

# Data Mining Beginner to Intermediate

Parameswari Ettiappan



High performance. Delivered.

consulting | technology | outsourcing

# Goals

---

- What is Machine Learning
- Use Cases
- Commonly used Terms
- Lifecycle of a ML Project
- Supervised Learning
- Unsupervised Learning

# Goals

---

- Data Acquisition
- Types of Data
- Data Types
- Exploratory Data Analysis
- Data Pre-processing
- Data Quality assessment
- Feature Scaling

# Goals

---

- Descriptive Statistics
- Methods to impute missing values
- Outlier/Anomaly Detection
- Data Visualization
- Histogram
- Bar Graph
- Scatter Plot
- Pie Chart

# Goals

---

- Box Plot
- Feature Selection
- Univariate Selection
- Feature Importance
- Correlation matrix and Heat map
- Underfitting vs Overfitting
- Bias-Variance Trade-off

# Goals

---

- Hypothesis Testing
- Statistical Assumptions
- Null Hypothesis
- Alternate Hypothesis
- One sample Z-test
- Z-test in Python
- T-test
- T-test in Python

# Goals

---

- Pearson's Chi Squared Test
- Confusion Matrix
- Absolute Error
- Relative Error
- RMSE
- Precision, Accuracy
- Recall
- Specificity
- F-Score
- ROC/AUC

# Goals

---

- Cost Function
- Gradient Descent
- What is Regression
- Basic Idea
- Linear Regression Applications
- Linear Regression

# Goals

---

- Types of Errors
- Better Regression Models
- Correlation is not Causation
- Polynomial Linear Regression
- Regularization
- Ridge Regression
- LASSO Regression

# Goals

---

- Types of Classification Algorithms
- Applications of Classification Algorithms
- Logistic Function
- Logistic Regression
- Application of Logistic Regression
- Types of Logistic Regression
- Decision Trees
- Working of Decision Tree

# Goals

---

- Attribute Selection measure
- Gini Index
- Information Gain
- Random Forests
- Working of RF
- Advantages and Disadvantages of RF algorithm
- Application of RF
- XGBoost
-

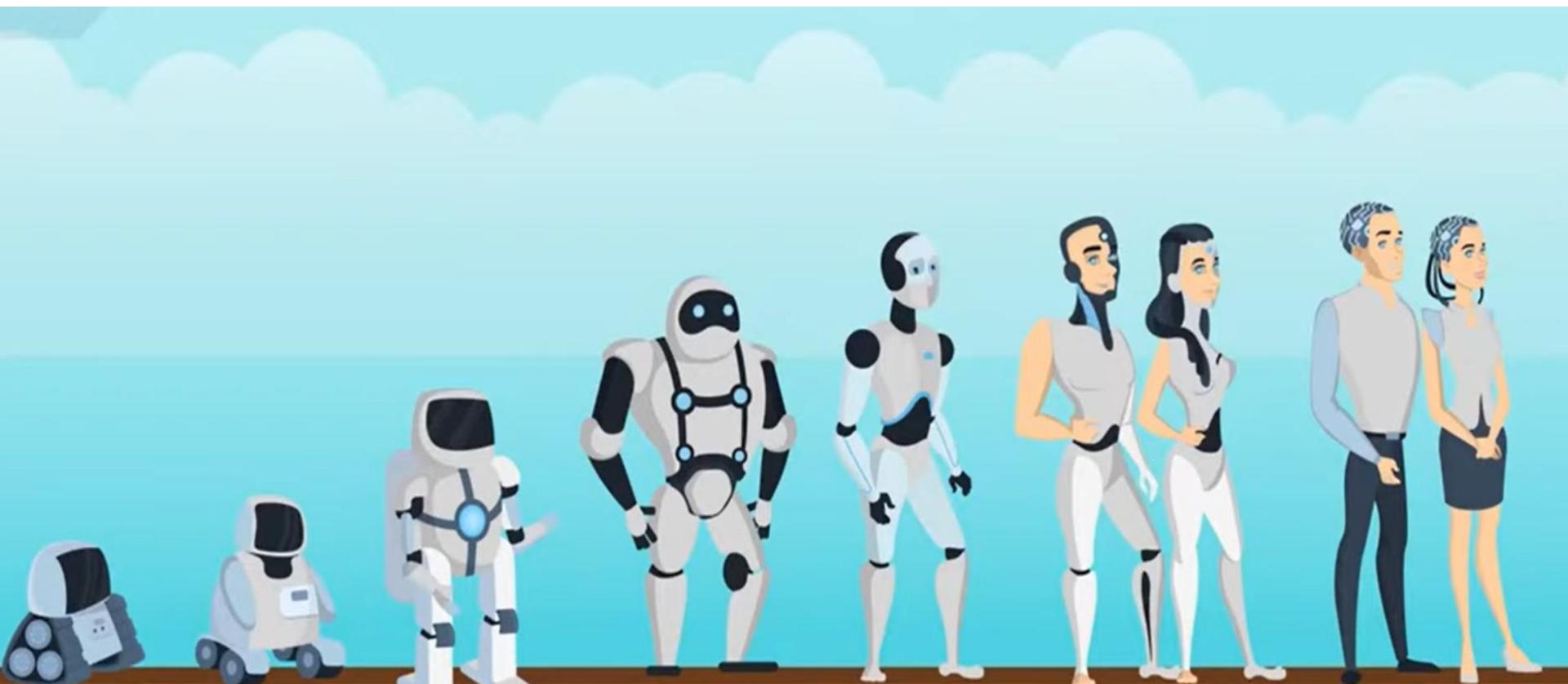
# Goals

---

- Why Unsupervised Learning
- Applications
- Clustering
- Types of Clustering
- Singular value Decomposition
- Independent Component Analysis
- Association Rules

# Machine Learning

## Evolution of Machines





DATA SCIENCE

# Machine Learning

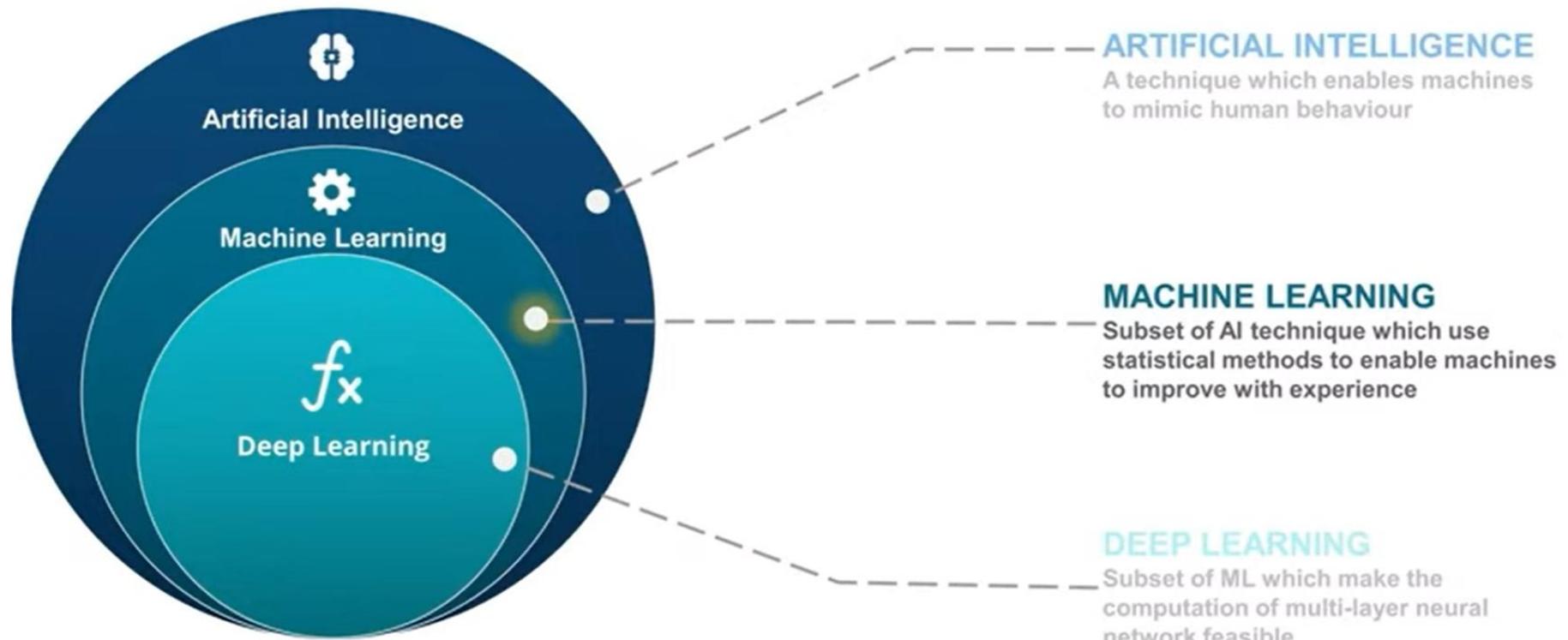
## What is Machine Learning?

Machine Learning is a subset of artificial intelligence. It focuses mainly on the designing of systems, thereby allowing them to learn and make predictions based on some experience which is data in case of machines.

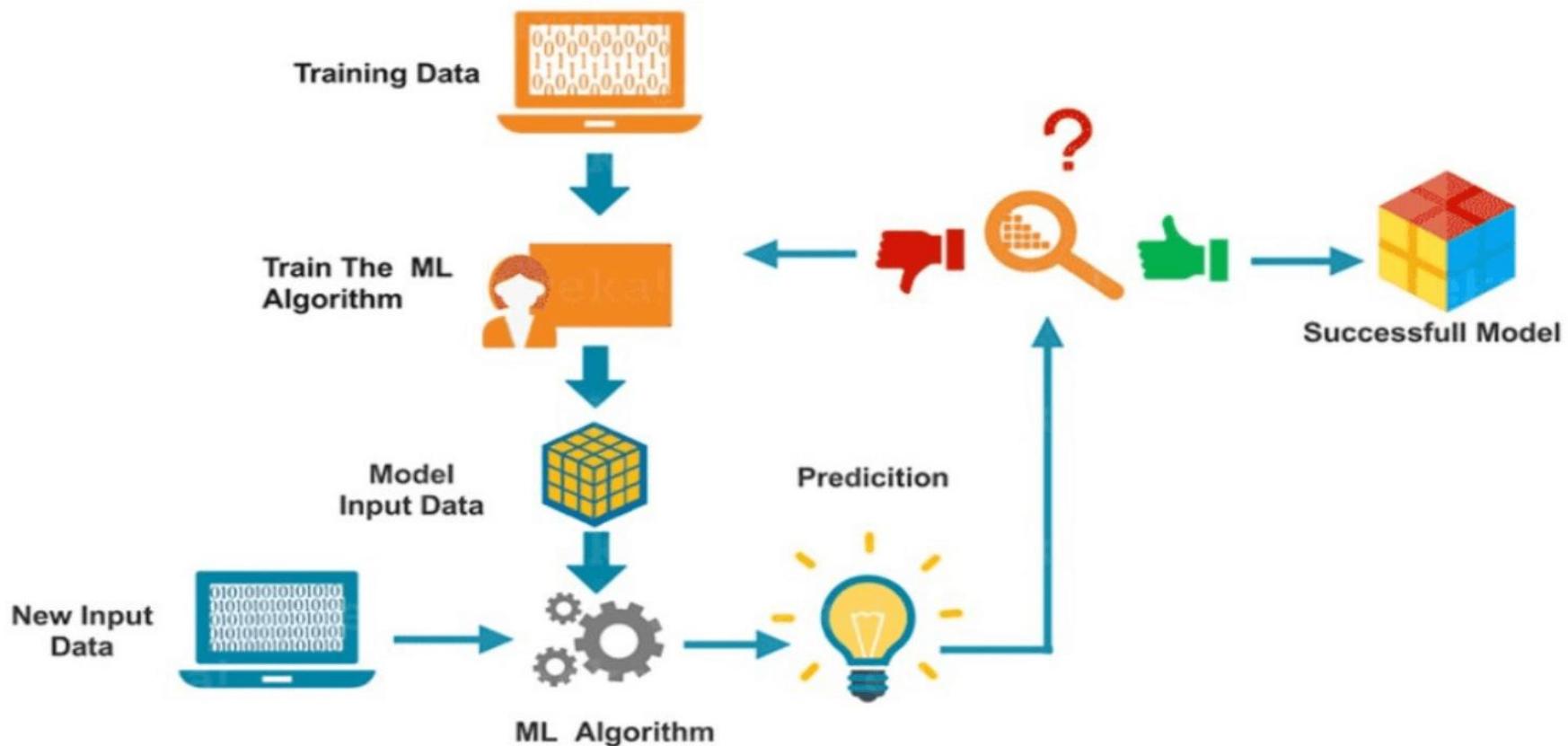
Learning?



# Big Confusion AI, ML and Deep Learning



# Machine Learning

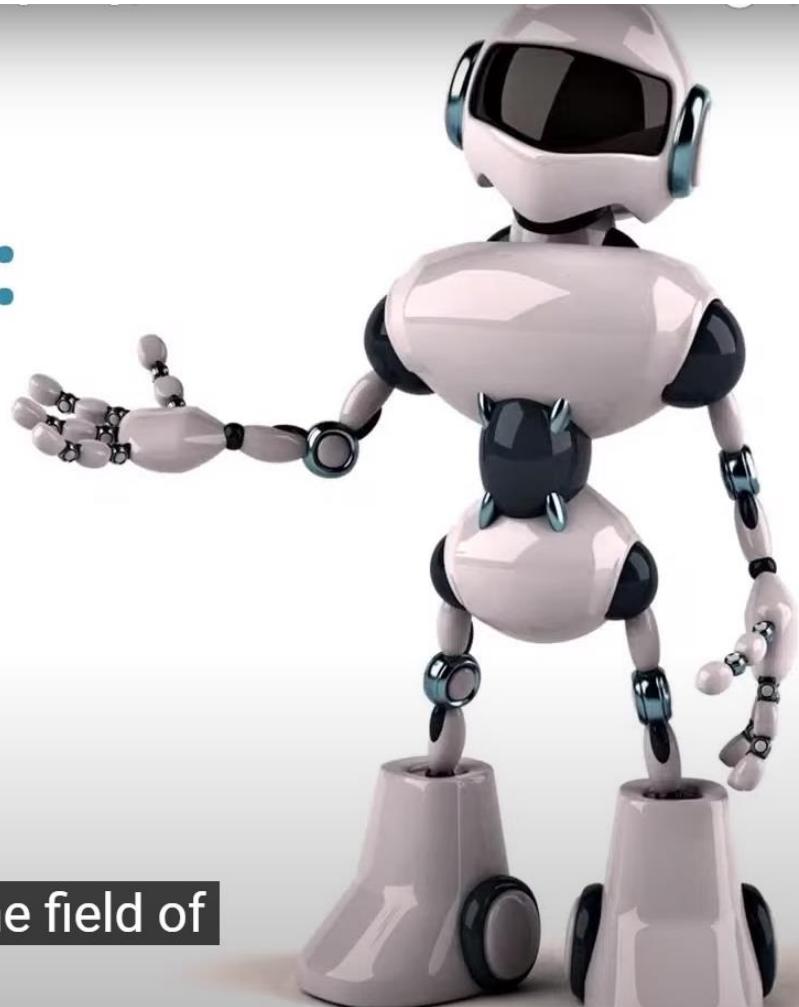


# Types of Machine Learning

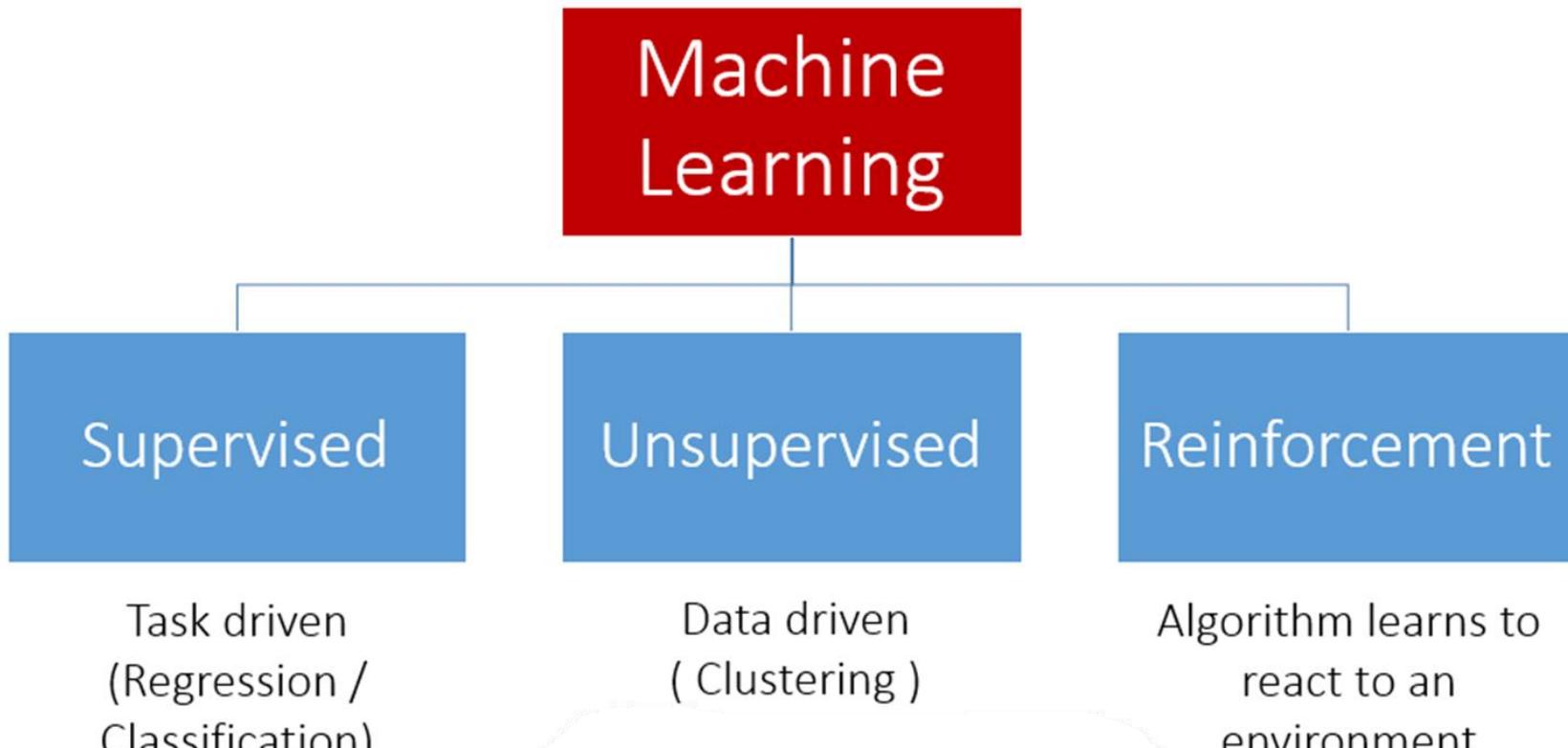
## Machine Learning Types:

- **Supervised**
- **Unsupervised**
- **Reinforcement**

how each of them is used in the field of  
banking

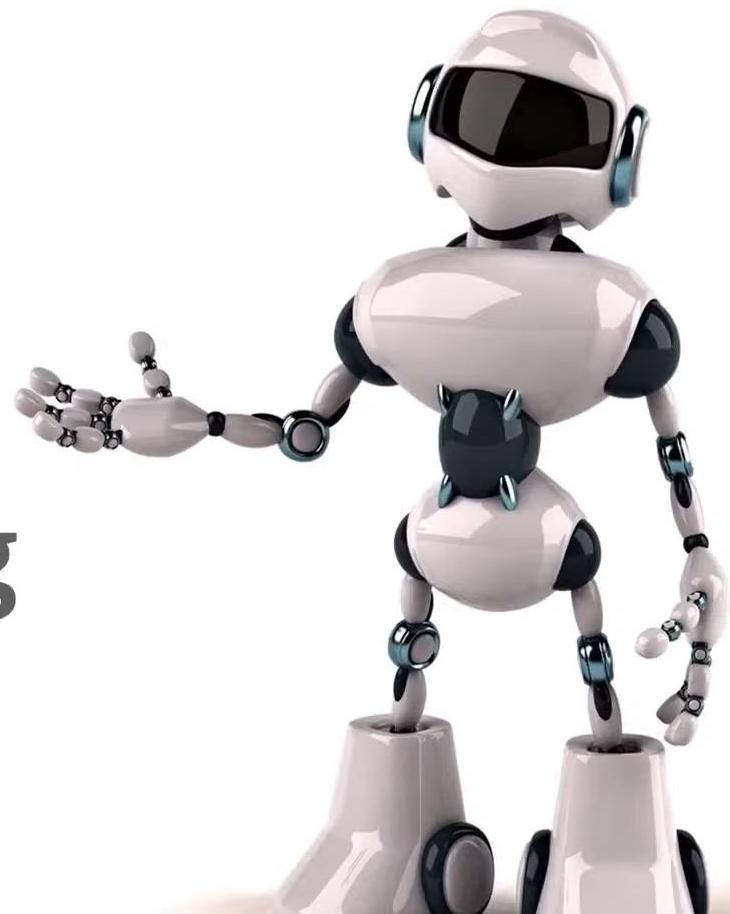


# Types of Machine Learning





# Supervised Machine Learning

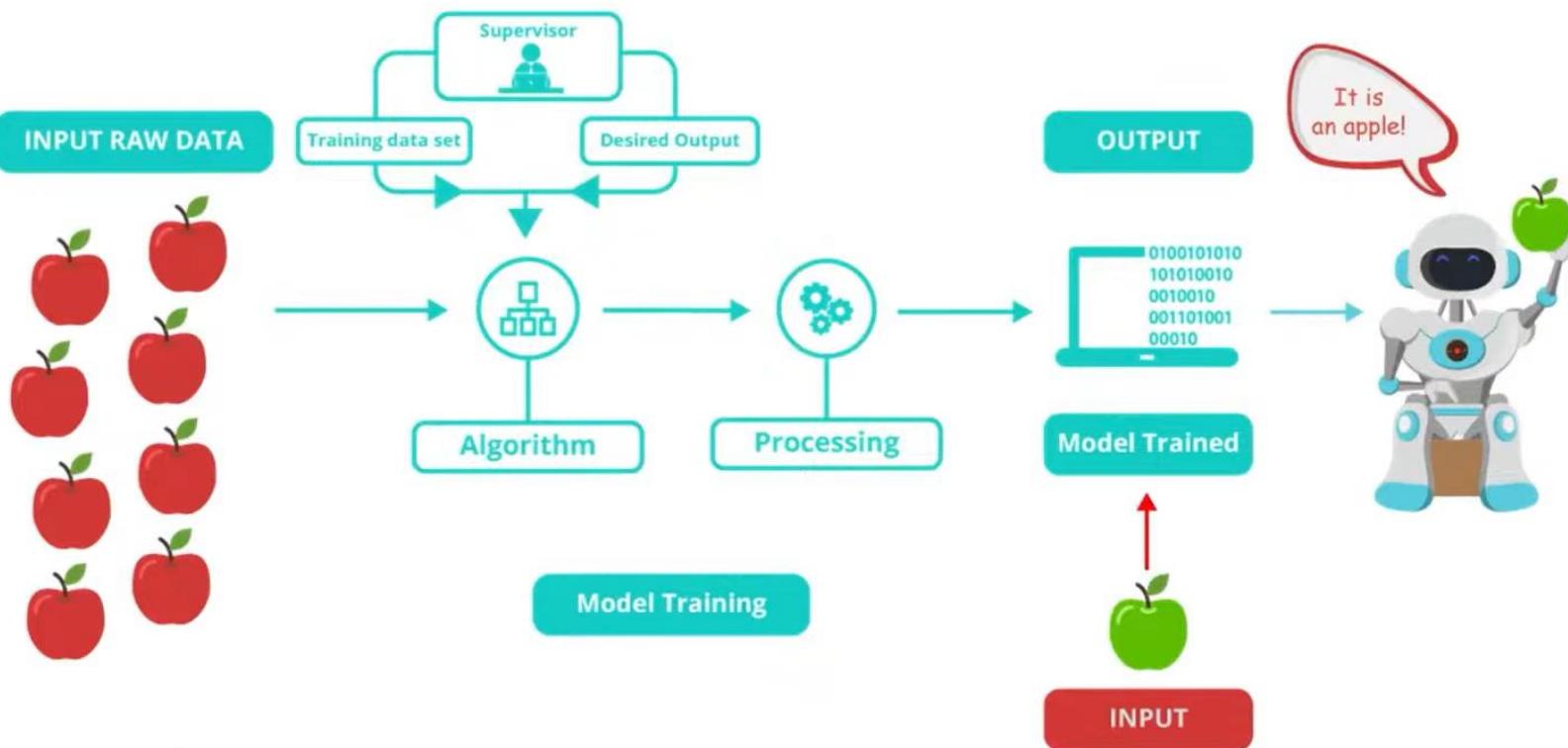


# Supervised Learning

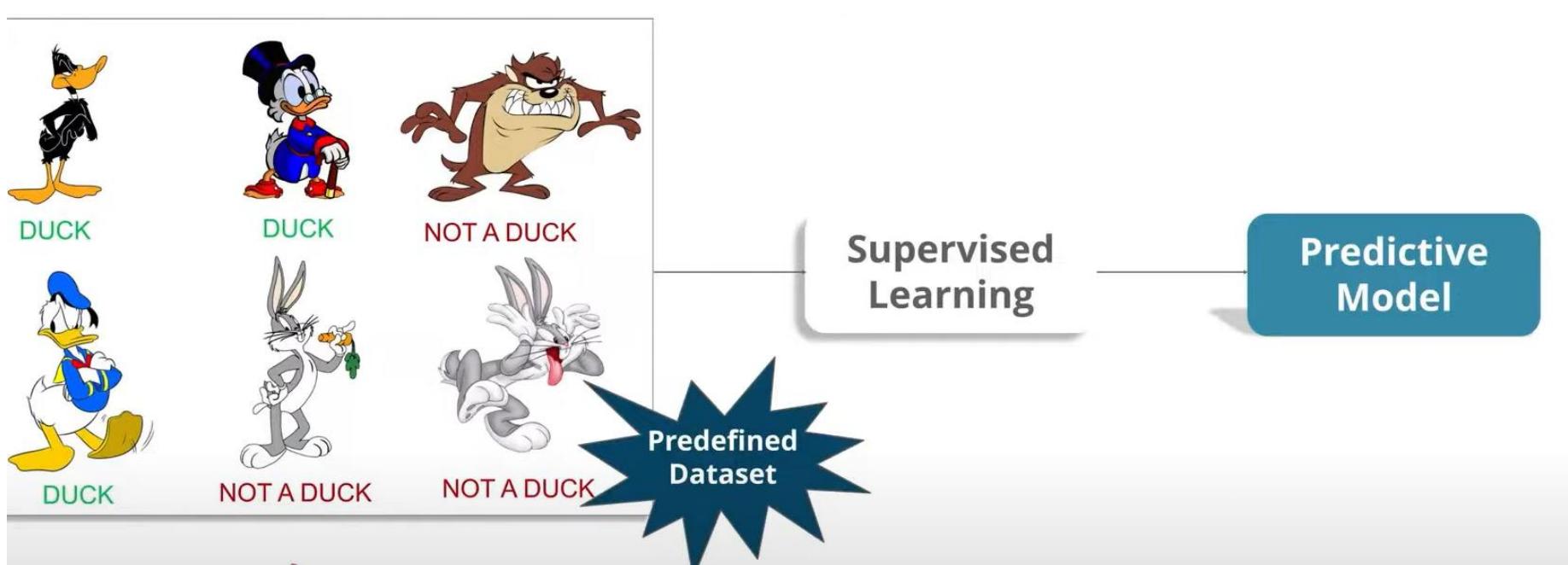
---

- Supervised learning categories and techniques
  - Linear classifier (numerical functions)
  - Parametric (Probabilistic functions)
  - Naïve Bayes, Gaussian discriminant analysis (GDA), Hidden Markov models (HMM), Probabilistic graphical models
  - Non-parametric (Instance-based functions)
  - K-nearest neighbors, Kernel regression, Kernel density estimation, Local regression
  - Non-metric (Symbolic functions)
  - Classification and regression tree (CART), decision tree
  - Aggregation
  - Bagging (bootstrap + aggregation), Adaboost, Random forest

# Supervised Learning



# Supervised Learning



# Supervised Learning

Linear Regression



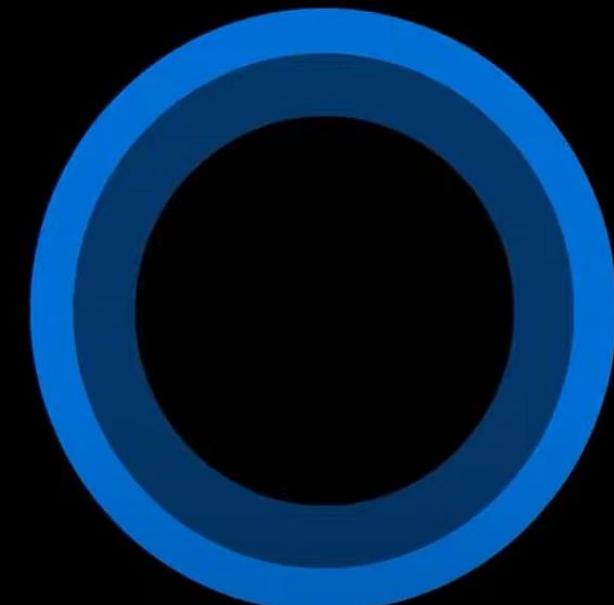
Support Vector

Random Forest





# Supervised Learning Use Cases



Hi. I'm Cortana.  
Ask me a question!

