UggadÃ³ttir		1	1
Aadish Keluskar		1	1
A. S. D. L.		1	1



# Try by yourself

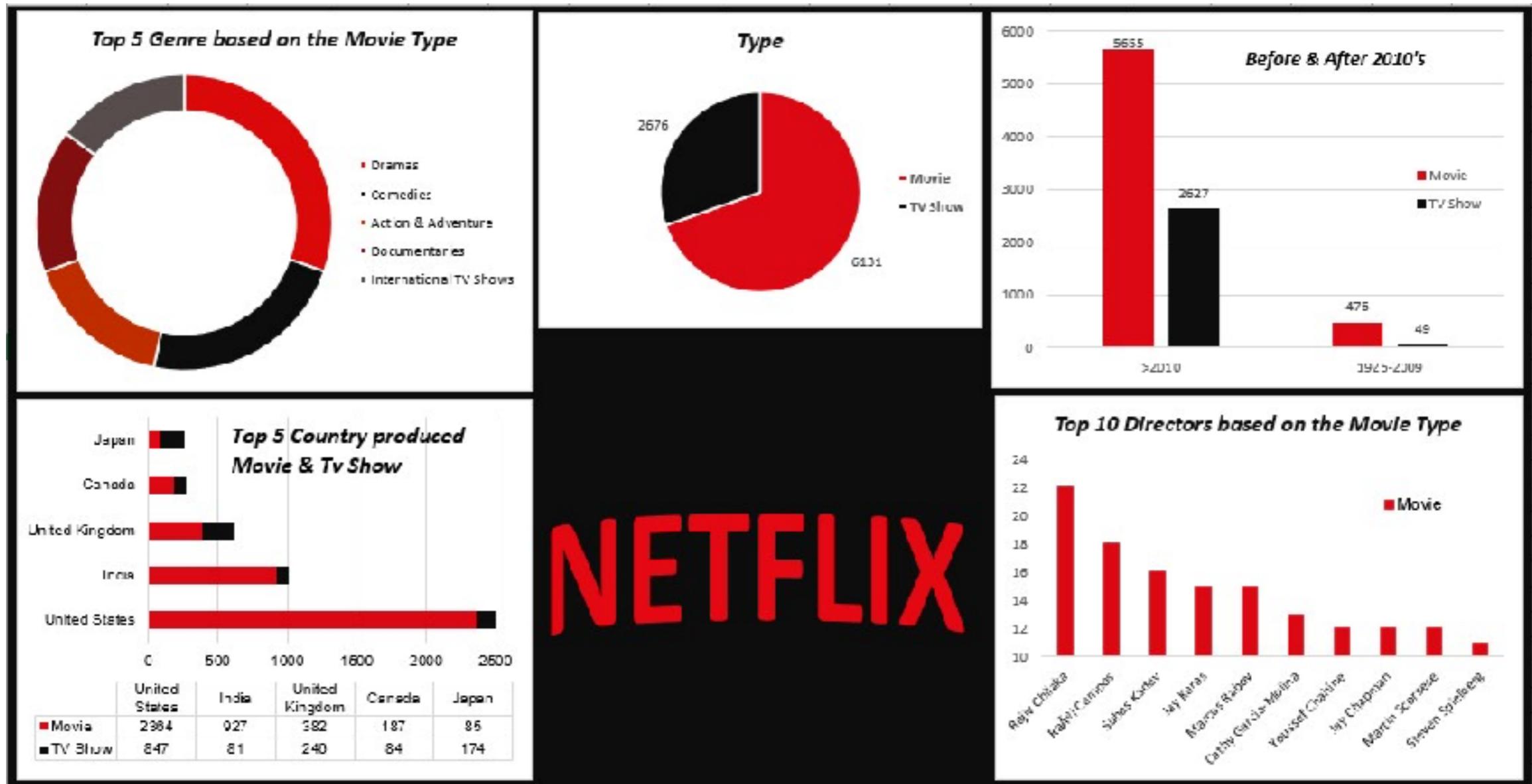
Top 5 country has produced more movies and TV shows ?

Which Genre has more in Movie Type ?

Count the movies & TV Shows before & after the year 2010 ?



# Step 6 :: Data Visualization



# **Step 7 ::**

# **Communicate**

# **the insight/findings to business**