# Supervised Learning Use Cases



# Supervised Learning Use Cases





# Supervised Learning Use Cases



**Supervised Learning in Bank**  
Predict Credit Worthiness of Credit Card Holders



# Supervised Learning Use Cases

Supervised Learning in Healthcare  
Predict Patient Readmission Rates

A large, semi-transparent teal rectangular box is positioned in the center of the slide. It contains the title "Supervised Learning in Healthcare" in a large, white, sans-serif font. Below it, in a slightly smaller font, is the subtitle "Predict Patient Readmission Rates". The background of the slide shows a close-up photograph of a healthcare worker wearing blue scrubs and a stethoscope around their neck. Their hands are clasped in front of them. The overall composition suggests a professional and medical context.

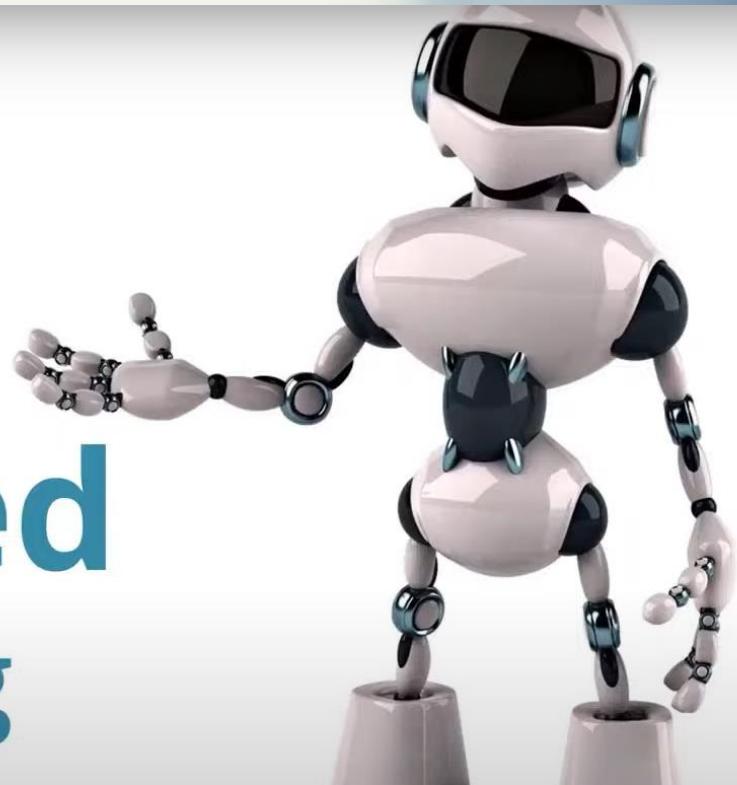
# Supervised Learning Use Cases

## Supervised Learning in Retail

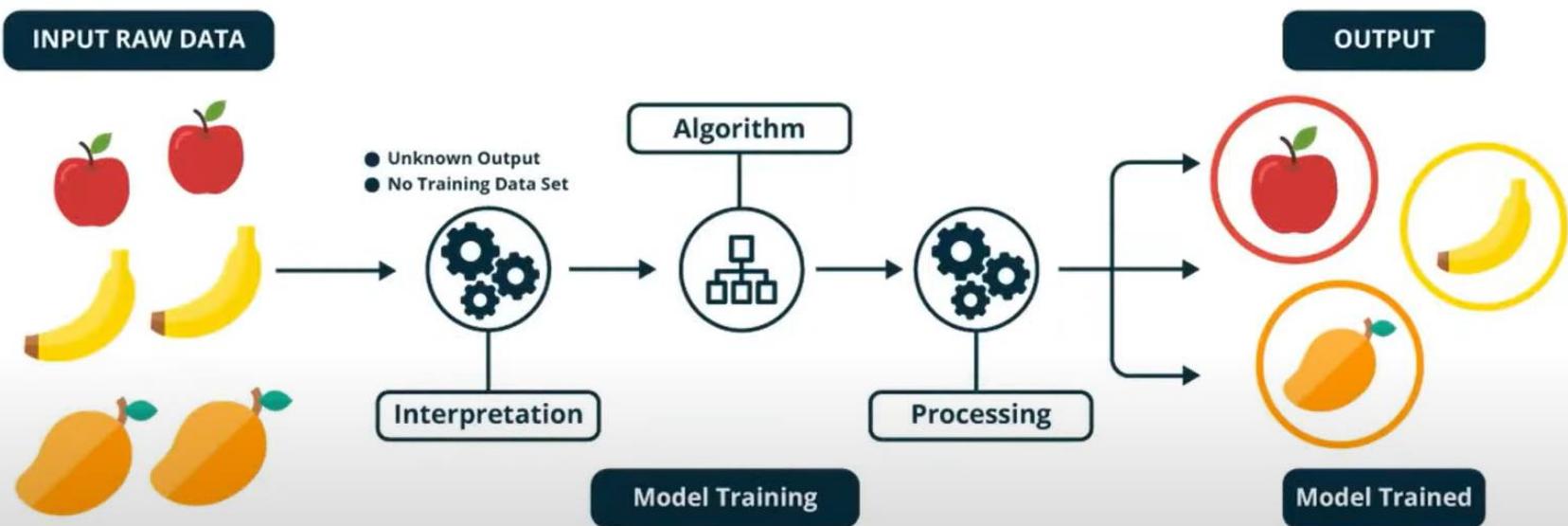
Analyse Products Customers Buy Together

variables that best correlate with  
re-admission next comes the retail

# Unsupervised Machine Learning



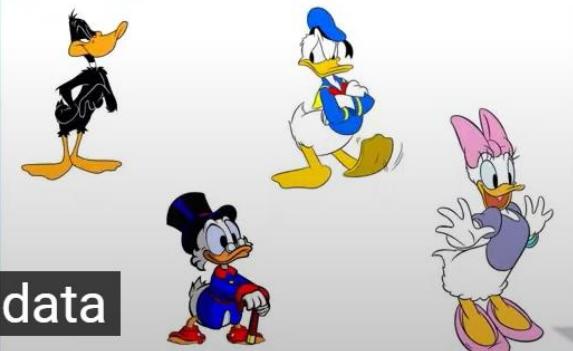
# Unsupervised Learning



# Unsupervised Learning



Unsupervised  
Learning



algorithm only knows which data

# Types of Unsupervised Learning

Apriori Algorithm



K- Means Algorithm



# Types of Unsupervised Learning

---

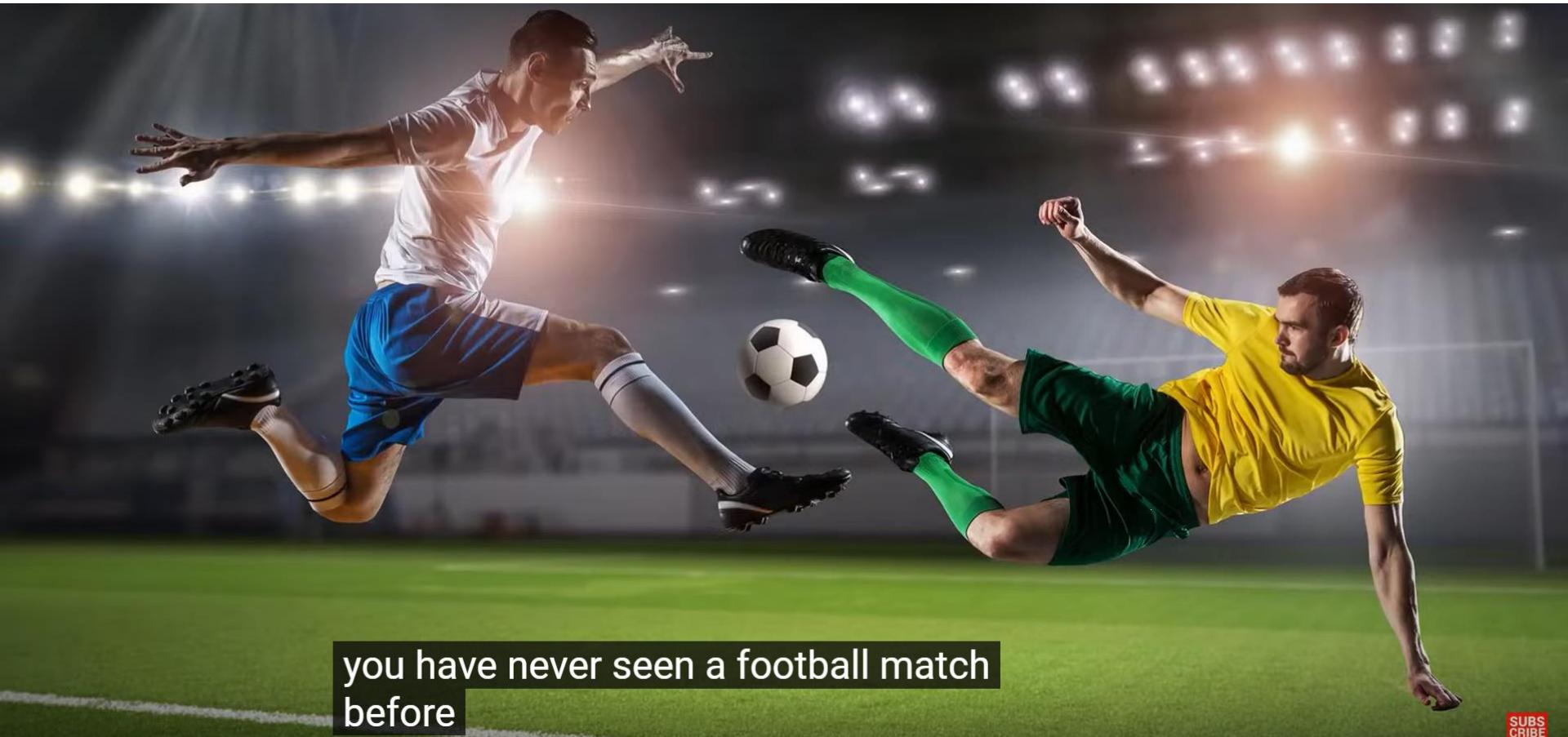
- Unsupervised learning categories and techniques
  - Clustering
  - K-means clustering
  - Spectral clustering
  - Density Estimation
  - Gaussian mixture model (GMM)
  - Graphical models
  - Dimensionality reduction
  - Principal component analysis (PCA)
  - Factor analysis

# Use cases - Unsupervised Learning



**Friend Invites You to a Party Full of Strangers and You are going to group them by age, qualifications and behaviour**

# Use cases - Unsupervised Learning



you have never seen a football match  
before

SUBS  
CRIBE

You start classifying them by the way you observe either by jersey, playing style or by defender etc.,



DATA SCIENCE

# Use cases - Unsupervised Learning

## Unsupervised Learning in Bank Segment Customers by Behavioural Characteristics





# Use cases - Unsupervised Learning

## Unsupervised Learning in Healthcare

Categorize MRI Data by Normal or Abnormal

A photograph showing a close-up of a healthcare professional's hands, wearing blue scrubs and a stethoscope, gently holding a patient's arm. The background is slightly blurred.

sector it is used to categorize the MRI  
data by normal or abnormal images

SUBSCRIBE

# Use cases - Unsupervised Learning

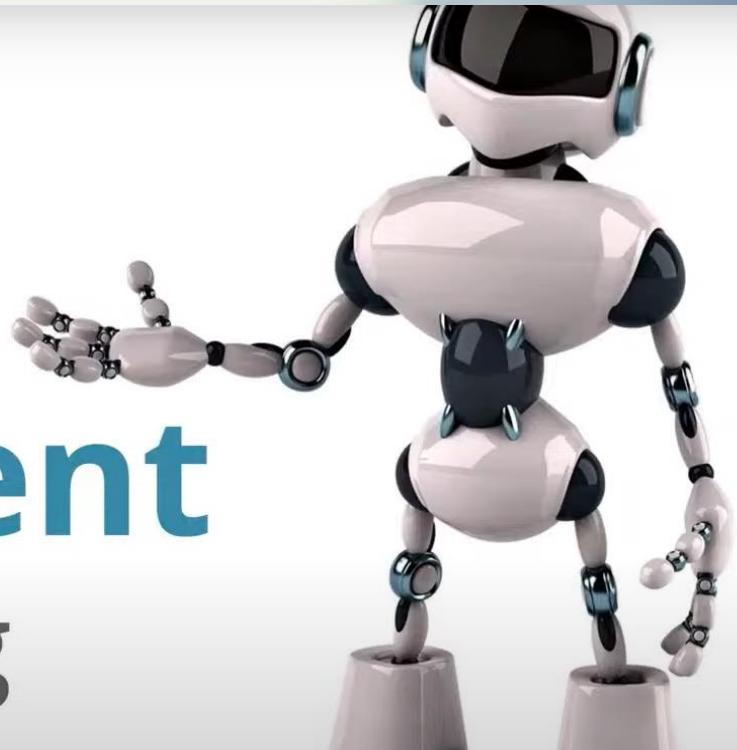
## Unsupervised Learning in Retail

Recommend Products to Customers Based on Past Purchase



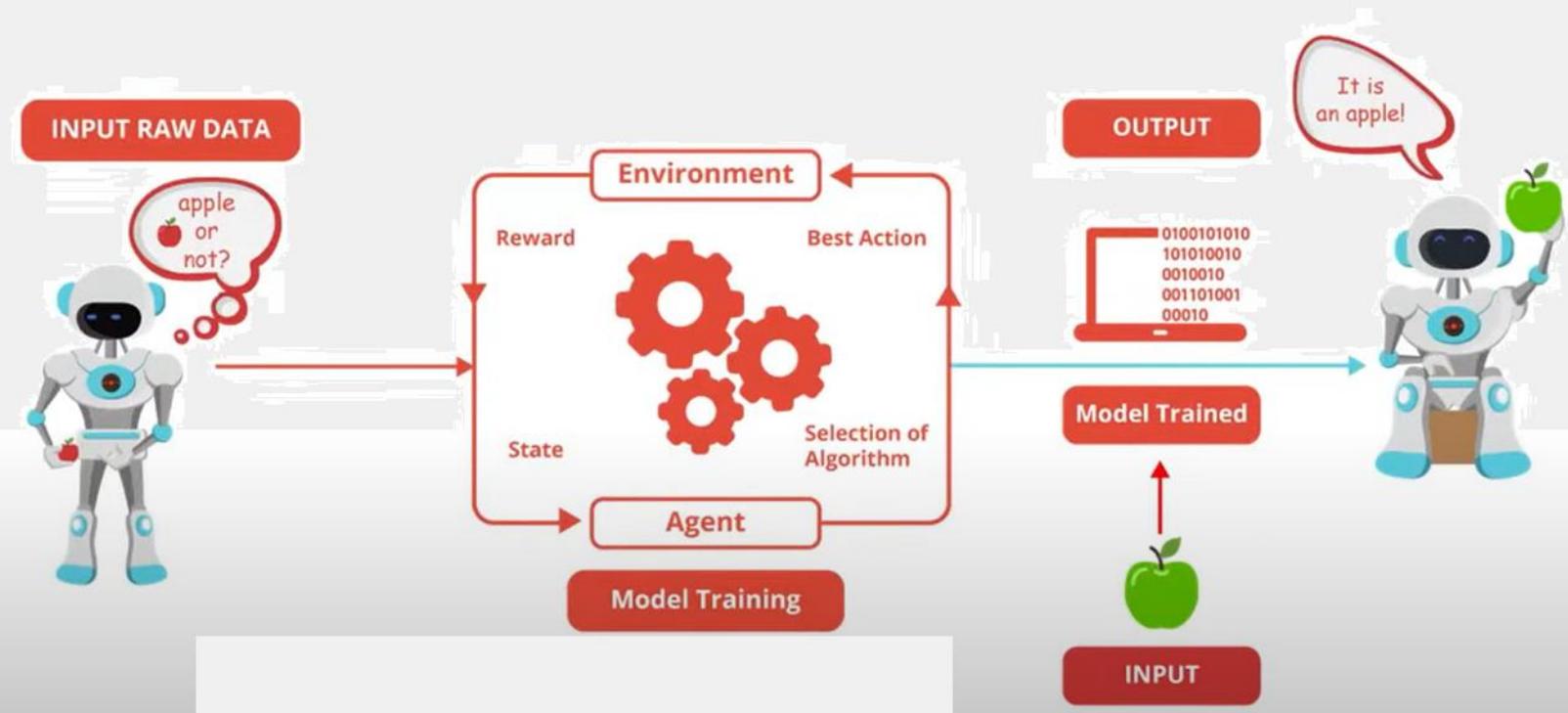


# Reinforcement Machine Learning



# Reinforcement Learning

## Exploration and Exploitation



# Reinforcement Learning





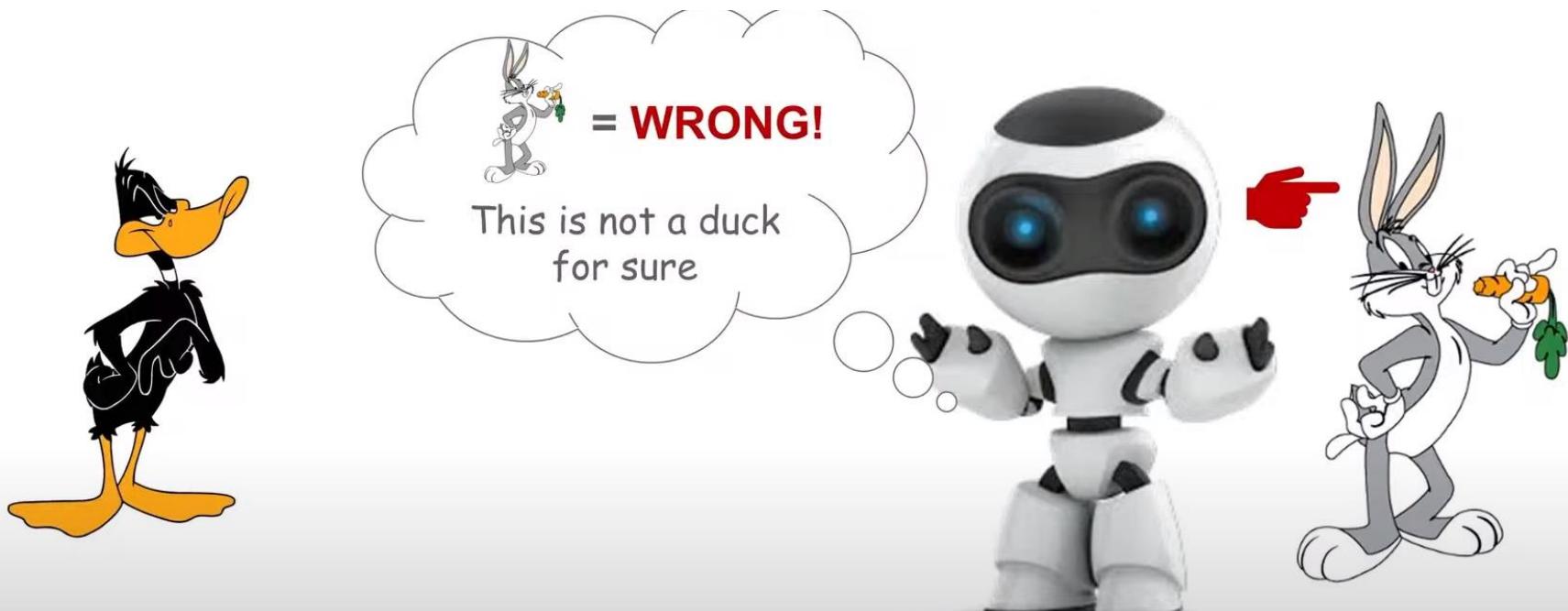
DATA SCIENCE

# Reinforcement Learning

Machine made wrong decision by choosing bunny as duck so penalized



# Reinforcement Learning



# Reinforcement Learning



## Reinforcement Examples & Use Cases



# Reinforcement Learning





DATA SCIENCE

# Reinforcement Learning Use cases

## Reinforcement Learning in Bank

Create 'Next Best Offer' Model for the Call Centre



# Reinforcement Learning Use cases

## Reinforcement Learning in Healthcare

Allocate Scarce Medical Resources to Handle Different ER Cases



# Reinforcement Learning Use cases



**Reinforcement Learning in Retail**  
Reduce Excess Stock with Dynamic Pricing

# Data Science

*Data Science is the process of extracting useful insights from data by using a variety of tools, algorithms and Machine Learning fundamentals.*

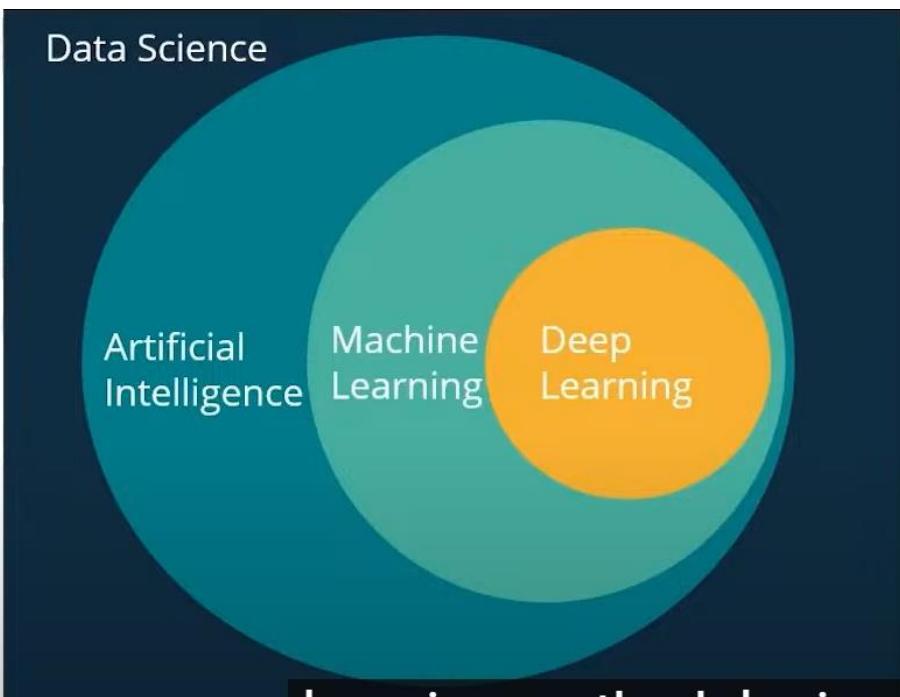


**Netflix analyses the viewers interest based on that movies, web series are produced**



# Fields Of Data Science

# Data Science

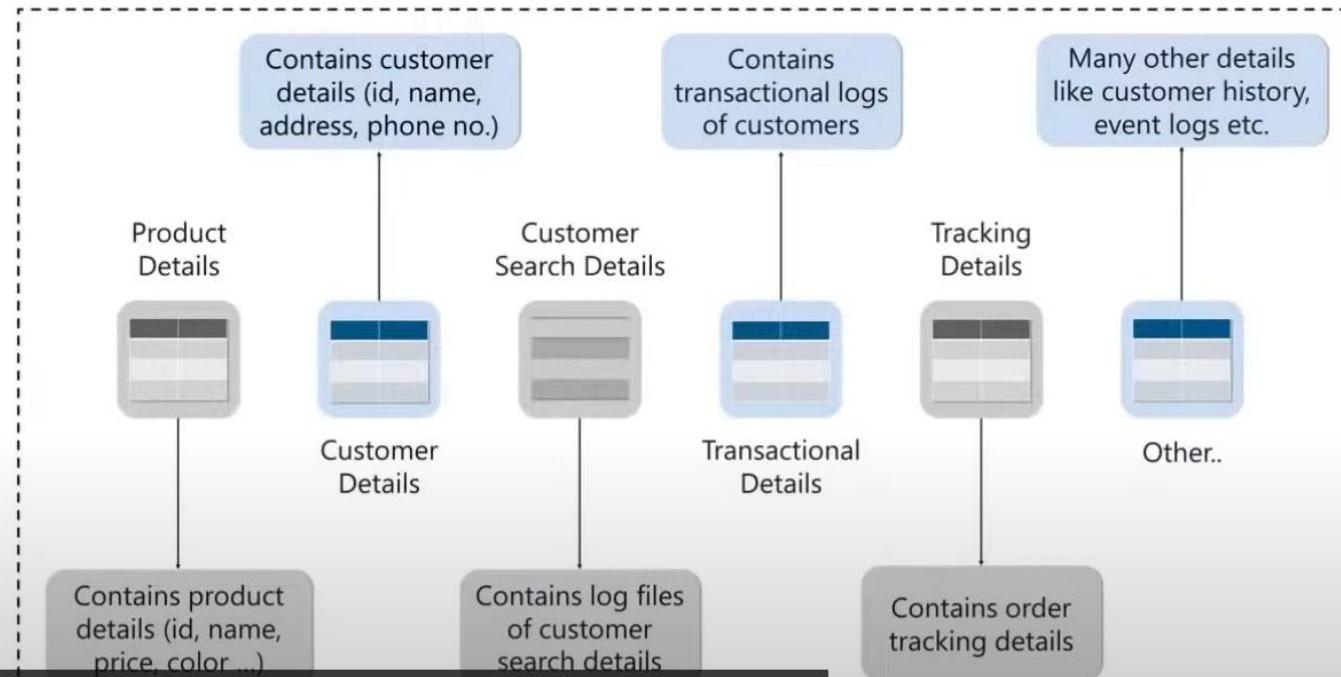


learning methodologies in order to analyze data and then extract useful

- *Data science* is the extraction of knowledge from data by using different techniques and algorithms
- *Artificial Intelligence* is a technique which enables machines to mimic human behaviour
- *Machine Learning* is a subset of AI technique which uses statistical methods to enable machines to improve with experience
- *Deep learning* is a subset of ML which make the computation of multi-layer neural network feasible

# Data Science

## Recommendation System which suggest laptop bag when you buy laptop



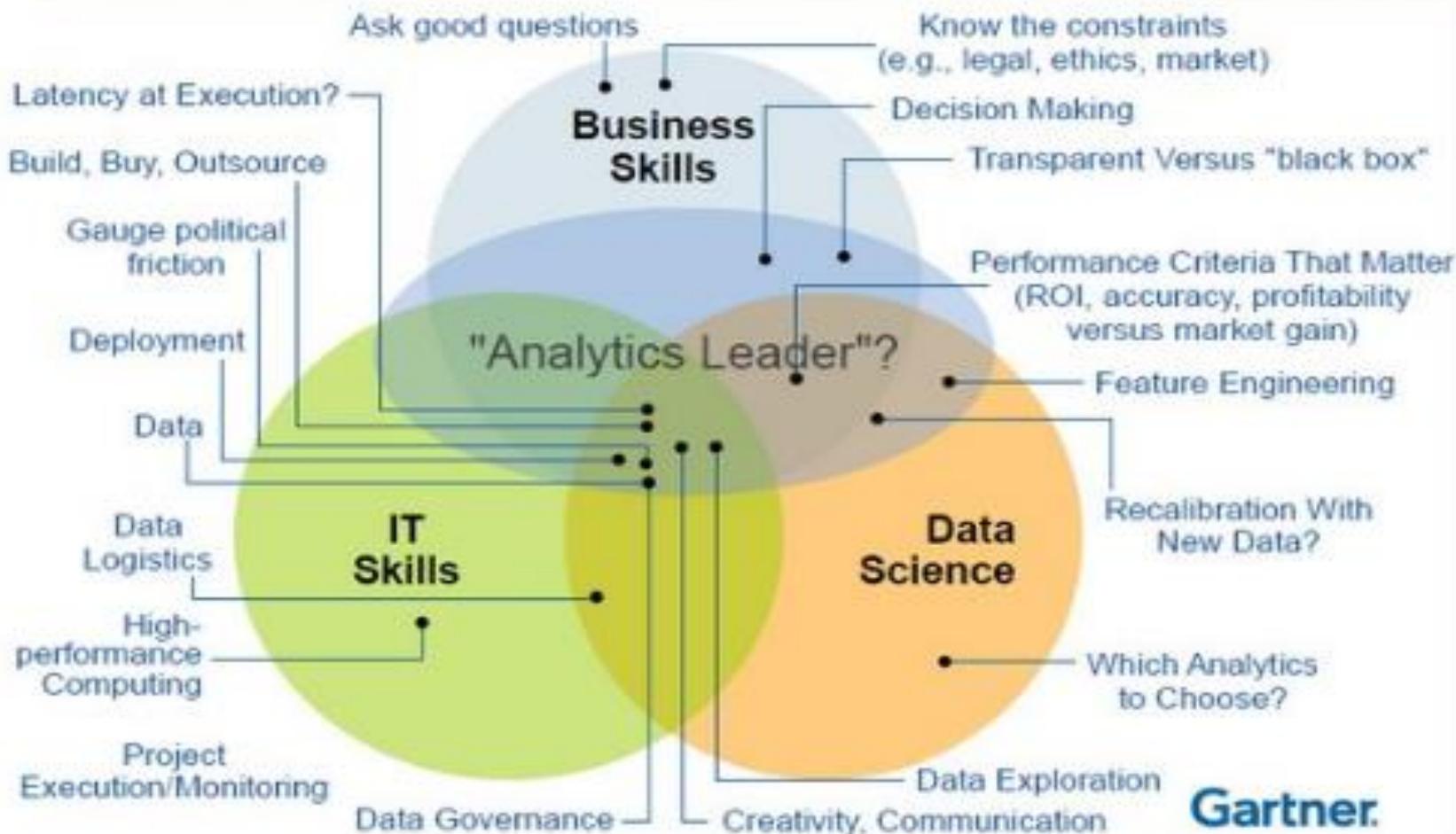
why companies like Amazon Walmart and Netflix are doing so well is because of



DATA SCIENCE

# Data Science

## Driving the Success of Data Science Solutions: Skills, Roles and Responsibilities ...



Gartner

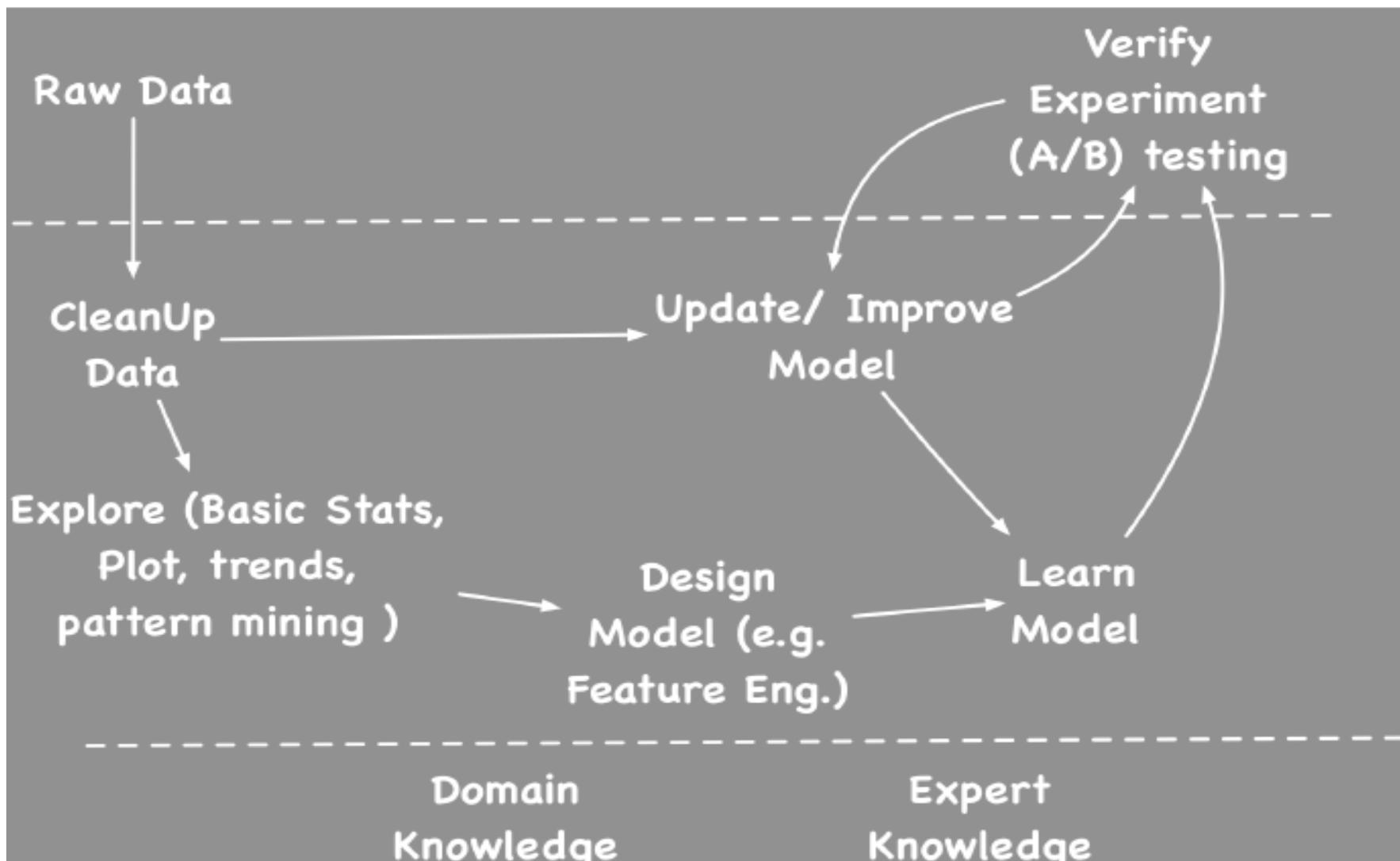
# Data Science

Extraction of knowledge from large volumes of data that are structured or unstructured.

It is a continuation of the fields **data mining** and **predictive analytics**



# Data Science Pipeline



## Examples of Big Data Analytics

---

- Using Big Data to win elections
- Big Data for finding a perfect match(Matrimony.com)
- Big Data for detecting water leakages(Bangalore Water supply and sewage system)
- Big Data for gaining insights into shopping behavior

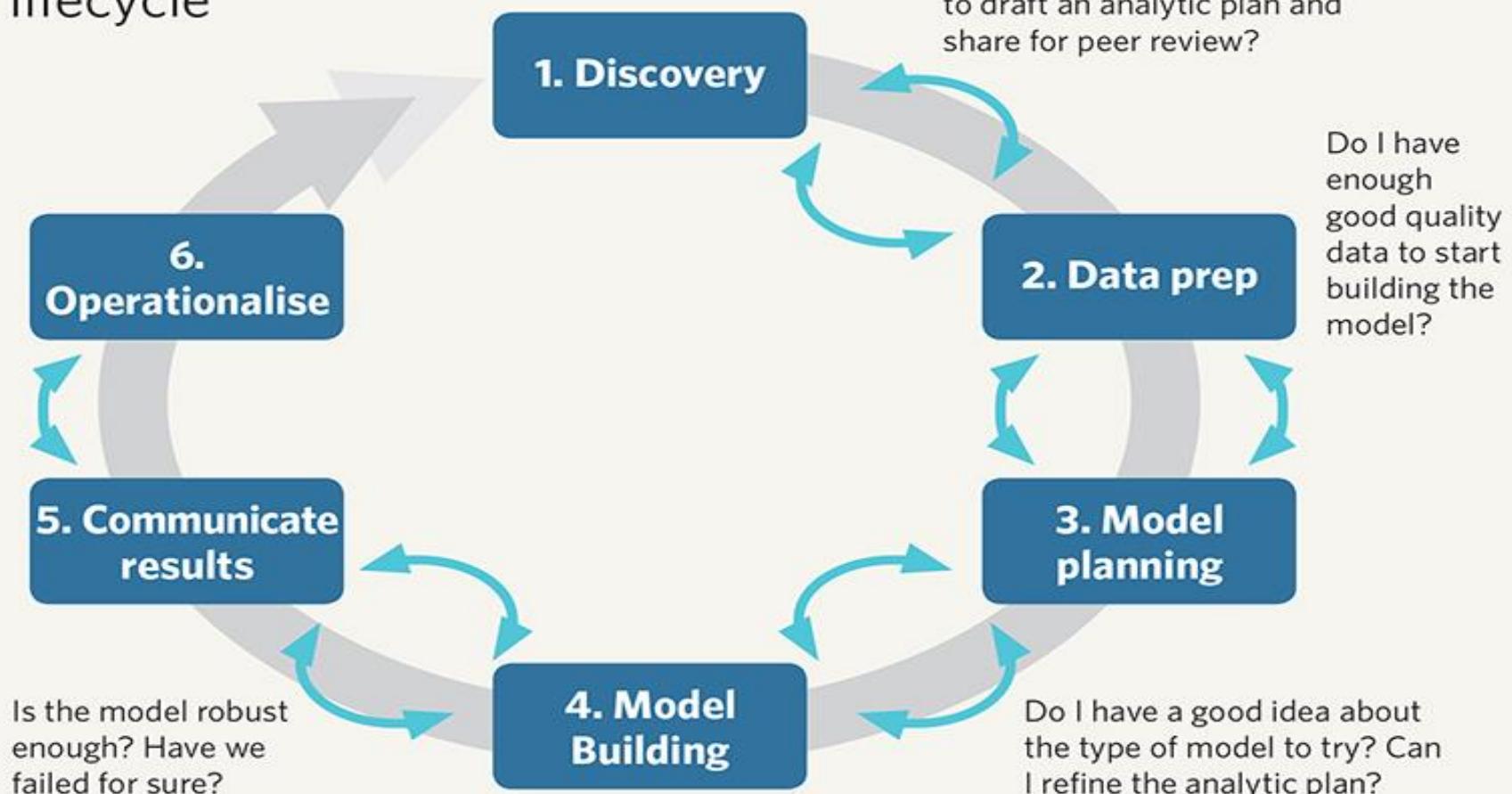
## Examples of Big Data Analytics

---

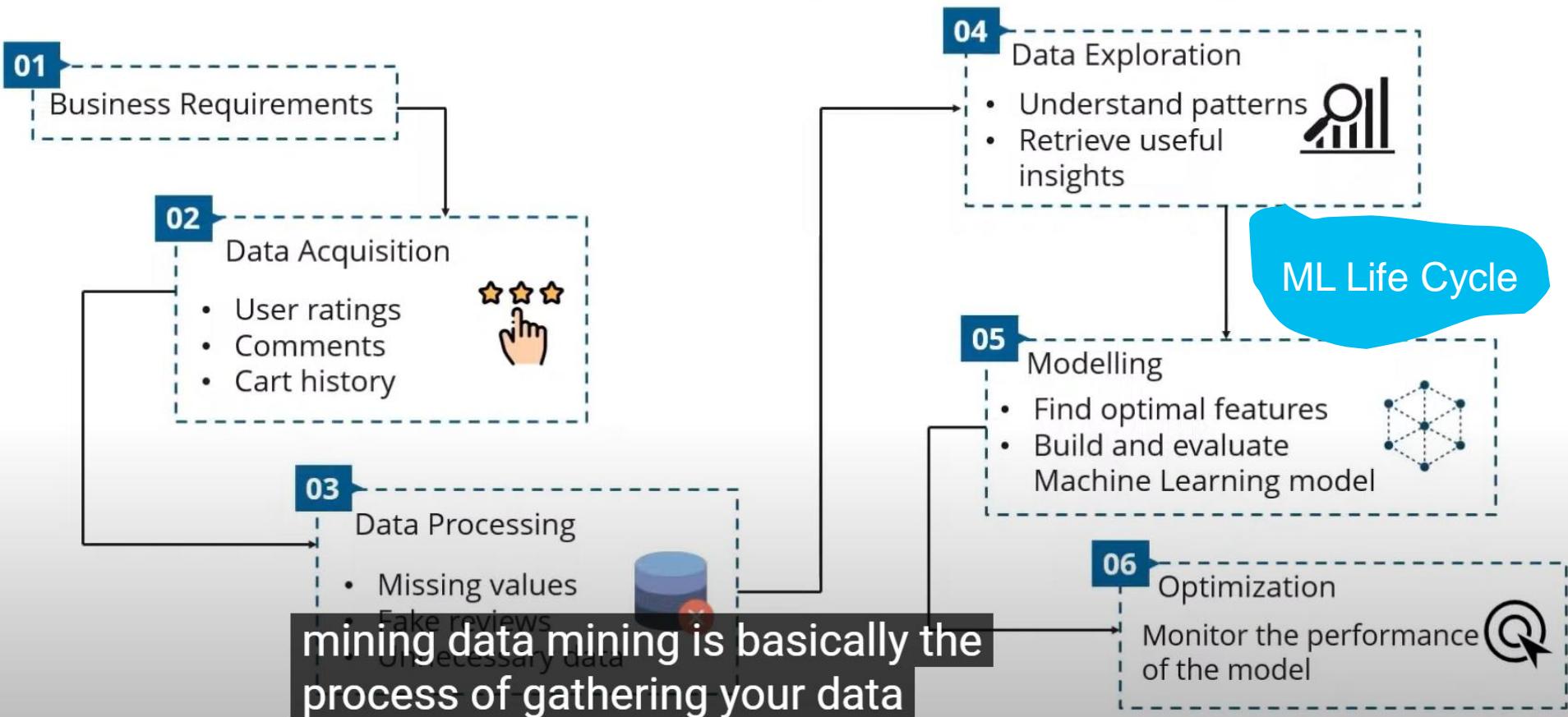
- Big Data for ensuring proper water supply(Kerala Water Supply System)
- Using Big Data to improve India's financial inclusion ratio(Micro-finance firm, Janalakshmi Financial Services)
- Using Big Data to improve product development(Reliance Games)
- Using Big Data to predict ticket confirmations for trains(PNR prediction)

# Data Science Lifecycle

## Data Analytics lifecycle

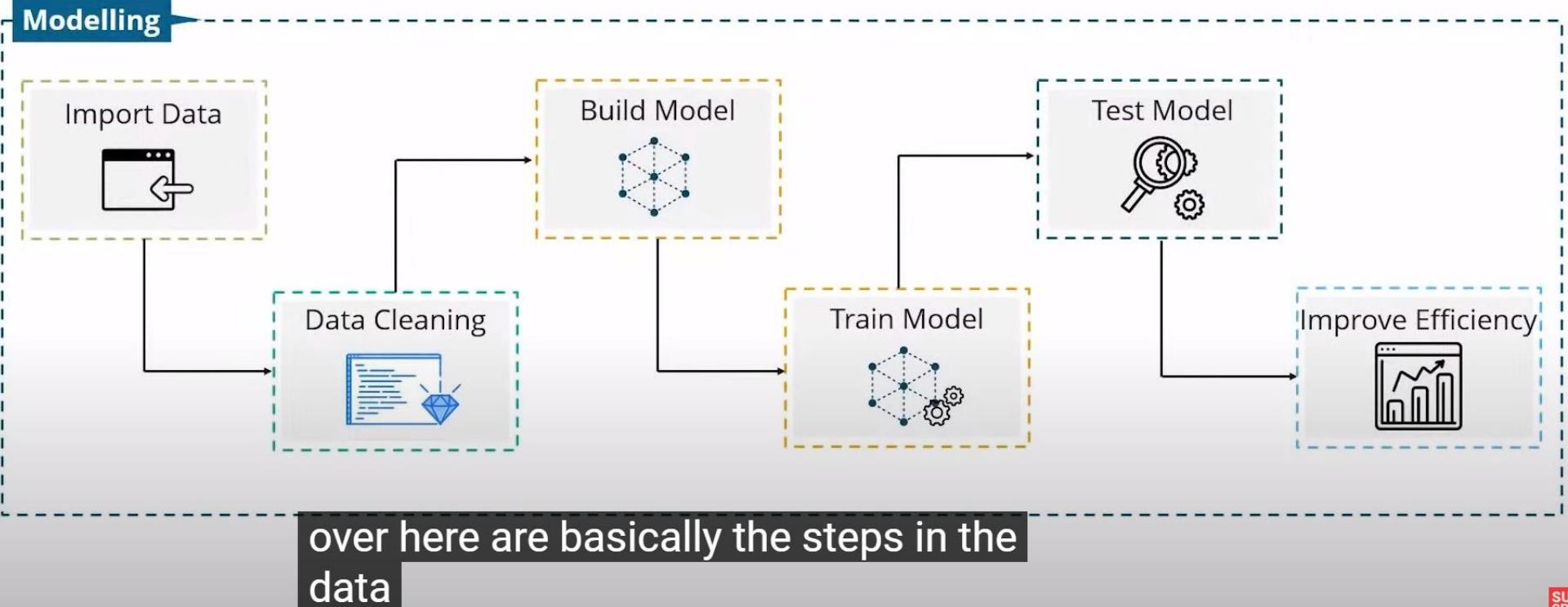


# Data Science Life Cycle

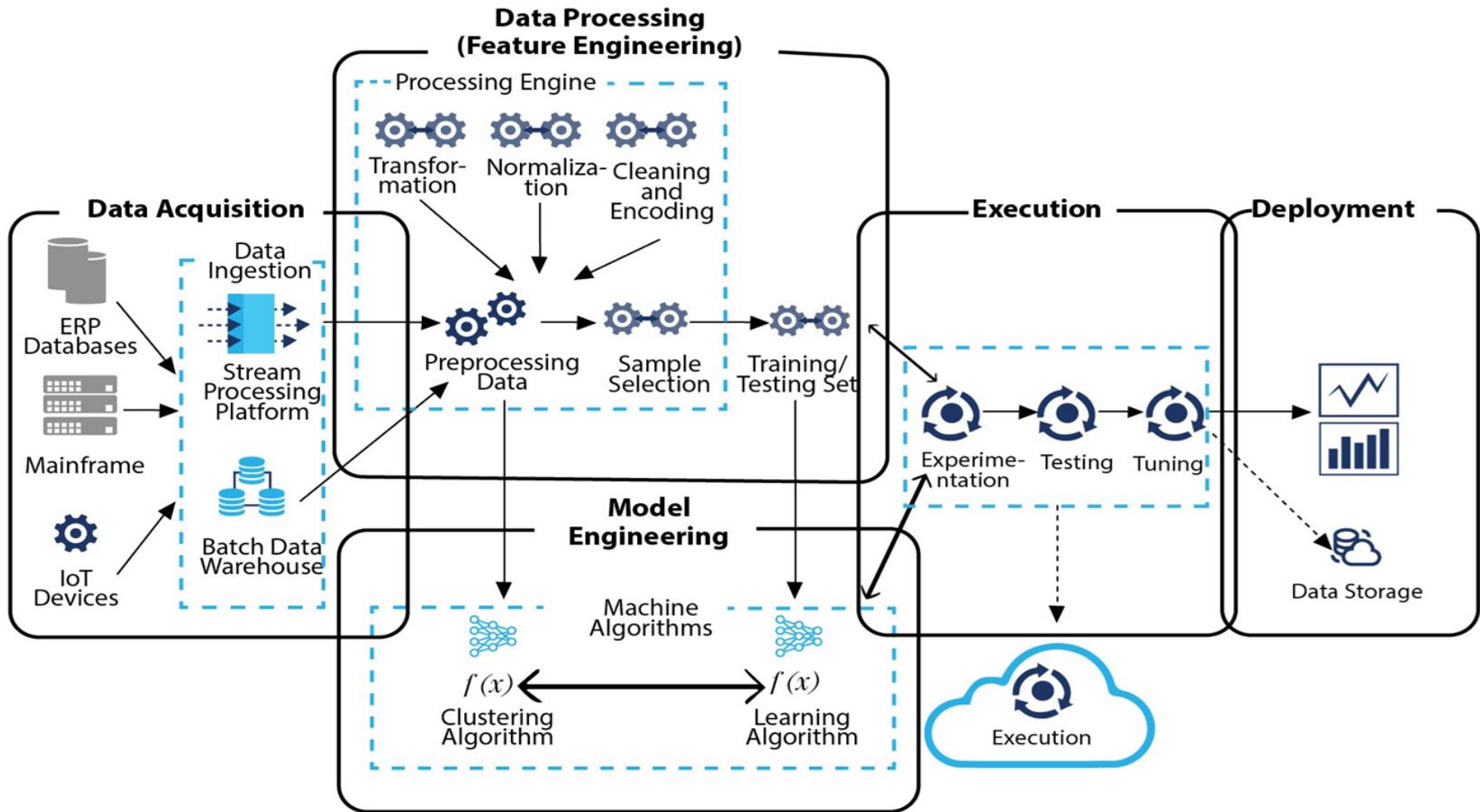


mining data mining is basically the process of gathering your data

# Machine Learning Lifecycle



# Machine Learning Lifecycle

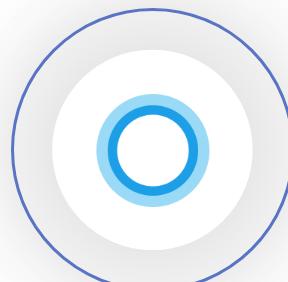


# AI and ML Use cases

## The Most Popular AI Solutions



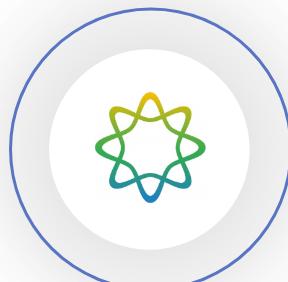
Siri



Cortana



Google Assistant



Elsa



FaceApp



Alexa

# AI and ML Use cases

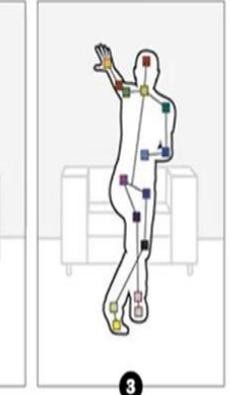
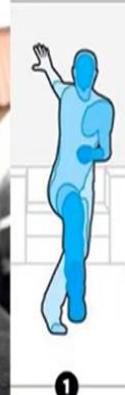
1



3



5





DATA SCIENCE

# AI and ML Use cases

6

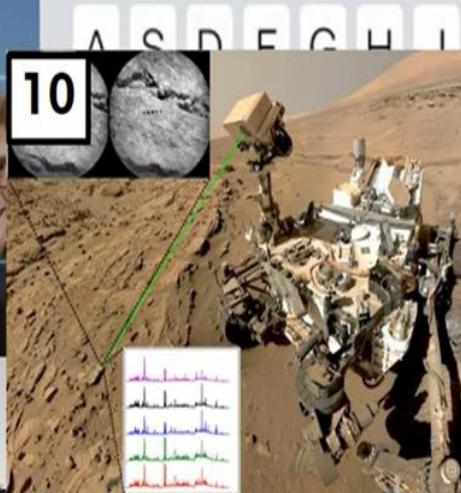


facebook  
Ads

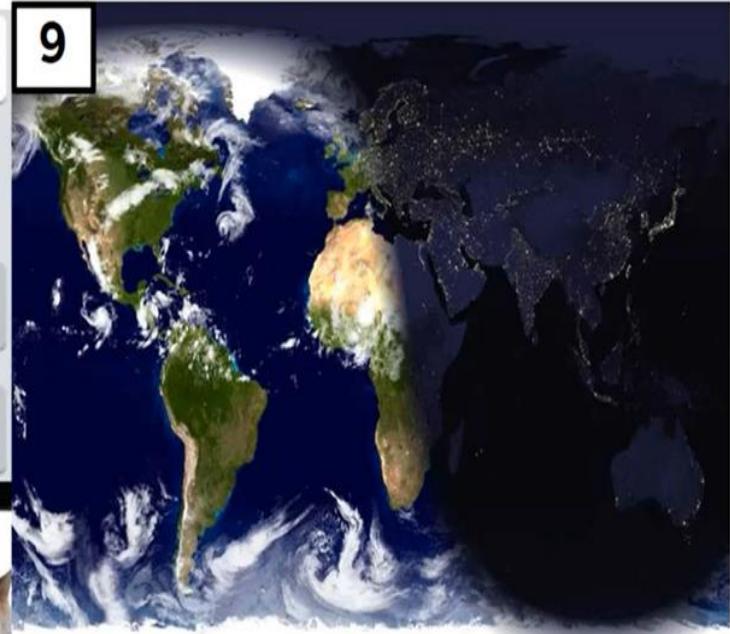


4

10



9



3



8

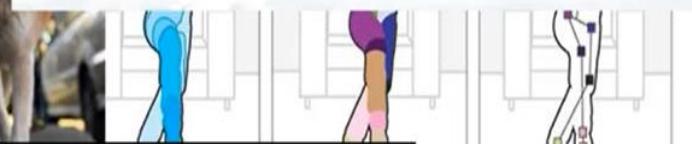


7

ama

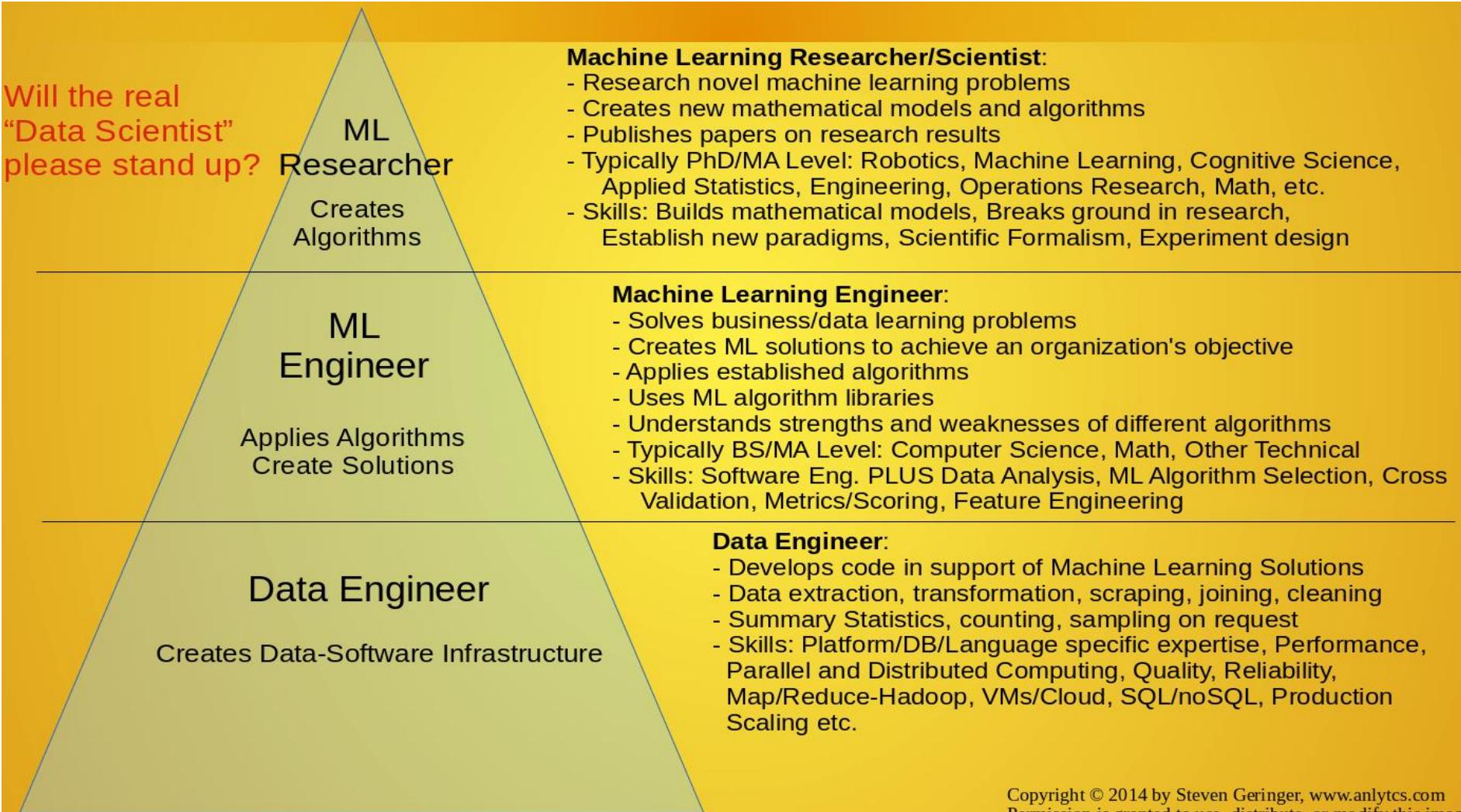
NETFLIX

audible



# ML Skill Pyramid

Will the real  
“Data Scientist”  
please stand up?





# How to become ML Engineer



# Machine Learning Engineer



## Who is a Machine Learning (ML) Engineer?



Machine learning engineers are sophisticated programmers who **develop** machines and systems that can learn and apply knowledge without specific direction

# Machine Learning Engineer



## What does an ML Engineer do?



Machine learning engineers are creators of the algorithms that allow a machine to find patterns in its own programming data, teaching it to understand commands and even think for itself

# Machine Learning Engineer



## Skills Needed to Become an ML Engineer

### Programming Skills

- **R**: Used for developing Statistical Software and Data Analysis
- **Python**: Lets you Create, Analyse and Organize large chunks of Data with ease
- **Java Programming**: Data Description

# Machine Learning Engineer

## Roles & Responsibilities of an ML Engineer

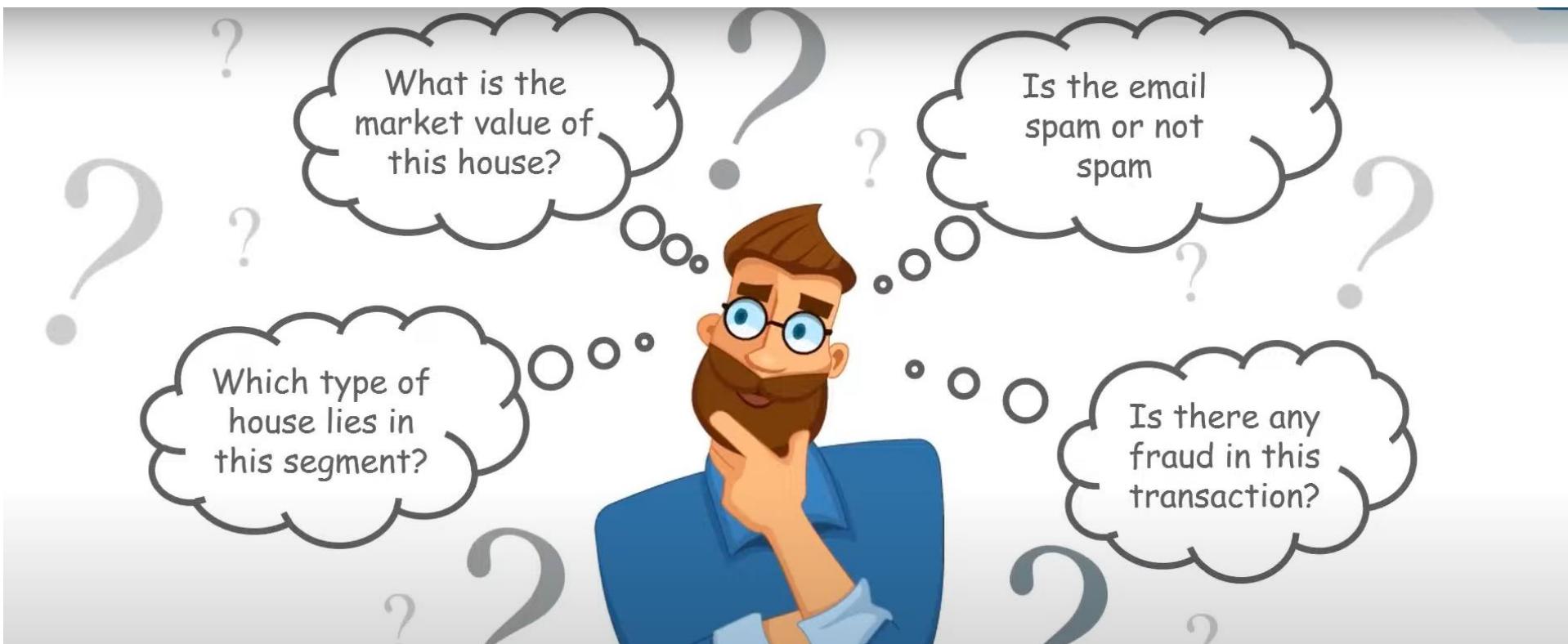
- ✓ Study some Prototypes & Transform them into Applications
  - ✓ Design and Build Machine Learning Systems
  - ✓ Find Appropriate Algorithms & Tools
  - ✓ Develop Machine Learning Applications
  - ✓ Select the right dataset & find the correct Data Representation Methods
  - ✓ Run Machine Learning Tests and Experiments
  - ✓ Train the Systems for Top-notch Accuracy
- lastly we need to train the systems for top Notch**





DATA SCIENCE

# Machine Learning Engineer





DATA SCIENCE

## ML is the Future

---

**2005 – 130 EXABYTES**

**2010 – 1,200 EXABYTES**

**2015 – 7,900 EXABYTES**

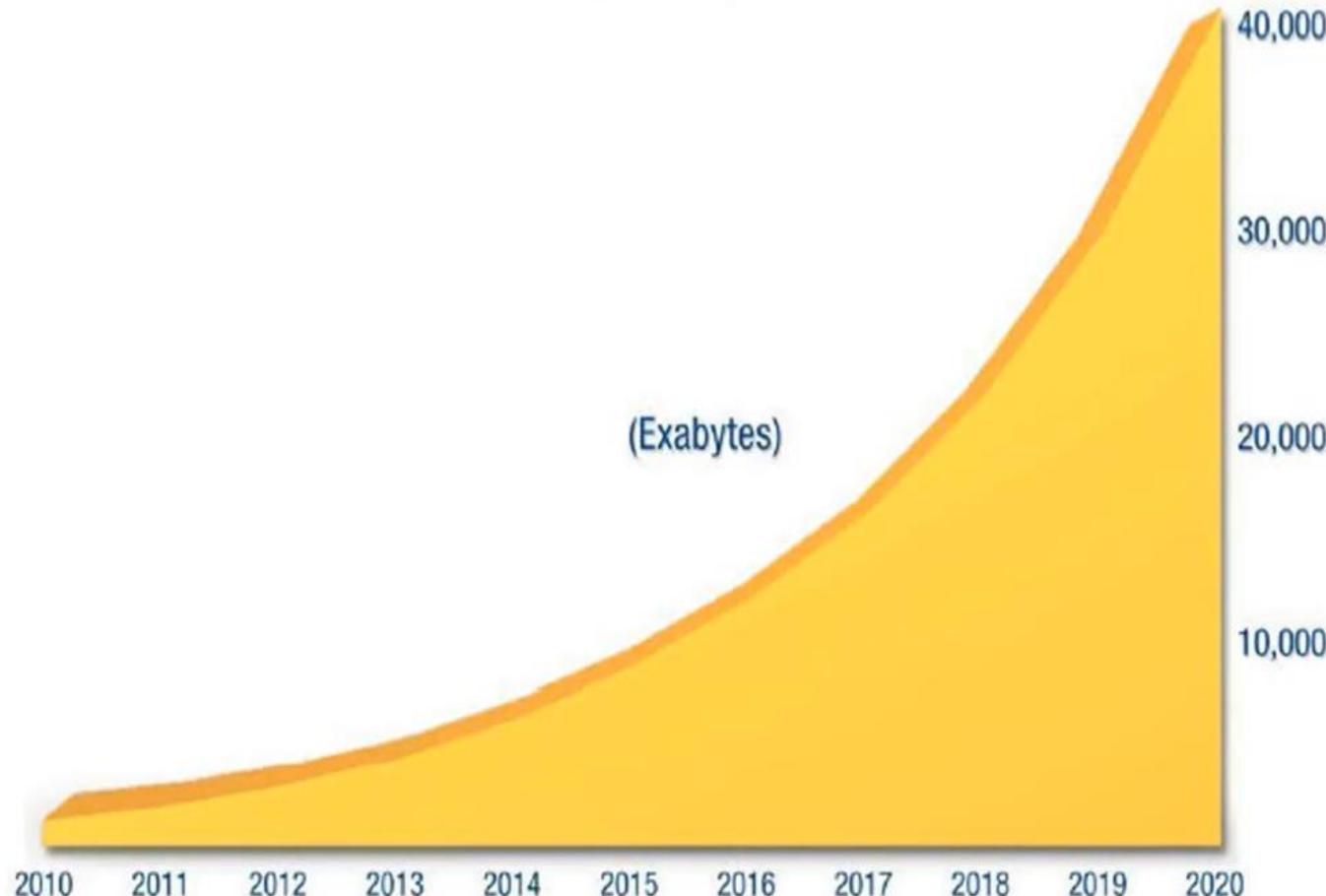
**2020 – 40,900 EXABYTES**



DATA SCIENCE

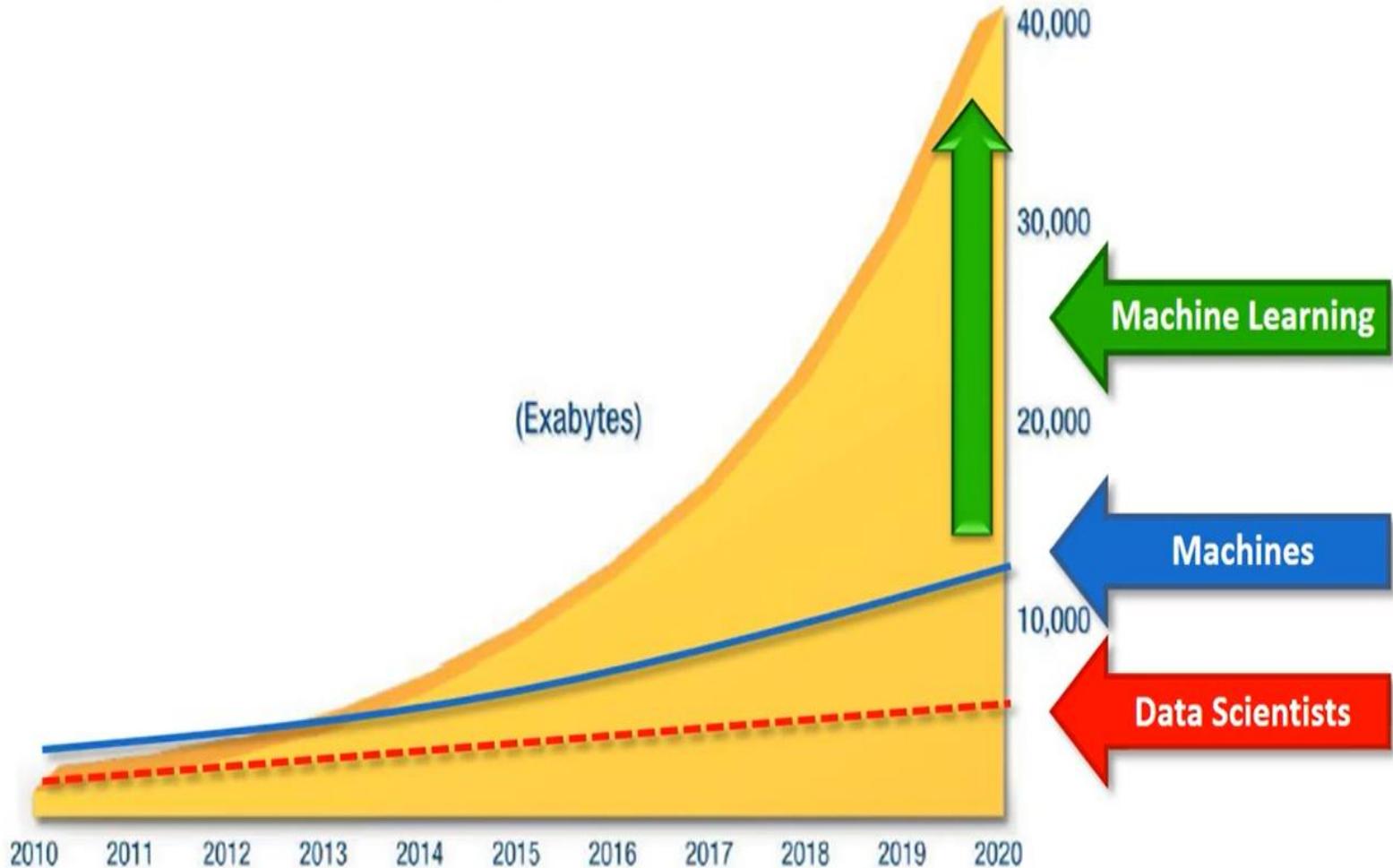
# ML is the Future

50-Fold Growth from the Beginning of 2010 to the end of 2020



# ML is the Future

50-Fold Growth from the Beginning of 2010 to the end of 2020



# Machine Learning Algorithms How to Choose

## Classification Algorithm

Category is predicted using the data

- Is the person a male or a female?
- Is the mail spam or non-spam?

these category of question would fall under the classification algorithm



# Machine Learning Algorithms How to Choose

## Anomaly Detection Algorithm

Identify unusual data points

in business like intrusion detection  
like identifying strange patterns in the

- Is there any fraud in this transaction?
- Is someone trying to hack our network?



# Machine Learning Algorithms How to Choose

## Clustering Algorithm

Groups data based on some condition

- Which type of house lies in this segment?
- What type of customer buys this product?

?

# Machine Learning Algorithms How to Choose

## Regression Algorithm

Data itself is predicted

- What is the market value of this house?
- Is it going to rain tomorrow?

important and broadly used machine learning and statistics tool it

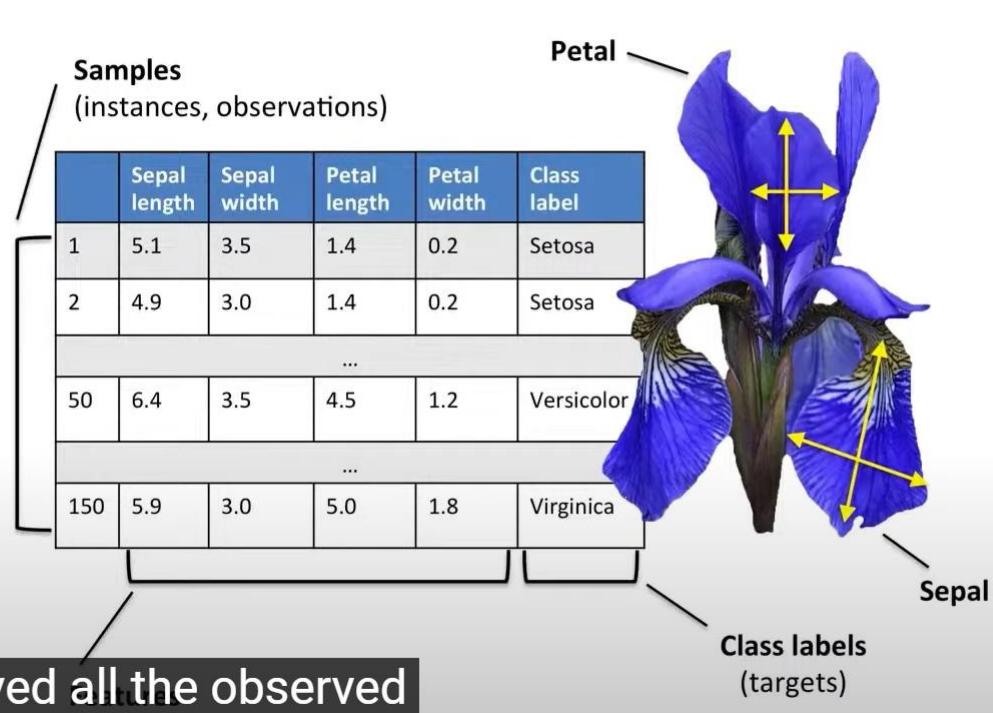


# Machine Learning Algorithms How to Choose



## Iris Dataset

the flower observed all the observed  
flowers belong to one of



# Machine Learning Algorithms How to Choose



## ANACONDA NAVIGATOR

[Sign in to Anaconda Cloud](#)



Home



Environments



Projects (beta)



Learning



Community

Documentation

Developer Blog

Feedback



Applications on

base (root)

Channels

Refresh



jupyterlab

0.31.4

An extensible environment for interactive and reproducible computing, based on the Jupyter Notebook and Architecture.

[Launch](#)



notebook

5.4.0

Web-based, interactive computing notebook environment. Edit and run human-readable docs while describing the data analysis.

[Launch](#)



qtconsole

4.3.1

PyQt GUI that supports inline figures, proper multiline editing with syntax highlighting, graphical calltips, and more.

[Launch](#)



spyder

3.2.6

Scientific Python Development EnviRonment. Powerful Python IDE with advanced editing, interactive testing, debugging and introspection features

[Launch](#)



glueviz

0.12.0

Multidimensional visualizations of scientific files. Explore relationships within and across multiple datasets.



orange3

3.3.1

Data selection and data analysis. Orange3 makes you more productive with R. Includes R essentials and notebooks.



rstudio

1.1.383

Integrated tools designed to help you be more productive with R. Includes R essentials and notebooks.



vscode

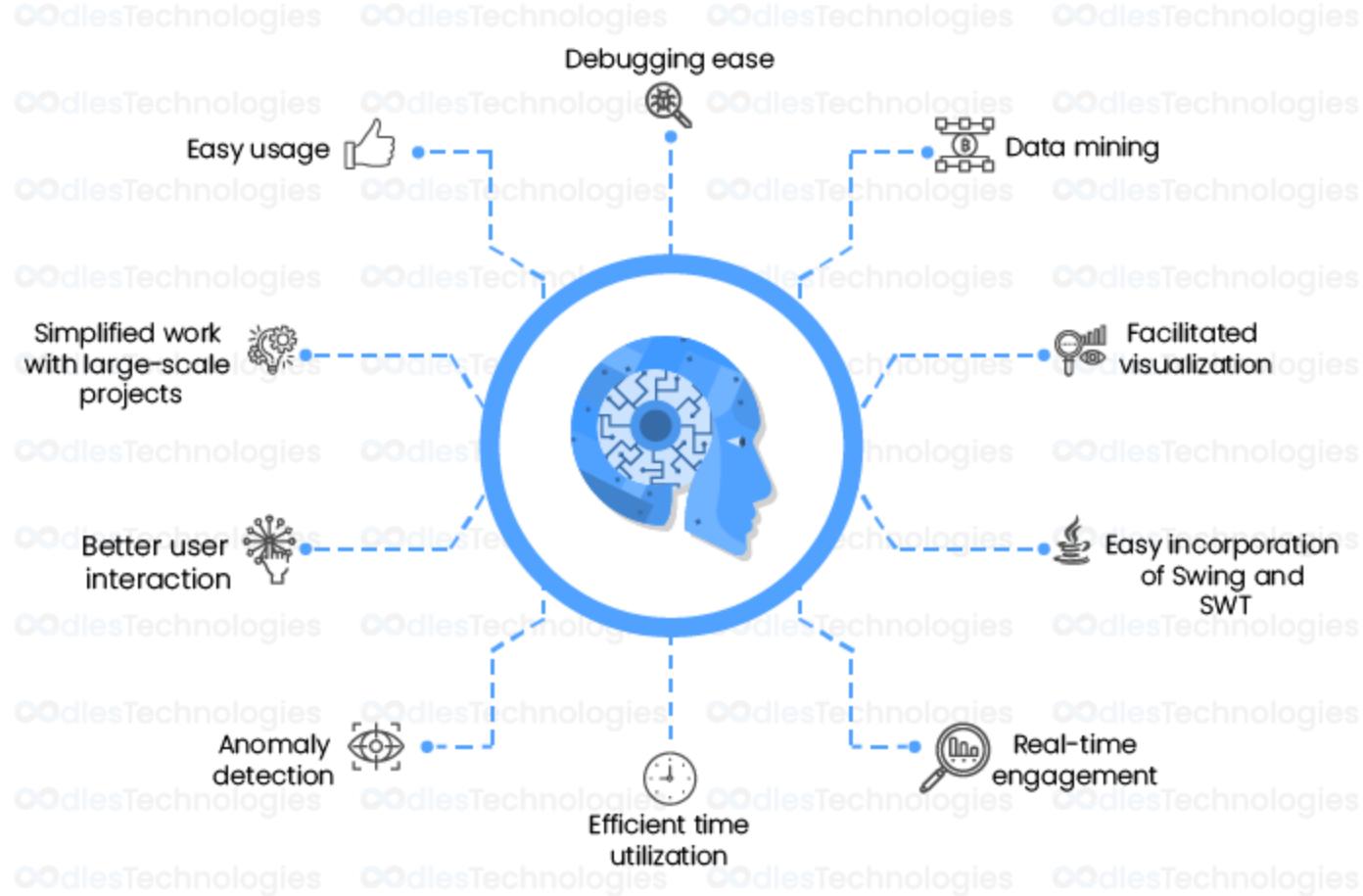
1.22.2

Streamlined code editor with support for development operations like debugging, task running and version control.

which is a web-based interactive Computing notebook

SUBSCRIBE

# Benefits of Programming AI in Java



## For Expert Systems

---

- Apache Jena — a framework to create web and linked data applications.
- PowerLoom — a platform to build knowledge-based applications and reasoning systems.
- d3web — a reasoning engine with many algorithms to solve given problems.
- Eye — a reasoning engine to perform semi-backward reasoning.
- Tweety — a collection of frameworks for logical aspects of AI and knowledge representation.

# For Neural Networks

---

- Neuroph — an open-source framework for neural network creation.
- Deeplearning4j — a deep learning library for JVM that also provides API for neural network creation.
- For Natural Language Processing
  - Apache OpenNLP — a toolkit to process the natural language text.
  - Stanford CoreNLP — a framework to perform NLP tasks.

## For ML

---

- Java-ML — a collection of machine learning algorithms.
- RapidMiner — a data science platform that provides machine learning algorithms through GUI and Java API.
- Weka — a collection of machine learning algorithms.
- Encog — a collection of advanced algorithms.
- For Genetic Algorithms
  - Jenetics — an advanced genetic algorithm.

## For ML

---

- Watchmaker — a framework for implementing genetic algorithms.
- ECJ 23 — a research framework with support for genetic algorithms.
- JGAP (Java Genetic Algorithms Package) — a genetic programming component.
- Eva — a simple OOP evolutionary algorithm framework.
- For automatic programming
  - Spring Roo — a lightweight developer tool.
  - Acceleo — a code generator for Eclipse which creates code from EMF models.

# Data Acquisition

---

- A data acquisition system is a system that comprises sensors, measurement devices, and a computer.
- A data acquisition system is used for processing acquired data, which involves collecting the information required to understand electrical or physical phenomena.



DATA SCIENCE

# Data Acquisition

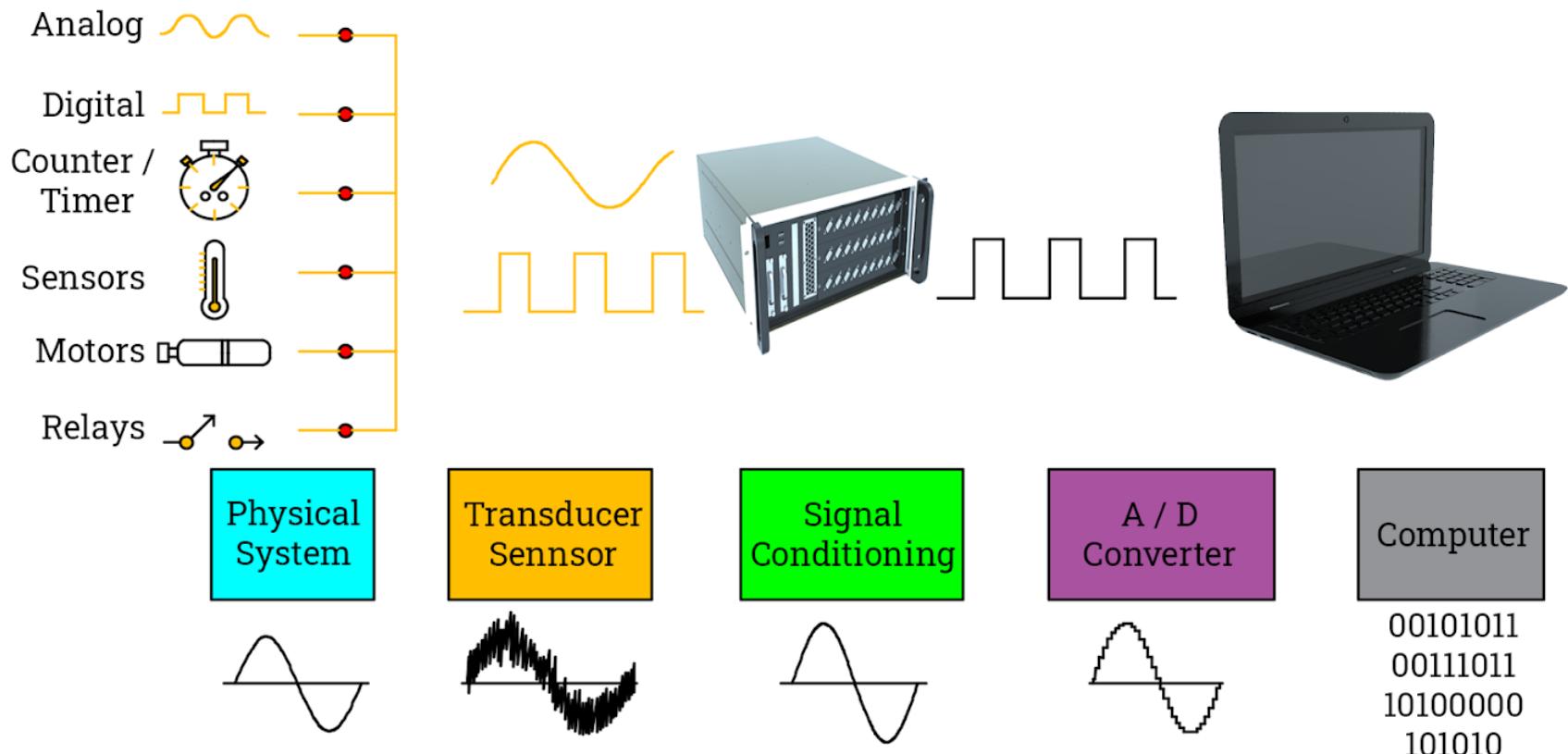


**Data Acquisition Systems**

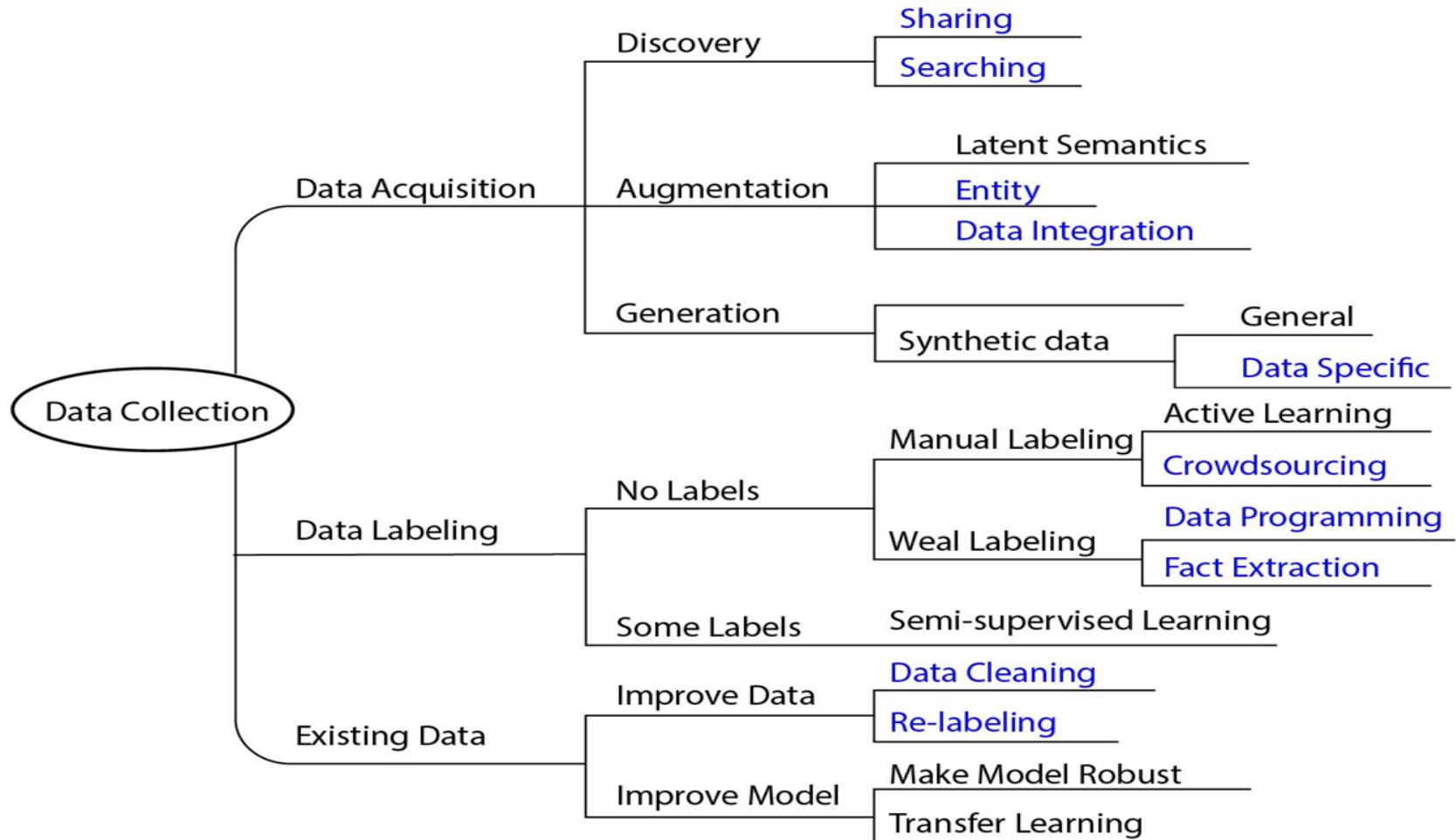


# Basic Components of Data Acquisition System

## Data Acquisition System Components



# Basic Components of Data Acquisition System



# Data Ingestion

---

- Data ingestion is the transportation of data from assorted sources to a storage medium where it can be accessed, used, and analyzed by an organization.
- The destination is typically a data warehouse, data mart, database, or a document store.
- Sources may be almost anything — including SaaS data, in-house apps, databases, spreadsheets, or even information scraped from the internet.

# Data Ingestion

- In a broader sense, data ingestion can be understood as a directed dataflow between two or more systems that result in a smooth, and independent, operation (a definition which already implies some independence or automation).
- Ingestion can occur in real-time, as soon as the source produces it, or in batches, when data is input in specific chunks at set periods.
- Generally, three steps occur within data ingestion:
- Data extraction – retrieving data from sources
- Data transformation – validating, cleaning, and normalizing data to ensure accuracy and reliability (sometimes known as trustworthiness)
- Data loading – routing or placing the data in its correct silo or database for analysis



# Reasons to automate data ingestion

---

- Automating data ingestion improves time-to-market goals.
- Automating data ingestion increases scalability.
- Automating data ingestion refocuses on necessary work.
- Automating data ingestion mitigates risk.

# Data Ingestion Tools

---

- Apache NiFi (a.k.a. Hortonworks DataFlow)
- Stream Sets Data Collector (SDC)
- Gobblin
- Sqoop
- Flume
- Kafka
- Airflow
- Spark
- petl

# Data Ingestion Tools

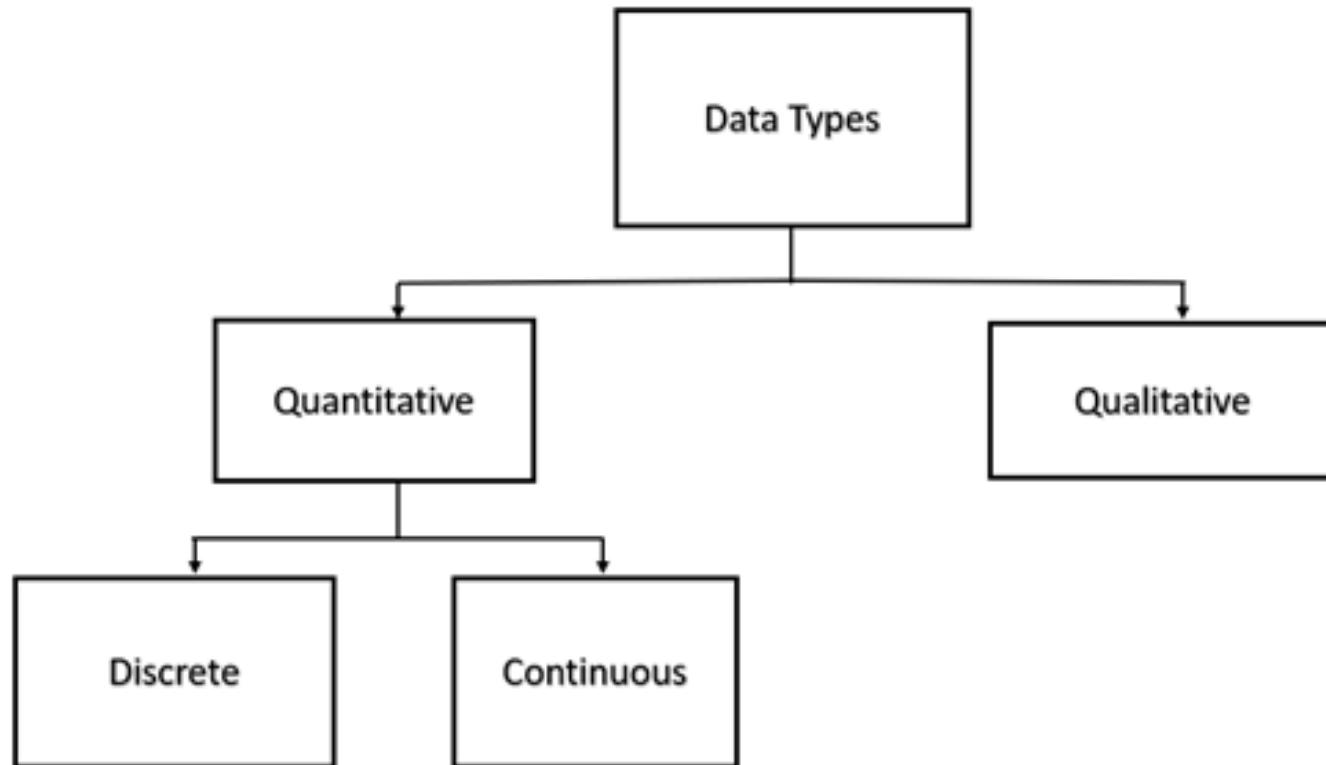
---

- Panoply
- Pandas
- Bubbles
- Bonobo
- Luigi
- etlalchemy



## Data Types In Machine Learning

# Data Types in Machine Learning



# Data Types in Machine Learning

---

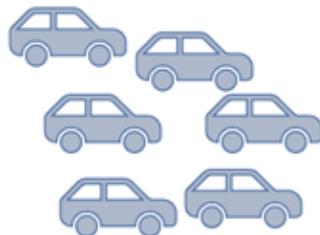
- Quantitative Data Type: –
  - This type of data type consists of numerical values.  
Anything which is measured by numbers.
  - ex., Profit, Quantity Sold, Height, Weight, Temperature, etc.

# Data Types in Machine Learning

- A.) Discrete Data Type: –
- The numeric data which have discrete values or whole numbers.
- This type of variable value if expressed in decimal format will have no proper meaning.
- Their values can be counted.
- ex: – No. Of Cars You Have, No. Of Marbles In Containers, Students In A Class, etc.



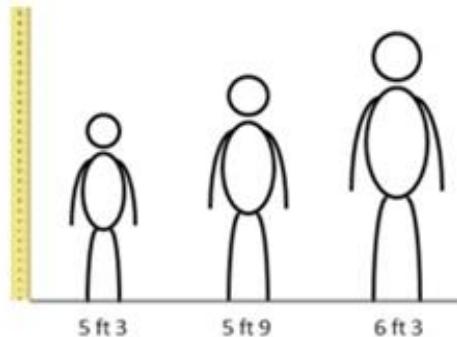
No. of Laptops



No. of Cars

# Data Types in Machine Learning

- B.) Continuous Data Type: –
- The numerical measures which can take the value within A certain range.
- This type of variable value if expressed in decimal format has true meaning.
- Their values can not be counted but measured. The value can be infinite
- ex: – Height, Weight, Time, Area, Distance, Measurement Of Rainfall, etc.



Height



Time

# Data Types in Machine Learning

---

- Qualitative Data Type: –
- These are the data types that cannot be expressed in numbers.
- This describes categories or groups and is hence known as the categorical data type.

# Data Types in Machine Learning

- Structured Data:
- This type of data is either number or words.
- This can take numerical values, but mathematical operations cannot be performed on it.
- This type of data is expressed in tabular format.
- ex) Sunny=1, Cloudy=2, Windy=3 Or Binary Form Data Like 0 Or 1, Good Or Bad, Etc.

ID	Name	Age	Degree
1	John	18	B.Sc.
2	David	31	Ph.D.
3	Robert	51	Ph.D.
4	Rick	26	M.Sc.
5	Michael	19	B.Sc.

# Data Types in Machine Learning

- Unstructured Data:
- This type of data does not have the proper format and therefore known as unstructured data.
- This comprises textual data, sounds, images, videos, etc.

			
Text files and documents	Server, website and application logs	Sensor data	Images
			
Video files	Audio files	Emails	Social media data



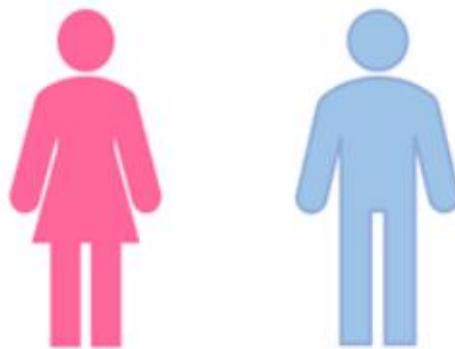
# Data Types in Machine Learning

---

- This, there are also other types refer as data types preliminaries or data measures:-
- Nominal
- Ordinal
- Interval
- Ratio

# Data Types in Machine Learning

- Nominal Data Type:
- This is in use to express names or labels which are not order or measurable.
- ex., Male or female (gender), race, country, etc.



*Fig: Gender (Female, Male),  
An Example Of Nominal  
Data Type*

# Data Types in Machine Learning

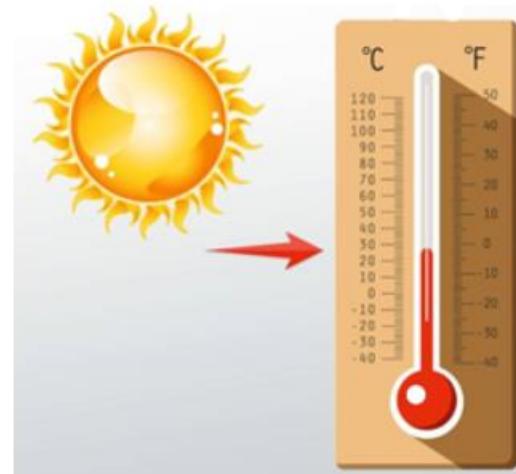
- Ordinal Data Type:
- This is also A categorical data type like nominal data but has some natural ordering associated with it.
- Ex., Likert rating scale, shirt sizes, ranks, grades, etc.,



*Fig: Rating (Good, Average, Poor), An Example Of Ordinal Data Type*

# Data Types in Machine Learning

- Interval Data Type:
- This is numeric data which has proper order and the exact zero means the true absence of A value attached.
- Here zero means not A complete absence but has some value. This is the local scale.
- ex., Temperature measured in degree Celsius, time, sat score, credit score, PH, etc. Difference between values is familiar. In this case, there is no absolute zero. Absolute

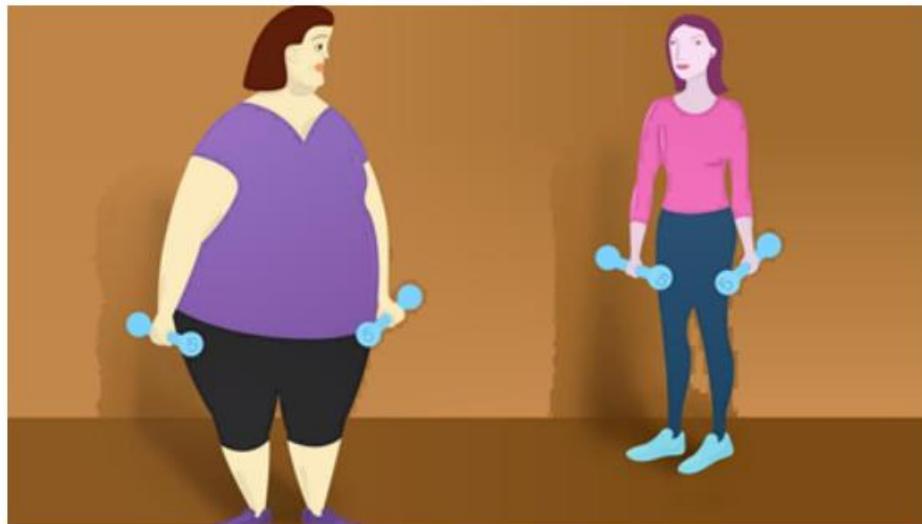


*Fig: Temperature, An Example Of*

*Interval Data Type*

# Data Types in Machine Learning

- Ratio Data Type:
- This quantitative data type is the same as the interval data type but has the absolute zero. Here zero means complete absence and the scale starts from zero. This is the global scale.
- ex., Temperature in kelvin, height, weight, etc.



*Fig: Weight, An Example Of Ratio Data Type*

# Difference between Structured, Semi-structured and Unstructured data



PROPERTIES	STRUCTURED DATA	SEMI-STRUCTURED DATA	UNSTRUCTURED DATA
Technology	It is based on Relational database table	It is based on XML/RDF	It is based on character and binary data
Transaction management	Matured transaction and various concurrency technique	Transaction is adapted from DBMS not matured	No transaction management and no concurrency
Version management	Versioning over tuples, row, tables	Versioning over tuples or graph is possible	Versioned as whole
Flexibility	It is schema dependent and less flexible	It is more flexible than structured data but less than flexible than unstructured data	It is very flexible and there is absence of schema
Scalability	It is very difficult to scale DB schema	Its scaling is simpler than structured data	It is very scalable
Robustness	Very robust	New technology, not very spread	—
Query performance	Structured query allow complex joining	Queries over anonymous nodes are possible	Only textual query are possible

# Data from Various Sources

---

- CSV file
- Flat File (tab, space, or any other separator)
- Text File (In a single file — reading data all at once)
- ZIP file
- Multiple Text Files (Data is split over multiple text files)
- Download File from Internet (File hosted on a server)
- Webpage (scraping)
- APIs (JSON)
- Text File (Reading data line by line)
- RDBMS (SQL Tables)

# Common Data Platforms

---

- Inconsistent column names
- Missing data
- Outliers
- Duplicate rows
- Untidy
- Need to process columns
- Column types can signal unexpected data values



# Exploratory Data Analysis (EDA)

- Exploratory Data Analysis (EDA) refers to the method of studying and exploring record sets.
- To apprehend their predominant traits, discover patterns, locate outliers, and identify relationships between variables.
- EDA is normally carried out as a preliminary step before undertaking extra formal statistical analyses or modeling.



# The Foremost Goals of EDA

---

- Data Cleaning: EDA involves examining the information for errors, lacking values, and inconsistencies.
- It includes techniques including records imputation, managing missing statistics, and figuring out and getting rid of outliers.



# The Foremost Goals of EDA

- Descriptive Statistics: EDA utilizes precise records to recognize the important tendency, variability, and distribution of variables.
- Measures like suggest, median, mode, preferred deviation, range, and percentiles are usually used.

# The Foremost Goals of EDA

---

- Data Visualization: EDA employs visual techniques to represent the statistics graphically.
- Visualizations consisting of histograms, box plots, scatter plots, line plots, heatmaps, and bar charts assist in identifying styles, trends, and relationships within the facts.



# The Foremost Goals of EDA

- Feature Engineering: EDA allows for the exploration of various variables and their adjustments to create new functions or derive meaningful insights.
- Feature engineering can contain scaling, normalization, binning, encoding express variables, and creating interplay or derived variables.



## The Foremost Goals of EDA

---

- Correlation and Relationships: EDA allows discover relationships and dependencies between variables.
- Techniques such as correlation analysis, scatter plots, and pass-tabulations offer insights into the power and direction of relationships between variables.

# The Foremost Goals of EDA

---

- Data Segmentation: EDA can contain dividing the information into significant segments based totally on sure standards or traits.
- This segmentation allows advantage insights into unique subgroups inside the information and might cause extra focused analysis.

# The Foremost Goals of EDA

---

- Hypothesis Generation: EDA aids in generating hypotheses or studies questions based totally on the preliminary exploration of the data.
- It facilitates form the inspiration for in addition evaluation and model building.

# The Foremost Goals of EDA

---

- Data Quality Assessment: EDA permits for assessing the nice and reliability of the information.
- It involves checking for records integrity, consistency, and accuracy to make certain the information is suitable for analysis.

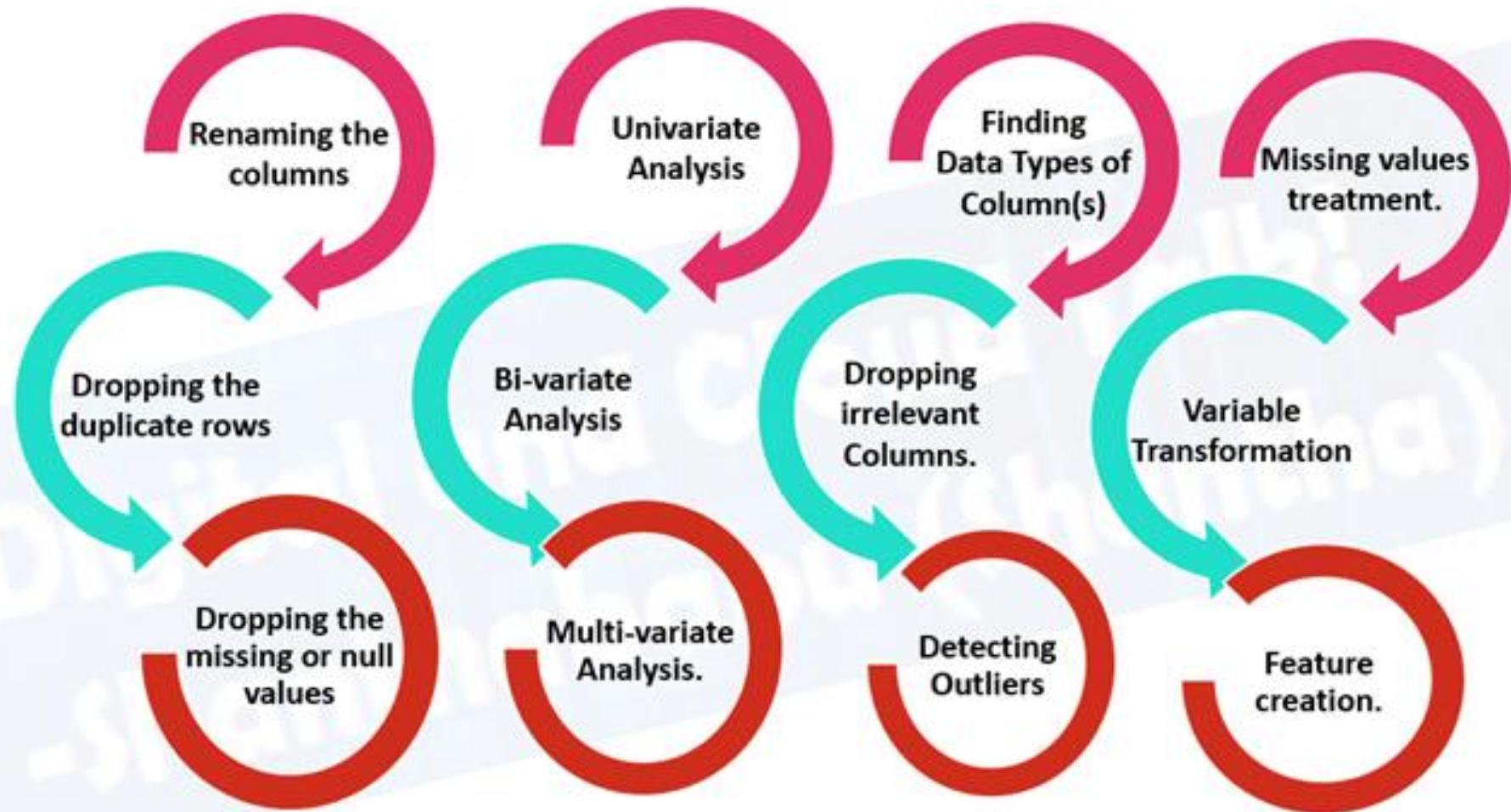
## Outcome of EDA

---

- Understanding the given dataset and helps clean up the given dataset.
- It gives you a clear picture of the features and the relationships between them.
- Providing guidelines for essential variables and leaving behind/removing non-essential variables.
- Handling Missing values or human error.
- Identifying outliers.
- EDA process would be maximizing insights of a dataset.



# Types of EDA



# Types of EDA

---

- Univariate Analysis: This sort of evaluation makes a specialty of analyzing character variables inside the records set.
- It involves summarizing and visualizing unmarried variable at a time to understand its distribution, relevant tendency, unfold, and different applicable records.
- Techniques like histograms, field plots, bar charts, and precis information are generally used in univariate analysis.

# Types of EDA

---

- Bivariate Analysis: Bivariate evaluation involves exploring the connection between variables.
- It enables find associations, correlations, and dependencies between pairs of variables.
- Scatter plots, line plots, correlation matrices, and move-tabulation are generally used strategies in bivariate analysis.

## Types of EDA

---

- Multivariate Analysis: Multivariate analysis extends bivariate evaluation to encompass greater than variables.
- It ambitions to apprehend the complex interactions and dependencies among more than one variables in a records set.
- Techniques inclusive of heatmaps, parallel coordinates, aspect analysis, and primary component analysis (PCA) are used for multivariate analysis.

## Types of EDA

---

- Time Series Analysis: This type of analysis is mainly applied to statistics sets that have a temporal component.
- Time collection evaluation entails inspecting and modeling styles, traits, and seasonality inside the statistics through the years.
- Techniques like line plots, autocorrelation analysis, transferring averages, and ARIMA (AutoRegressive Integrated Moving Average) fashions are generally utilized in time series analysis.

## Types of EDA

---

- Missing Data Analysis: Missing information is a not unusual issue in datasets, and it may impact the reliability and validity of the evaluation.
- Missing statistics analysis includes figuring out missing values, know-how the patterns of missingness, and using suitable techniques to deal with missing data.
- Techniques along with lacking facts styles, imputation strategies, and sensitivity evaluation are employed in lacking facts evaluation.

## Types of EDA

---

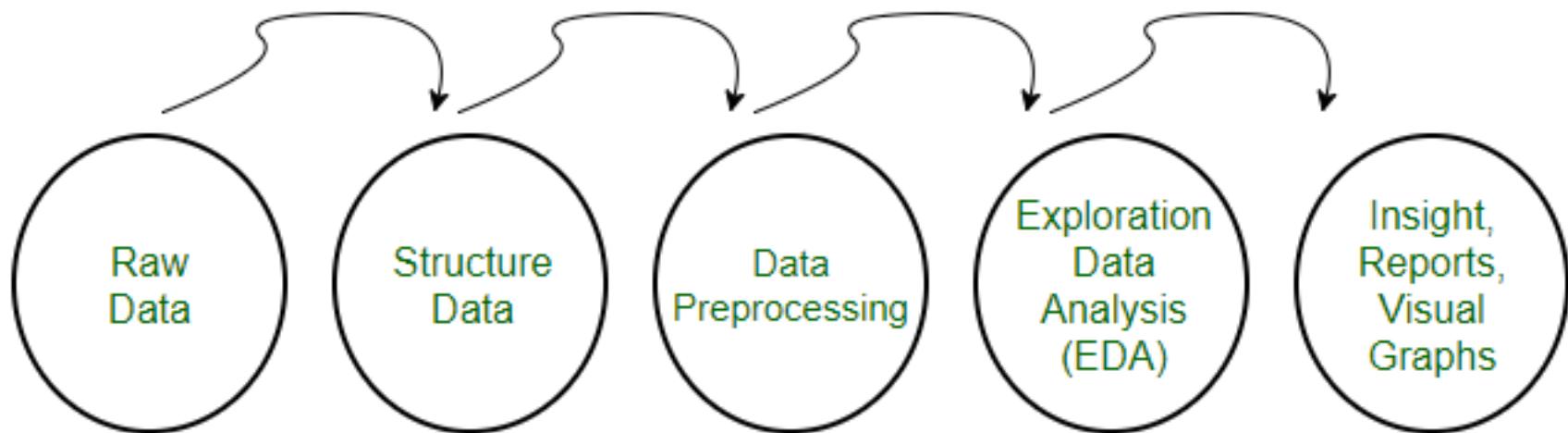
- Outlier Analysis: Outliers are statistics factors that drastically deviate from the general sample of the facts.
- Outlier analysis includes identifying and knowledge the presence of outliers, their capability reasons, and their impact at the analysis.
- Techniques along with box plots, scatter plots, z-rankings, and clustering algorithms are used for outlier evaluation.

## Types of EDA

---

- Data Visualization: Data visualization is a critical factor of EDA that entails creating visible representations of the statistics to facilitate understanding and exploration.
- Various visualization techniques, inclusive of bar charts, histograms, scatter plots, line plots, heatmaps, and interactive dashboards, are used to represent exclusive kinds of statistics.

# Data preprocessing



# Steps in Data Preprocessing

---

- Step 1: Import the necessary libraries
- Step 2: Load the dataset
- Step 3: Statistical Analysis
- Step 4: Check the outliers
- Step 5: Correlation
- Step 6: Separate independent features and Target Variables
- Step 7: Normalization or Standardization

# Data Quality



Measuring  
Knowing  
Improving





# Data Quality assessment

- Serviceability
- Meta-Data
- Reliability
- Relevance
- Structure
- Richness of Content

Dimensions

What is DQA

- A scientific way to determine if data is usable
- If usable, to what degree?
- Level of Standardization
- Data Cleansing Strategy

DQA

Process

- Extract
- Prepare
- Discover
- Evaluate
- Report

Objectives

- Completeness
- Correctness
- Conformance

# Our 5-step DQA Process

## Extract

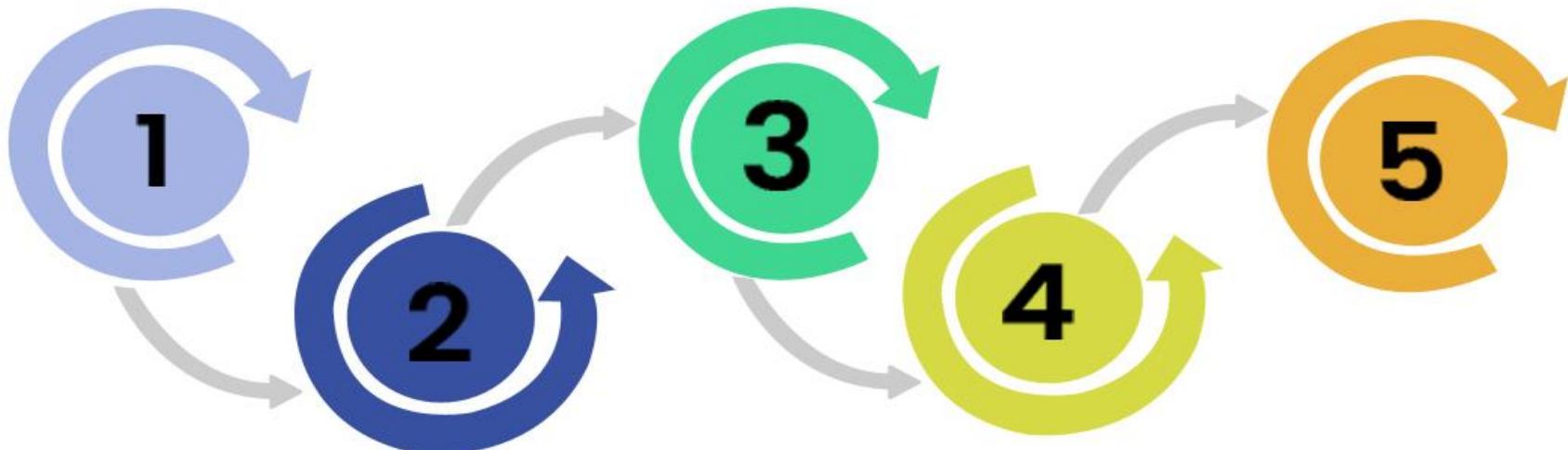
- Identify all data sources
- Identify data size – volume, duration, no. of samples
- Convert / Format data

## Discover

- Perform EDA (Exploratory Data Analysis)
- Data visualization
- Identify external data that could be useful

## Report

- Comprehensiveness Quotient
- Accuracy
- Richness of content



## Prepare

- Clean the data – Outliers and Missing values
- Feature Engineering
- Feature scaling – Standardization & Normalization

## Evaluate

- Feature importance using ML models
- Feature correlations and optimization
- Measure performance of predictive models



DATA SCIENCE

# Benefits of DQA

## Benefits of DQA

### Data Comprehensiveness Scores

- Missing Values
- Accuracy
- Conformance
- Integrity Check



### Dark Data Identification

Detect parts of data that can't be used effectively, often because of data quality problems



### Improved Predictive & Prescriptive Capability

Improve the ability of the data to deliver prescriptive insights and predictive outcomes



### Campaign Effectiveness

Improve efficacy of the business use case through better data quality



# Feature Scaling

- Feature scaling is a method used to normalize the range of independent variables or features of data.
- In data processing, it is also known as data normalization and is generally performed during the data preprocessing step.
- Just to give you an example — if you have multiple independent variables like age, salary, and height; With their range as (18–100 Years), (25,000–75,000 Euros), and (1–2 Meters) respectively.
- Feature scaling would help them all to be in the same range, for example- centered around 0 or in the range (0,1) depending on the scaling technique.

# Why Feature Scaling

#	Emp	Age	Salary
1	Emp1	44	73000
2	Emp2	27	47000
3	Emp3	30	53000
4	Emp4	38	62000
5	Emp5	40	57000
6	Emp6	35	53000
7	Emp7	48	78000

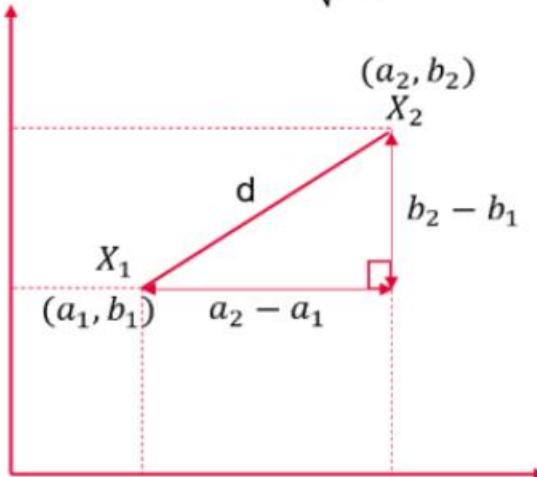
Age Range: 27-48

Salary Range: 47000-78000



# Why Feature Scaling

$$\text{Euclidean Distance } d = \sqrt{(a_2 - a_1)^2 + (b_2 - b_1)^2}$$



- $X_1$  and  $X_2$  are two data points means two different rows in data set. And Data set is two dimensional feature space
- For three dimensional feature space suppose two data points are  $X_1(a_1, b_1, c_1)$  and  $X_2(a_2, b_2, c_2)$  then distance can be calculated as below:

$$d = \sqrt{(a_2 - a_1)^2 + (b_2 - b_1)^2 + (c_2 - c_1)^2}$$

- In case of multi-dimensional vector space, More generic formula will be

$$\text{Euclidean Distance } d = \sum_{i=1}^k \sqrt{(X1_i - X2_i)^2}$$



# Why Feature Scaling

#	Emp	Age	Salary
1	Emp1	44	73000
2	Emp2	27	47000
3	Emp3	30	53000
4	Emp4	38	62000
5	Emp5	40	57000
6	Emp6	35	53000
7	Emp7	48	78000

$$\text{Distance between Emp2 and Emp1} = \sqrt{(27 - 44)^2 + (47000 - 73000)^2} = 31.06$$

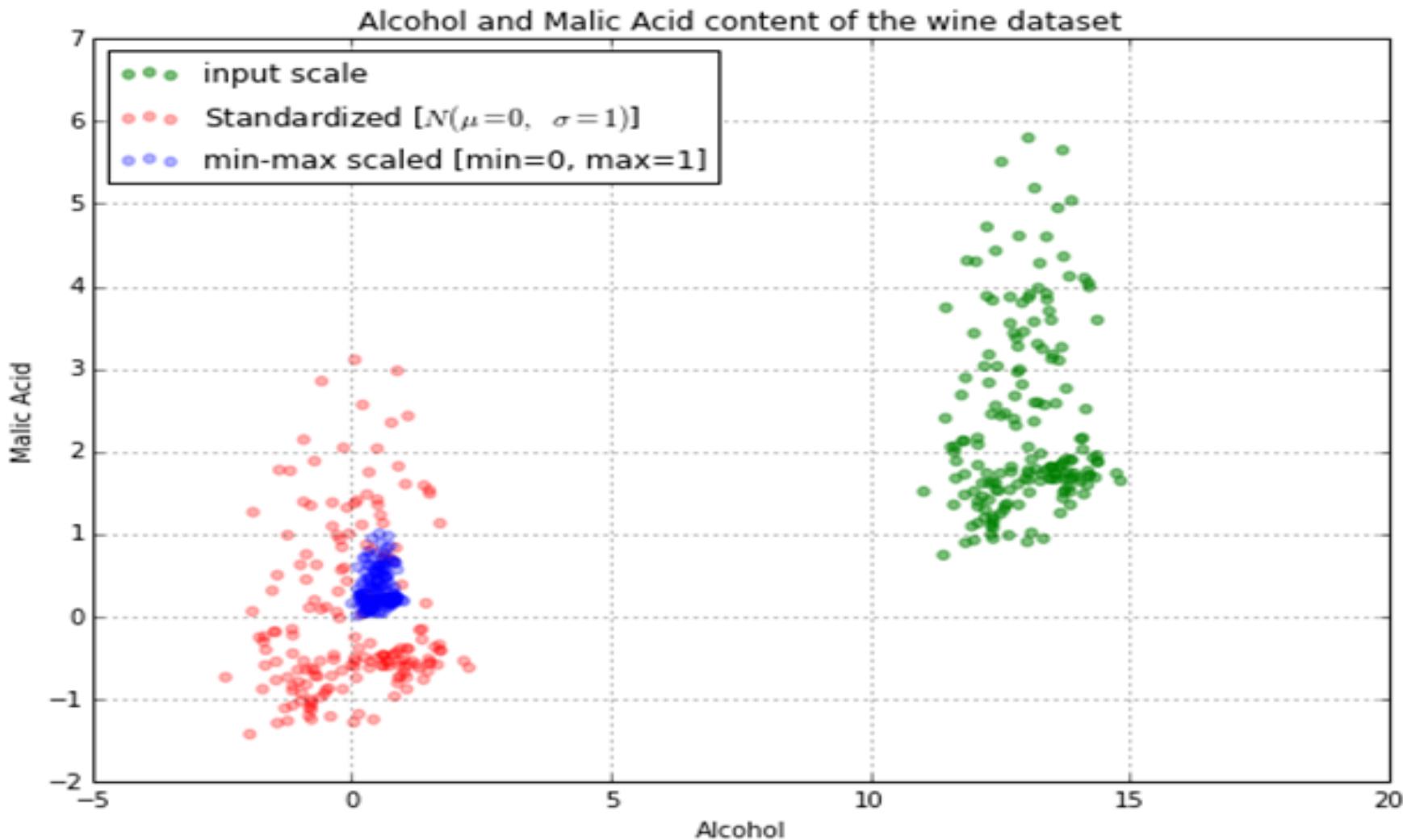
$$\text{Distance between Emp2 and Emp3} = \sqrt{(30 - 27)^2 + (53000 - 47000)^2} = 6.70$$

# Feature Scaling

---

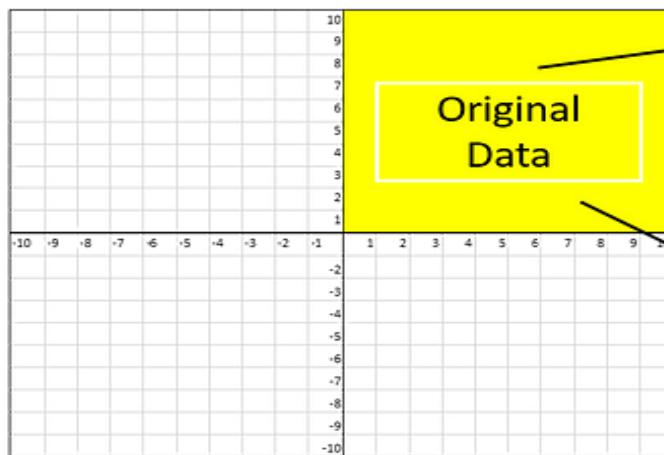
- To visualize the above, let us take an example of the independent variables of alcohol and Malic Acid content in the wine dataset from the “Wine Dataset” that is deposited on the UCI machine learning repository.
- Below you can see the impact of the two most common scaling techniques (Normalization and Standardization) on the dataset.

# Feature Scaling

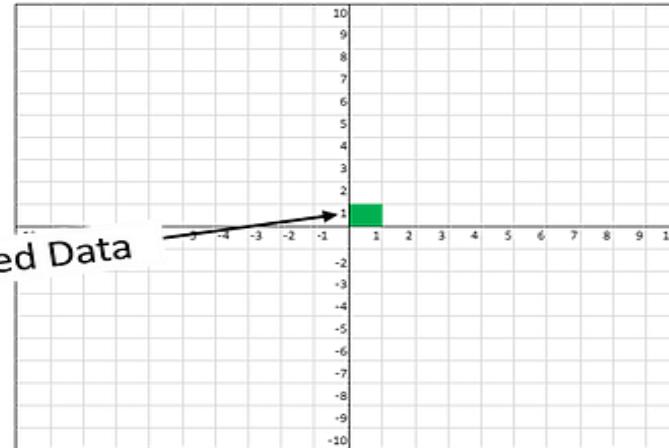




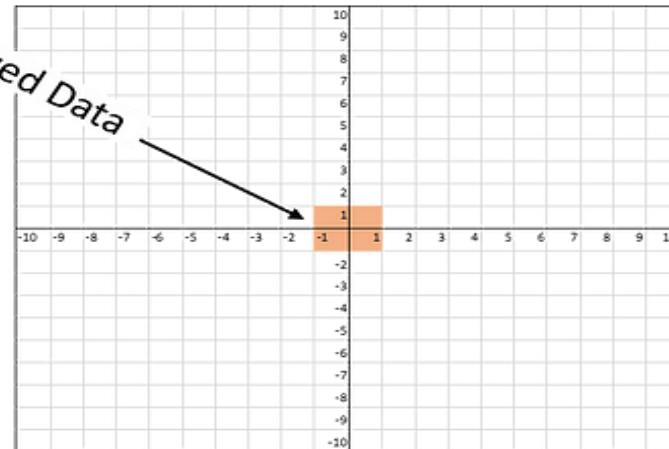
# Methods for Scaling



Original Data



Normalized Data



Standardized Data

# Methods for Scaling

## Normalization

Also known as min-max scaling or min-max normalization, it is the simplest method and consists of rescaling the range of features to scale the range in [0, 1]. The general formula for normalization is given as:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Here,  $\max(x)$  and  $\min(x)$  are the maximum and the minimum values of the feature respectively.

We can also do a normalization over different intervals, e.g. choosing to have the variable laying in any  $[a, b]$  interval,  $a$  and  $b$  being real numbers. To rescale a range between an arbitrary set of values  $[a, b]$ , the formula becomes:

# Methods for Scaling

---

We can also do a normalization over different intervals, e.g. choosing to have the variable laying in any  $[a, b]$  interval,  $a$  and  $b$  being real numbers. To rescale a range between an arbitrary set of values  $[a, b]$ , the formula becomes:

$$x' = a + \frac{(x - \min(x))(b - a)}{\max(x) - \min(x)}$$



# Methods for Scaling

$$x_{\text{norm}} = \frac{x - \min(x)}{\max(x) - \min(x)}$$

#	Emp	Age	Salary
1	Emp1	44	73000
2	Emp2	27	47000
3	Emp3	30	53000
4	Emp4	38	62000
5	Emp5	40	57000
6	Emp6	35	53000
7	Emp7	48	78000

Normalization

Age	Normalized Age	Salary	Normalized Salary
44	0.80952381	73000	0.838709677
27	0	47000	0
30	0.142857143	53000	0.193548387
38	0.523809524	62000	0.483870968
40	0.619047619	57000	0.322580645
35	0.380952381	53000	0.193548387
48	1	78000	1

Range 0-1

Range 0-1

How to calculate Normalized value?

X = 35, min = 27, max = 48 for column Age.

$$X_{\text{norm}}(\text{for } 35) = \frac{35-27}{48-27} = 0.3809$$

## After Normalization

$$\text{Distance between Emp2 and Emp1} = \sqrt{(0 - .80)^2 + (0 - .83)^2} = 1.15$$

$$\text{Distance between Emp2 and Emp3} = \sqrt{(.14 - 0)^2 + (.19 - 0)^2} = 0.23$$

Comparison will be  
more significant

# Methods for Scaling

## Standardization

Feature standardization makes the values of each feature in the data have zero mean and unit variance. The general method of calculation is to determine the distribution mean and standard deviation for each feature and calculate the new data point by the following formula:

$$x' = \frac{x - \bar{x}}{\sigma}$$

Here,  $\sigma$  is the standard deviation of the feature vector, and  $\bar{x}$  is the average of the feature vector.



# Methods for Scaling

$$x_{\text{stand}} = \frac{x - \text{mean}(x)}{\text{standard deviation } (x)}$$

#	Emp	Age	Salary
1	Emp1	44	73000
2	Emp2	27	47000
3	Emp3	30	53000
4	Emp4	38	62000
5	Emp5	40	57000
6	Emp6	35	53000
7	Emp7	48	78000

Mean =  
37.42857  
Std. Dev. =  
6.883876

Mean =  
60428.5714  
Std. Dev. =  
10499.7570

Standardization

How to calculate Standardized value?  
 $X = 35$ , mean = 37.42, Std. Dev. = 6.88  
for column Age.  
 $X_{\text{std}}(\text{for } 35) = \frac{35 - 37.42}{6.88} = -0.3527$

Age	Standardized Age	Salary	Standardized Salary
44	0.954611636	73000	1.197306616
27	-1.514927162	47000	-1.278941158
30	-1.079126198	53000	-0.707499364
38	0.083009708	62000	0.149663327
40	0.373543684	57000	-0.326538168
35	-0.352791257	53000	-0.707499364
48	1.535679589	78000	1.673508111

Mean = 0  
Std. dev. = 1

Mean = 0  
Std. dev. = 1



# Methods for Scaling

## After Standardization

$$\text{Distance between Emp2 and Emp1} = \sqrt{(-1.51 - 0.95)^2 + (-1.27 - 1.19)^2} = 3.47$$

$$\text{Distance between Emp2 and Emp3} = \sqrt{(-1.07 + 1.51)^2 + (-0.70 + 1.27)^2} = 0.71$$

Comparison will be  
more significant

# Normalization vs Standardization

- If you have outliers in your feature (column), normalizing your data will scale most of the data to a small interval, which means all features will have the same scale and hence it will not handle outliers well.
- Standardization is more robust to outliers, and in many cases, it is preferable over Max-Min Normalization.
- Normalization is good to use when your data does not follow a Normal distribution.
- This can be useful in algorithms that do not assume any distribution of the data like K-Nearest Neighbors and Neural Networks.
- Standardization, on the other hand, can be helpful in cases where the data follows a Normal distribution.
- However, this does not have to be necessarily true.
- Also, unlike normalization, standardization does not have a bounding range.
- So, even if you have outliers in your data, they will not be affected by standardization.

# Types of Feature Scaling

- Standardization:
- Standardization is also known as z-score Normalization.
- In standardization, features are scaled to have zero-mean and one-standard-deviation.
- It means after standardization features will have mean = 0 and standard deviation = 1.
  - Standard Scaler
- Normalization:
- Normalization is also known as min-max normalization or min-max scaling.
- Normalization re-scales values in the range of 0-1
  - Min Max Scaling
  - Mean Normalization
  - Max Absolute Scaling
  - Robust Scaling etc.

# Which Algorithms Need Feature Scaling

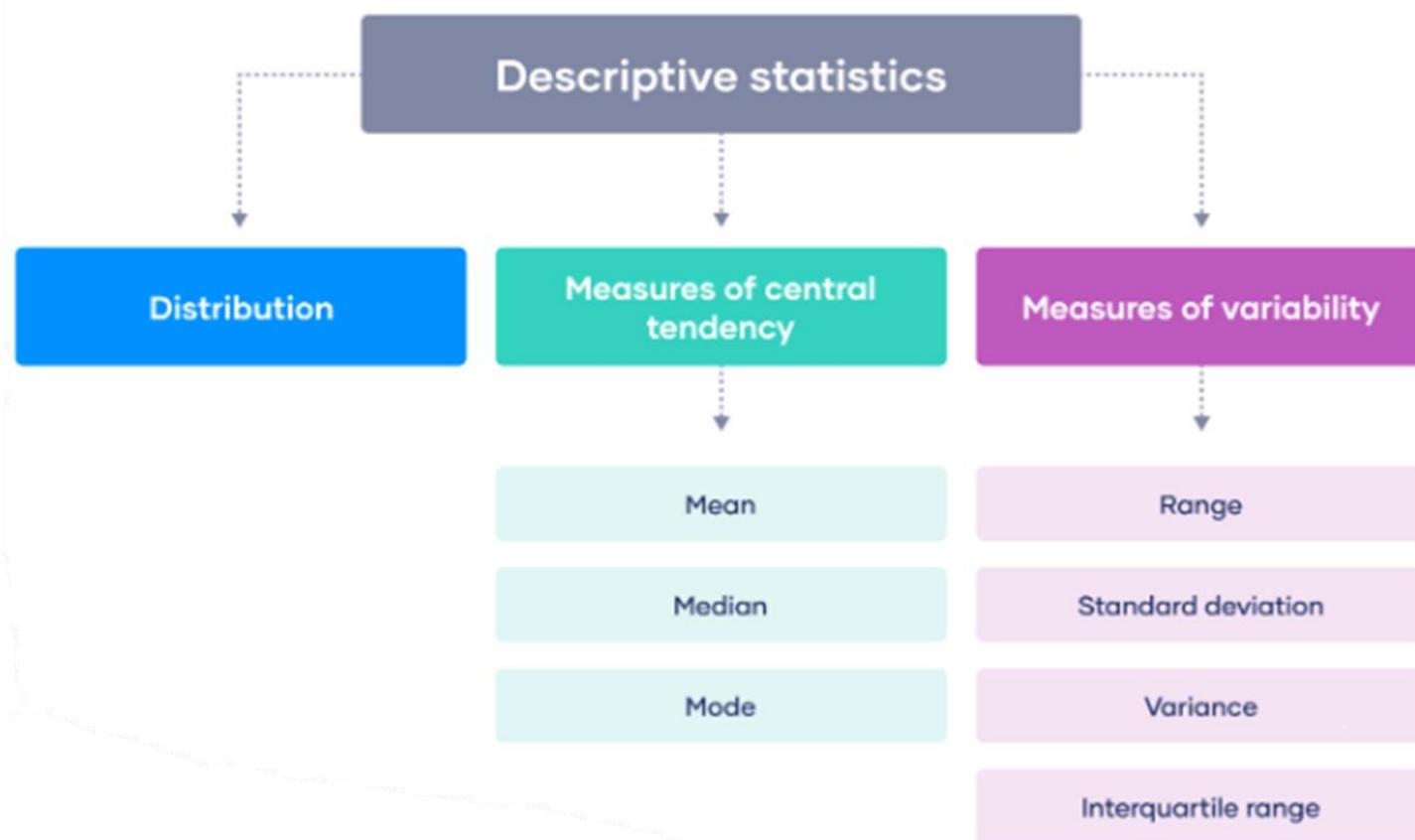
Algorithms	Reasons for applying feature scaling
K-means	Use Euclidean Distance measure
K-nearest neighbours	Measure the distance between pairs of samples and these distances are influenced by the measurement units
Principal Component Analysis (PCA)	Get the features with maximum variance
Artificial Neural Network	Apply gradient descent
Gradient Descent	Theta calculation becomes faster after feature scaling and the learning rate in the update equation of Stochastic gradient descent is the same for every parameter.

# Descriptive Statistics

---

- Descriptive statistics summarize and organize characteristics of a data set.
- A data set is a collection of responses or observations from a sample or entire population.
- In quantitative research, after collecting data, the first step of statistical analysis is to describe characteristics of the responses, such as the average of one variable (e.g., age), or the relation between two variables (e.g., age and creativity).
- The next step is inferential statistics, which help you decide whether your data confirms or refutes your hypothesis and whether it is generalizable to a larger population.

# Types of descriptive statistics



# Frequency distribution

A data set is made up of a distribution of values, or scores. In tables or graphs, you can summarize the frequency of every possible value of a variable in numbers or percentages. This is called a **frequency distribution**.

## Simple frequency distribution table

## Grouped frequency distribution table

For the variable of gender, you list all possible answers on the left hand column. You count the number or percentage of responses for each answer and display it on the right hand column.

Gender	Number
Male	182
Female	235
Other	27

From this table, you can see that more women than men or people with another gender identity took part in the study.

# Measures of central tendency

Measures of central tendency estimate the center, or average, of a data set. The mean, median and mode are 3 ways of finding the average.

Here we will demonstrate how to calculate the mean, median, and mode using the first 6 responses of our survey.

## Mean      Median      Mode

The **mean**, or  $M$ , is the most commonly used method for finding the average.

To find the mean, simply add up all response values and divide the sum by the total number of responses. The total number of responses or observations is called  $N$ .

### Mean number of library visits

<b>Data set</b>	15, 3, 12, 0, 24, 3
<b>Sum of all values</b>	$15 + 3 + 12 + 0 + 24 + 3 = 57$
<b>Total number of responses</b>	$N = 6$
<b>Mean</b>	Divide the sum of values by $N$ to find $M$ : $57/6 = 9.5$

# Measures of central tendency

Measures of central tendency estimate the center, or average, of a data set. The mean, median and mode are 3 ways of finding the average.

Here we will demonstrate how to calculate the mean, median, and mode using the first 6 responses of our survey.

## Mean      Median      Mode

The **mean**, or  $M$ , is the most commonly used method for finding the average.

To find the mean, simply add up all response values and divide the sum by the total number of responses. The total number of responses or observations is called  $N$ .

### Mean number of library visits

<b>Data set</b>	15, 3, 12, 0, 24, 3
<b>Sum of all values</b>	$15 + 3 + 12 + 0 + 24 + 3 = 57$
<b>Total number of responses</b>	$N = 6$
<b>Mean</b>	Divide the sum of values by $N$ to find $M$ : $57/6 = 9.5$

# Measures of variability

Measures of variability give you a sense of how spread out the response values are. The range, standard deviation and variance each reflect different aspects of spread.

## Range

The range gives you an idea of how far apart the most extreme response scores are. To [find the range](#), simply subtract the lowest value from the highest value.

Range of visits to the library in the past year

**Ordered data set:** 0, 3, 3, 12, 15, 24

**Range:**  $24 - 0 = 24$

# Standard deviation

- The standard deviation (s or SD) is the average amount of variability in our dataset.
- It tells us, on average, how far each score lies from the mean.
- The larger the standard deviation, the more variable the data set is.
- There are six steps for finding the standard deviation:
  - List each score and find their mean.
  - Subtract the mean from each score to get the deviation from the mean.
  - Square each of these deviations.
  - Add up all of the squared deviations.
  - Divide the sum of the squared deviations by  $N - 1$ .
  - Find the square root of the number you found.

# Standard deviation

Standard deviations of visits to the library in the past year

In the table below, you complete **Steps 1 through 4.**

Raw data	Deviation from mean	Squared deviation
15	$15 - 9.5 = 5.5$	30.25
3	$3 - 9.5 = -6.5$	42.25
12	$12 - 9.5 = 2.5$	6.25
0	$0 - 9.5 = -9.5$	90.25
24	$24 - 9.5 = 14.5$	210.25
3	$3 - 9.5 = -6.5$	42.25
$M = 9.5$	Sum = 0	Sum of squares = 421.5

**Step 5:**  $421.5/5 = 84.3$

**Step 6:**  $\sqrt{84.3} = 9.18$

From learning that  $s = 9.18$ , you can say that on average, each score deviates from the mean by 9.18 points.

# Variance

## Variance

The variance is the average of squared deviations from the mean. Variance reflects the degree of spread in the data set. The more spread the data, the larger the variance is in relation to the mean.

To find the variance, simply square the standard deviation. The symbol for variance is  $s^2$ .

Variance of visits to the library in the past year

**Data set:** 15, 3, 12, 0, 24, 3

$$s = 9.18$$

$$s^2 = 84.3$$

# Univariate descriptive statistics

Univariate descriptive statistics focus on only one variable at a time. It's important to examine data from each variable separately using multiple measures of distribution, central tendency and spread. Programs like SPSS and Excel can be used to easily calculate these.

Visits to the library	
<i>N</i>	6
Mean	9.5
Median	7.5
Mode	3
Standard deviation	9.18
Variance	84.3
Range	24

If you were to only consider the mean as a measure of central tendency, your impression of the "middle" of the data set can be **skewed** by outliers, unlike the median or mode.

Likewise, while the range is sensitive to **outliers**, you should also consider the standard deviation and variance to get easily comparable measures of spread.

# Bivariate descriptive statistics

- If we've collected data on more than one variable, we can use bivariate or multivariate descriptive statistics to explore whether there are relationships between them.
- In bivariate analysis, we simultaneously study the frequency and variability of two variables to see if they vary together.
- We can also compare the central tendency of the two variables before performing further statistical tests.
- Multivariate analysis is the same as bivariate analysis but with more than two variables.

# Contingency table

In a contingency table, each cell represents the intersection of two variables. Usually, an **independent variable** (e.g., gender) appears along the vertical axis and a dependent one appears along the horizontal axis (e.g., activities). You read “across” the table to see how the independent and **dependent variables** relate to each other.

		Number of visits to the library in the past year				
		0-4	5-8	9-12	13-16	17+
Group	Children	32	68	37	23	22
	Adults	36	48	43	83	25

Interpreting a contingency table is easier when the raw data is converted to percentages. Percentages make each row comparable to the other by making it seem as if each group had only 100 observations or participants. When creating a percentage-based contingency table, you add the  $N$  for each independent variable on the end.

# Contingency table

Visits to the library in the past year (Percentages)

Group	0-4	5-8	9-12	13-16	17+	N
Children	18%	37%	20%	13%	12%	182
Adults	15%	20%	18%	35%	11%	235

From this table, it is more clear that similar proportions of children and adults go to the library over 17 times a year. Additionally, children most commonly went to the library between 5 and 8 times, while for adults, this number was between 13 and 16.

# What is Data Imputation

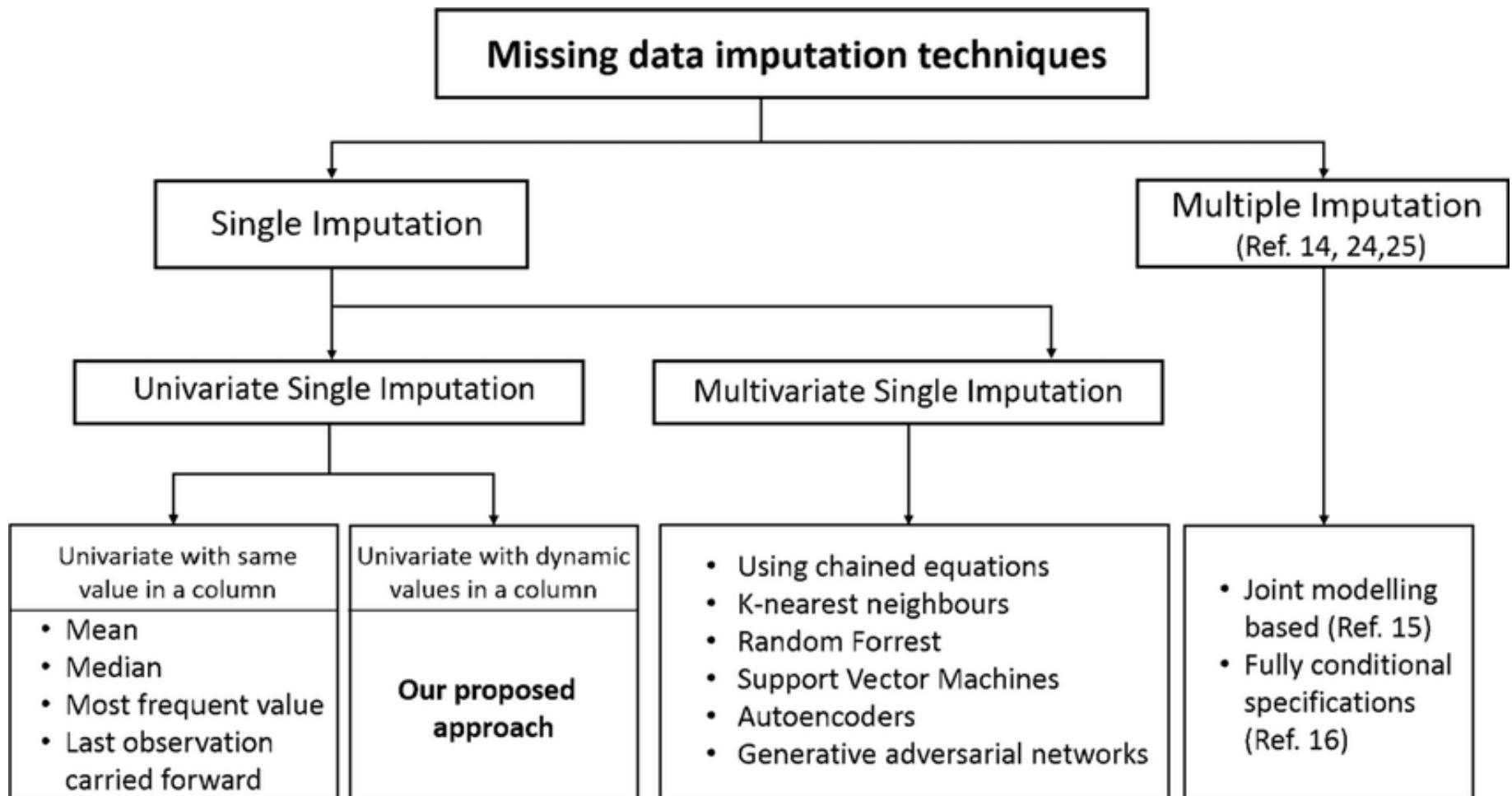
---

- Data imputation is a method for retaining the dataset's data and information by substituting missing data with a different value.
- These methods are employed because it would be impractical to remove data from a dataset each time.
- Additionally, doing so would substantially reduce the dataset's size, raising questions about bias and impairing analysis.

# Why Data Imputation

- Distorts Dataset: Large amounts of missing data can lead to anomalies in the variable distribution, which can change the relative importance of different categories in the dataset.
- Unable to work with machine learning-related Python libraries: When utilizing ML libraries (SkLearn is the most popular), mistakes may occur because there is no automatic handling of these missing data.
- Impacts on the Final Model: Missing data may lead to bias in the dataset, which could affect the final model's analysis.
- Desire to restore the entire dataset: This typically occurs when we don't want to lose any of the data in our dataset because all of it is crucial. Additionally, while the dataset is not very large, eliminating a portion of it could have a substantial effect on the final model.

# Missing Data Imputation Techniques





# Mark Missing Values

```
1 # load and summarize the dataset
2 from pandas import read_csv
3 # load the dataset
4 dataset = read_csv('pima-indians-diabetes.csv', header=None)
5 # summarize the dataset
6 print(dataset.describe())
```

Running this example produces the following output:

```
1          0         1         2   ...        6         7         8
2 count  768.000000  768.000000  768.000000  ...  768.000000  768.000000  768.000000
3 mean   3.845052   120.894531   69.105469   ...   0.471876   33.240885   0.348958
4 std    3.369578   31.972618   19.355807   ...   0.331329   11.760232   0.476951
5 min    0.000000   0.000000   0.000000   ...   0.078000   21.000000   0.000000
6 25%    1.000000   99.000000   62.000000   ...   0.243750   24.000000   0.000000
7 50%    3.000000  117.000000  72.000000   ...   0.372500   29.000000   0.000000
8 75%    6.000000  140.250000  80.000000   ...   0.626250   41.000000   1.000000
9 max   17.000000  199.000000 122.000000   ...   2.420000   81.000000   1.000000
10
11 [8 rows x 9 columns]
```

# Mark Missing Values

```
1 # load the dataset and review rows
2 from pandas import read_csv
3 # load the dataset
4 dataset = read_csv('pima-indians-diabetes.csv', header=None)
5 # print the first 20 rows of data
6 print(dataset.head(20))
```

# Mark Missing Values

Running the example, we can clearly see 0 values in the columns 2, 3, 4, and 5.

1	0	1	2	3	4	5	6	7	8	
2	0	6	148	72	35	0	33.6	0.627	50	1
3	1	1	85	66	29	0	26.6	0.351	31	0
4	2	8	183	64	0	0	23.3	0.672	32	1
5	3	1	89	66	23	94	28.1	0.167	21	0
6	4	0	137	40	35	168	43.1	2.288	33	1
7	5	5	116	74	0	0	25.6	0.201	30	0
8	6	3	78	50	32	88	31.0	0.248	26	1
9	7	10	115	0	0	0	35.3	0.134	29	0
10	8	2	197	70	45	543	30.5	0.158	53	1
11	9	8	125	96	0	0	0.0	0.232	54	1
12	10	4	110	92	0	0	37.6	0.191	30	0
13	11	10	168	74	0	0	38.0	0.537	34	1
14	12	10	139	80	0	0	27.1	1.441	57	0
15	13	1	189	60	23	846	30.1	0.398	59	1
16	14	5	166	72	19	175	25.8	0.587	51	1
17	15	7	100	0	0	0	30.0	0.484	32	1
18	16	0	118	84	47	230	45.8	0.551	31	1
19	17	7	107	74	0	0	29.6	0.254	31	1
20	18	1	103	30	38	83	43.3	0.183	33	0
21	19	1	115	70	30	96	34.6	0.529	32	1

# Count Column wise Missing Values

```
1 # example of summarizing the number of missing values for each variable
2 from pandas import read_csv
3 # load the dataset
4 dataset = read_csv('pima-indians-diabetes.csv', header=None)
5 # count the number of missing values for each column
6 num_missing = (dataset[[1,2,3,4,5]] == 0).sum()
7 # report the results
8 print(num_missing)
```

Running the example prints the following output:

1	1	5
2	2	35
3	3	227
4	4	374
5	5	11

# Replace Missing Values

Pandas provides the `fillna()` function for replacing missing values with a specific value.

For example, we can use `fillna()` to replace missing values with the mean value for each column, as follows:

```
1 # manually impute missing values with numpy
2 from pandas import read_csv
3 from numpy import nan
4 # load the dataset
5 dataset = read_csv('pima-indians-diabetes.csv', header=None)
6 # mark zero values as missing or NaN
7 dataset[[1,2,3,4,5]] = dataset[[1,2,3,4,5]].replace(0, nan)
8 # fill missing values with mean column values
9 dataset.fillna(dataset.mean(), inplace=True)
10 # count the number of NaN values in each column
11 print(dataset.isnull().sum())
```

Running the example provides a count of the number of missing values in each column, showing zero missing values.

```
1 0    0
2 1    0
3 2    0
4 3    0
5 4    0
6 5    0
7 6    0
8 7    0
9 8    0
10 dtype: int64
```

# Replace Missing Values

Pandas provides the `fillna()` function for replacing missing values with a specific value.

For example, we can use `fillna()` to replace missing values with the mean value for each column, as follows:

```
1 # manually impute missing values with numpy
2 from pandas import read_csv
3 from numpy import nan
4 # load the dataset
5 dataset = read_csv('pima-indians-diabetes.csv', header=None)
6 # mark zero values as missing or NaN
7 dataset[[1,2,3,4,5]] = dataset[[1,2,3,4,5]].replace(0, nan)
8 # fill missing values with mean column values
9 dataset.fillna(dataset.mean(), inplace=True)
10 # count the number of NaN values in each column
11 print(dataset.isnull().sum())
```

Running the example provides a count of the number of missing values in each column, showing zero missing values.

```
1 0    0
2 1    0
3 2    0
4 3    0
5 4    0
6 5    0
7 6    0
8 7    0
9 8    0
10 dtype: int64
```

# Replace Missing Values

The example below uses the SimpleImputer class to replace missing values with the mean of each column then prints the number of NaN values in the transformed matrix.

```
1 # example of imputing missing values using scikit-learn
2 from numpy import nan
3 from numpy import isnan
4 from pandas import read_csv
5 from sklearn.impute import SimpleImputer
6 # load the dataset
7 dataset = read_csv('pima-indians-diabetes.csv', header=None)
8 # mark zero values as missing or NaN
9 dataset[[1,2,3,4,5]] = dataset[[1,2,3,4,5]].replace(0, nan)
10 # retrieve the numpy array
11 values = dataset.values
12 # define the imputer
13 imputer = SimpleImputer(missing_values=nan, strategy='mean')
14 # transform the dataset
15 transformed_values = imputer.fit_transform(values)
16 # count the number of NaN values in each column
17 print(f'Missing: {isnan(transformed_values).sum()}')
```

Running the example shows that all NaN values were imputed successfully.

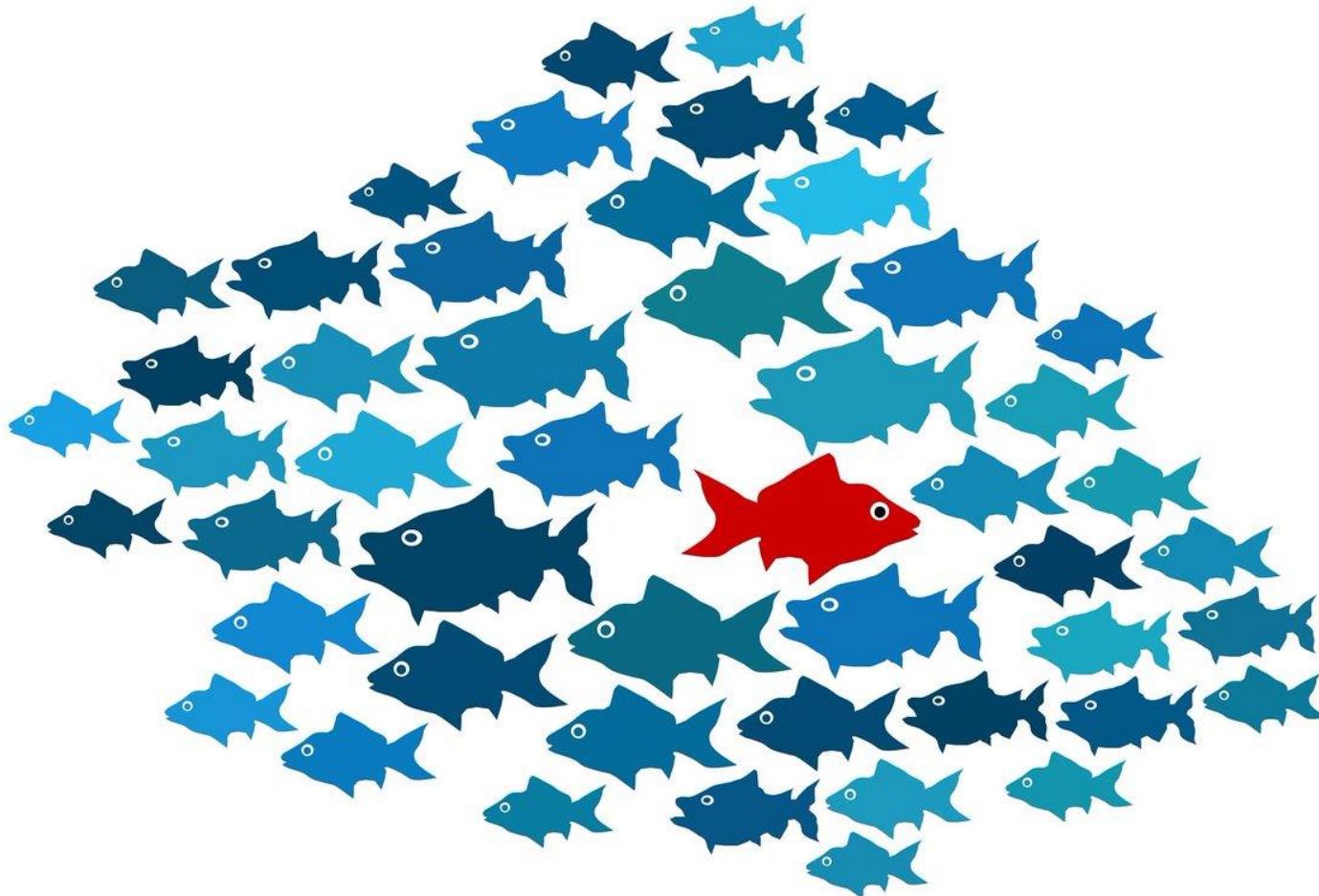
```
1 Missing: 0
```

# Outlier/Anomaly Detection

---

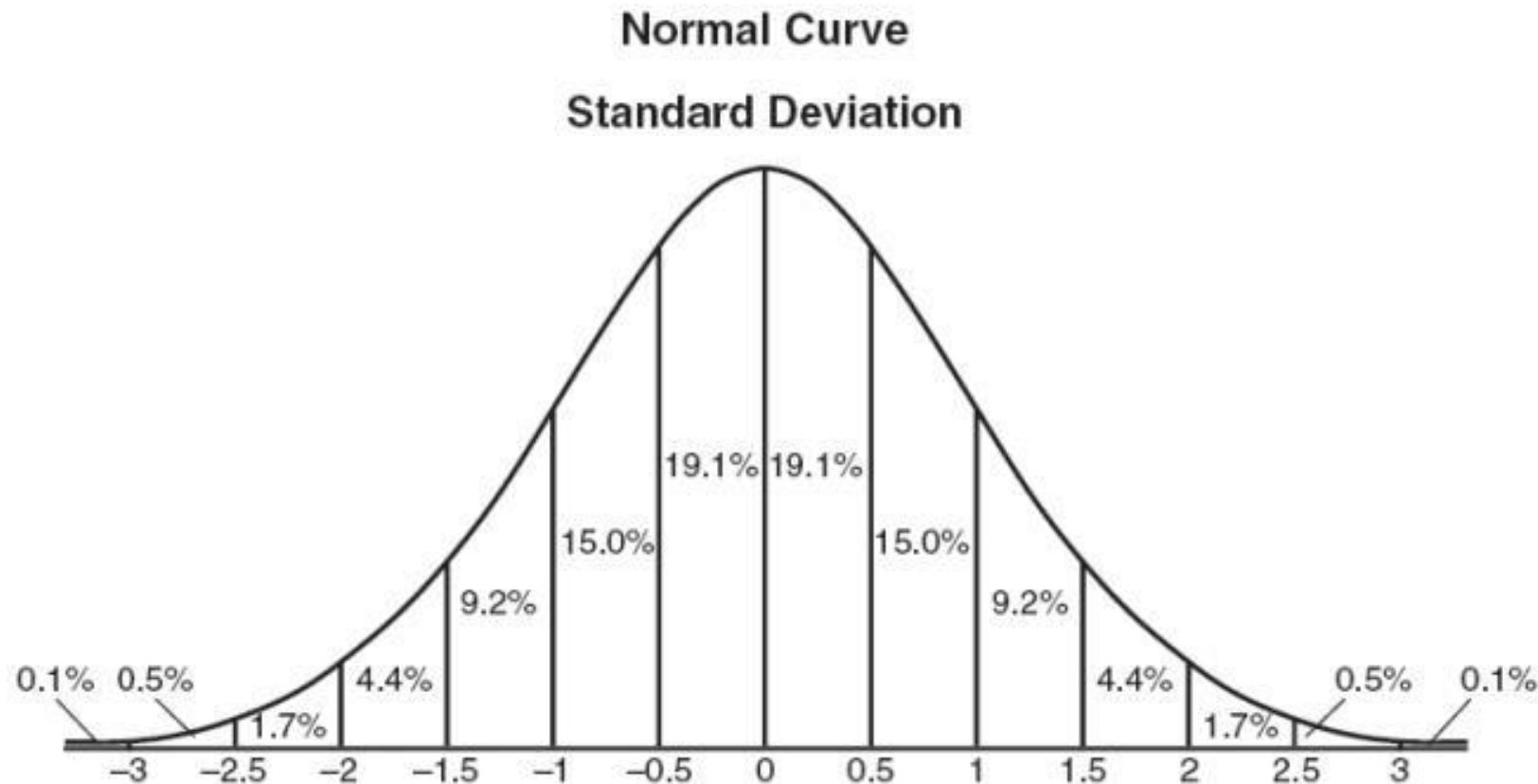
- Anomaly detection, sometimes called outlier detection, is a process of finding patterns or instances in a dataset that deviate significantly from the expected or “normal behavior.”
- The definition of both “normal” and anomalous data significantly varies depending on the context.

# Outlier/Anomaly Detection



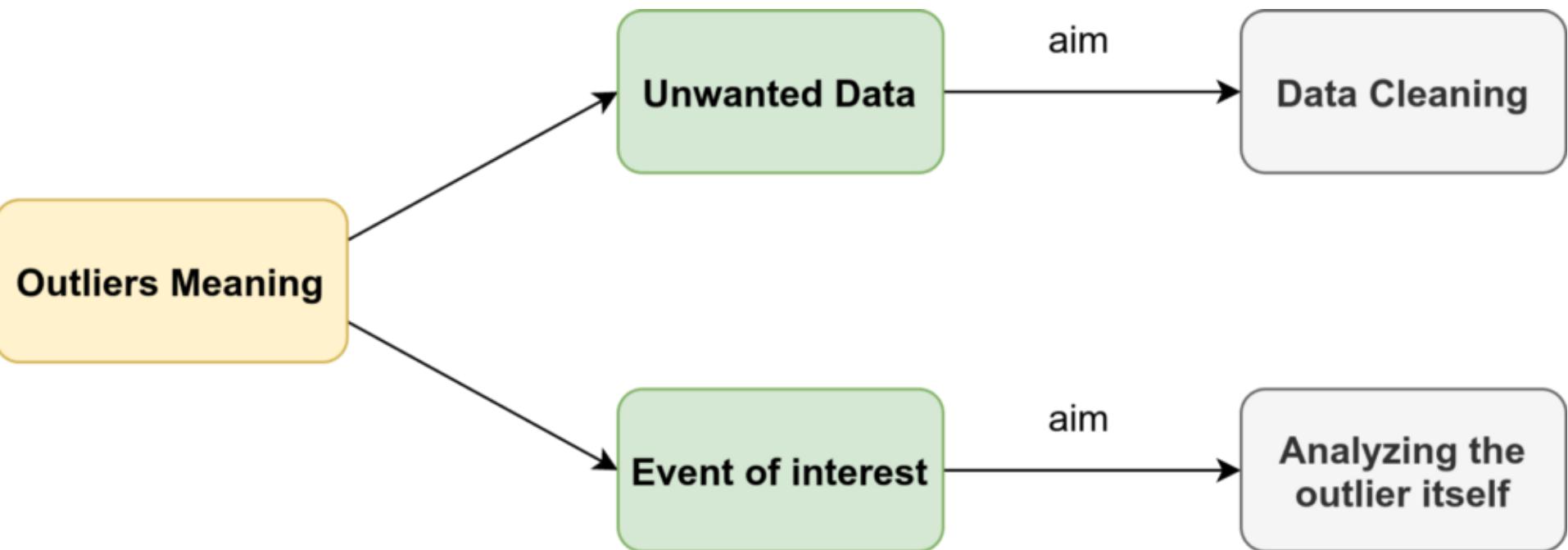
# Outlier/Anomaly Detection

The data points which fall below mean-3(sigma) or above mean+3(sigma) are outliers.



# Outlier/Anomaly Detection

- An observation which deviates so much from other observations as to arouse suspicions that it was generated by a different mechanism.



# Outlier/Anomaly Detection

---

- Financial transactions
  - Normal: Routine purchases and consistent spending by an individual in London.
  - Outlier: A massive withdrawal from Ireland from the same account, hinting at potential fraud.

# Outlier/Anomaly Detection

---

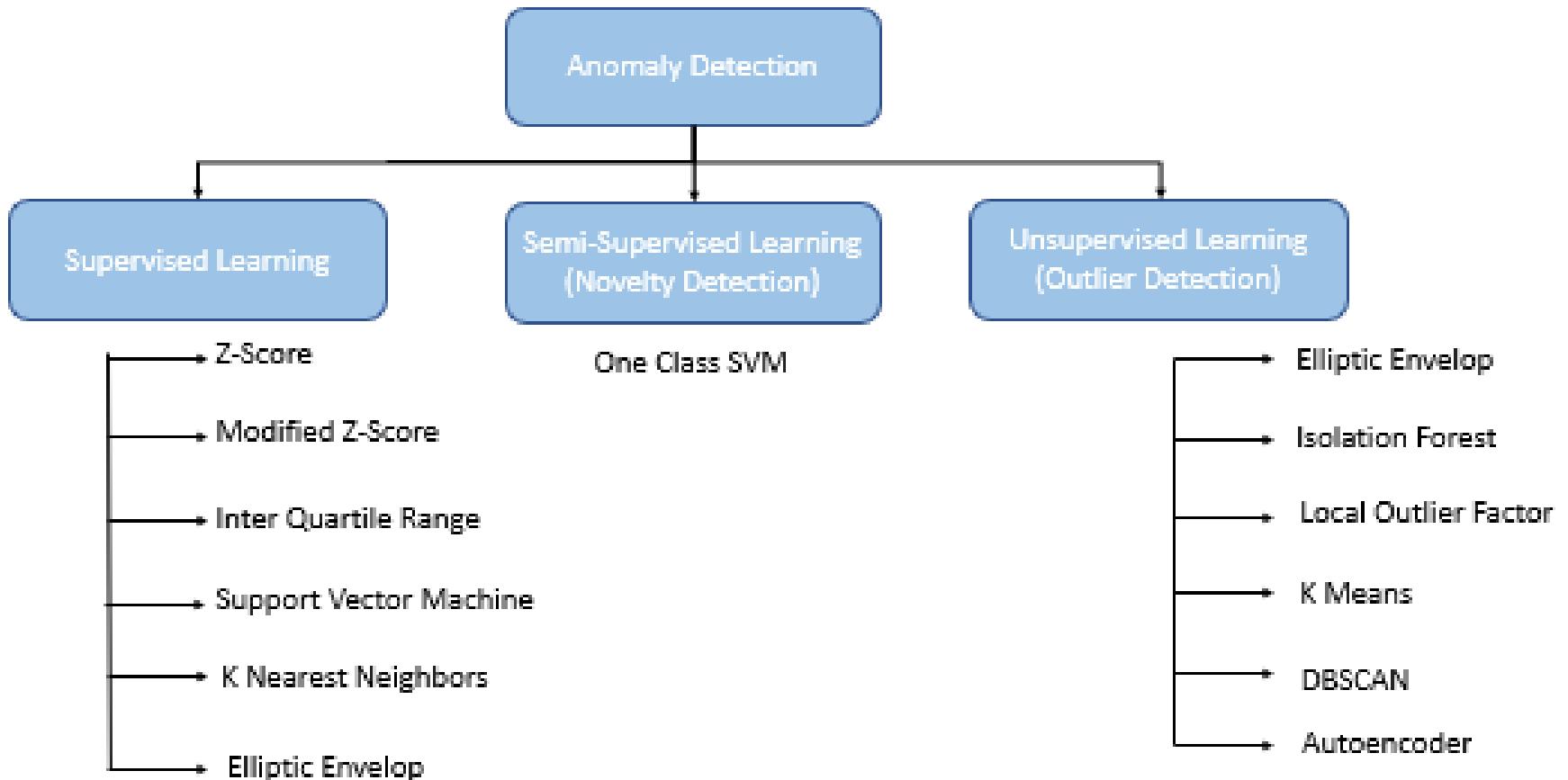
- Network traffic in cybersecurity
  - Normal: Regular communication, steady data transfer, and adherence to protocol.
  - Outlier: Abrupt increase in data transfer or use of unknown protocols signaling a potential breach or malware.

# Outlier/Anomaly Detection

---

- Patient vital signs monitoring
  - Normal: Stable heart rate and consistent blood pressure
  - Outlier: Sudden increase in heart rate and decrease in blood pressure, indicating a potential emergency or equipment failure.

# Types of Outliers



# Types of Outliers

---

- Univariate outliers exist in a single variable or feature in isolation.
- Univariate outliers are extreme or abnormal values that deviate from the typical range of values for that specific feature.
- Multivariate outliers are found by combining the values of multiple variables at the same time.

# Types of Outliers

---

- For example, consider a dataset of housing prices in a neighborhood.
- Most houses cost between \$200,000 and \$400,000, but there is House A with an exceptionally high price of \$1,000,000.
- When we analyze only the price, House A is a clear outlier.

# Types of Outliers

- Now, let's add two more variables to our dataset: the square footage and the number of bedrooms.
- When we consider the square footage, the number of bedrooms, and the price, it's House B that looks odd:
  - It has half the square footage of the mean house price.
  - It has only one bedroom.
  - It costs the top of the range \$380.000.
- When we look at these variables individually, they seem ordinary.
- Only when we put them together do we find out that House B is a clear multivariate outlier.

# Anomaly Detection Methods And When to Use Each One



- For univariate outlier detection, the most popular methods are:
  - Z-score (standard score): the z-score measures how many standard deviations a data point is away from the mean.
  - Generally, instances with a z-score over 3 are chosen as outliers.
  - Interquartile range (IQR): The IQR is the range between the first quartile (Q1) and the third quartile (Q3) of a distribution.
  - When an instance is beyond Q1 or Q3 for some multiplier of IQR, they are considered outliers. The most common multiplier is 1.5, making the outlier range  $[Q1 - 1.5 * \text{IQR}, Q3 + 1.5 * \text{IQR}]$ .

# Anomaly Detection Methods And When to Use Each One



- Modified z-scores: like z-scores, but modified z-scores use the median and a measure called Median Absolute Deviation (MAD) to find outliers.
- Since mean and standard deviation are easily skewed by outliers, modified z-scores are generally considered more robust.

## Z-Score

---

$$Z = (X - \mu)\sigma$$

where,

- $X$  = a single raw data value
- $\mu$  = population mean
- $\sigma$  = population standard deviation

Refer AnomalyOutliers.py

# Real-World Applications of Anomaly Detection



- Cybersecurity
- Healthcare
- Industrial equipment monitoring
- Network intrusion detection
- Energy grid monitoring
- E-commerce and user behavior analysis
- Quality control in manufacturing

# Data Visualization

- Data visualization provides a good, organized pictorial representation of the data which makes it easier to understand, observe, analyze.



**DATA  
VISUALIZATION  
WITH PYTHON**



# Data Visualization Libraries

---

- Matplotlib: low level, provides lots of freedom
- Pandas Visualization: easy to use interface, built on Matplotlib
- Seaborn: high-level interface, great default styles
- Plotnine: based on R's ggplot2, uses Grammar of Graphics
- Plotly: can create interactive plots

# Feature Selection

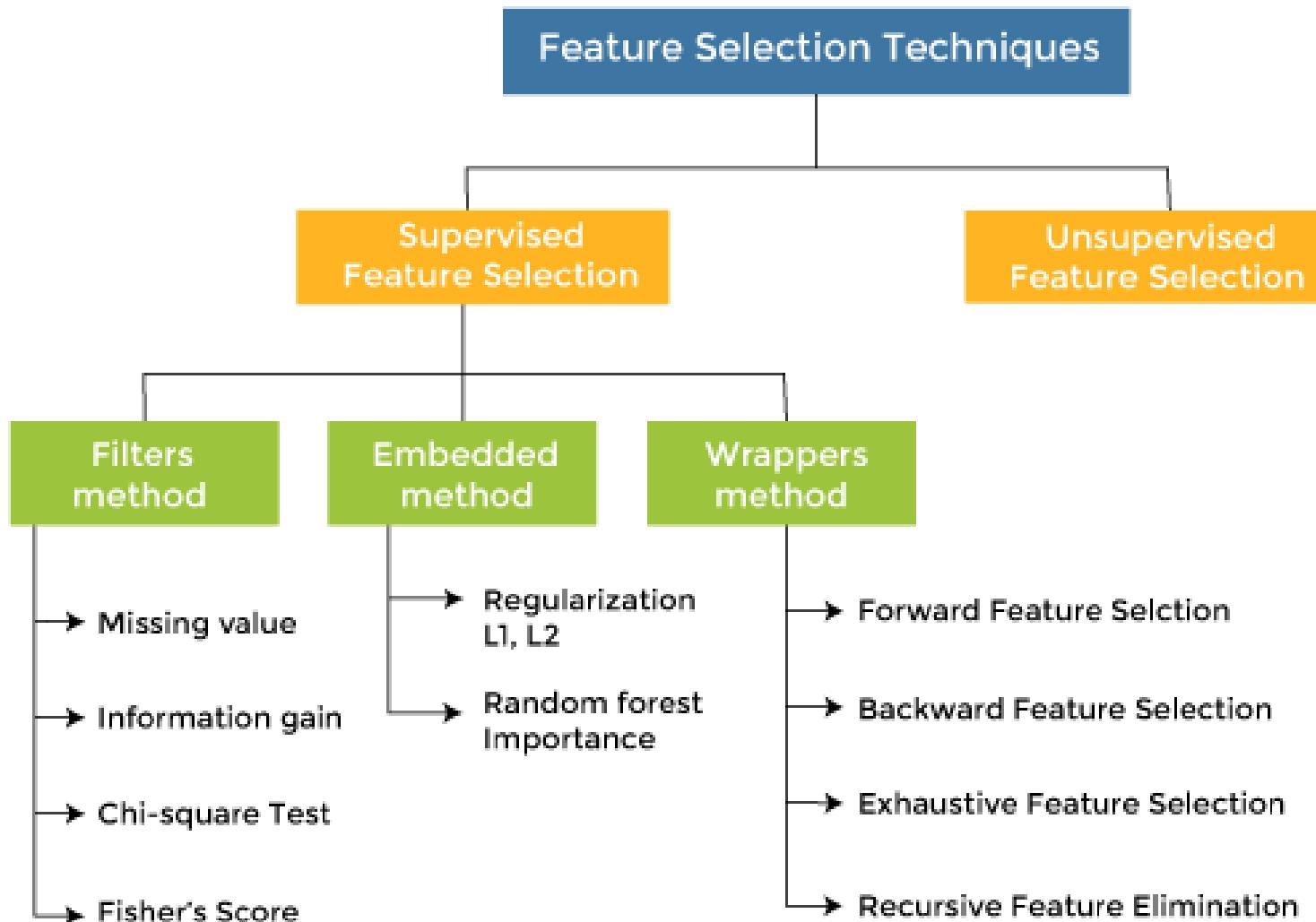
- Feature selection is the process of reducing the number of input variables when developing a predictive model.
- It is desirable to reduce the number of input variables to both reduce the computational cost of modeling and, in some cases, to improve the performance of the model.
- Statistical-based feature selection methods involve evaluating the relationship between each input variable and the target variable using statistics.
- Selecting those input variables that have the strongest relationship with the target variable

# Feature Selection

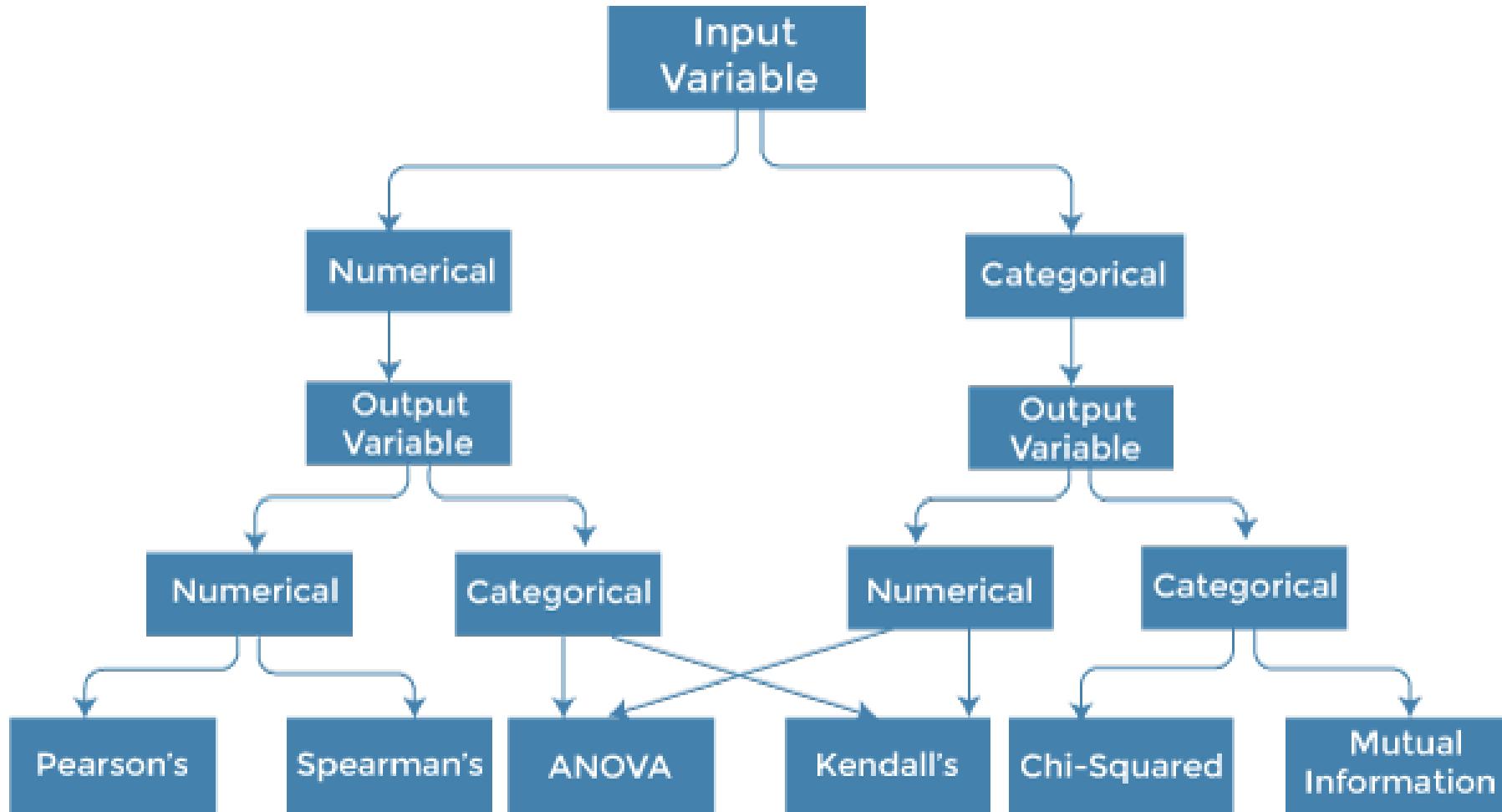
---

- These methods can be fast and effective, although the choice of statistical measures depends on the data type of both the input and output variables.

# Feature Selection Techniques



# Statistics for Filter-Based Feature Selection Methods



# Statistics for Filter-Based Feature Selection Methods



Input Variable	Output Variable	Feature Selection technique
Numerical	Numerical	<ul style="list-style-type: none"><li>○ Pearson's correlation coefficient (For linear Correlation).</li><li>○ Spearman's rank coefficient (for non-linear correlation).</li></ul>
Numerical	Categorical	<ul style="list-style-type: none"><li>○ ANOVA correlation coefficient (linear).</li><li>○ Kendall's rank coefficient (nonlinear).</li></ul>
Categorical	Numerical	<ul style="list-style-type: none"><li>○ Kendall's rank coefficient (linear).</li><li>○ ANOVA correlation coefficient (nonlinear).</li></ul>
Categorical	Categorical	<ul style="list-style-type: none"><li>○ Chi-Squared test (contingency tables).</li><li>○ Mutual Information.</li></ul>

# How to select features and what are Benefits of performing feature selection before modeling your data?



- Reduces Overfitting: Less redundant data means less opportunity to make decisions based on noise.
- Improves Accuracy: Less misleading data means modeling accuracy improves.
- Reduces Training Time: fewer data points reduce algorithm complexity and algorithms train faster.

# Univariate Selection

- Statistical tests can be used to select those features that have the strongest relationship with the output variable.
- The scikit-learn library provides the SelectKBest class that can be used with a suite of different statistical tests to select a specific number of features.
- Many different statistical test scan be used with this selection method.
- For example the ANOVA F-value method is appropriate for numerical inputs and categorical data, as we see in the Pima dataset.

# Feature Importance

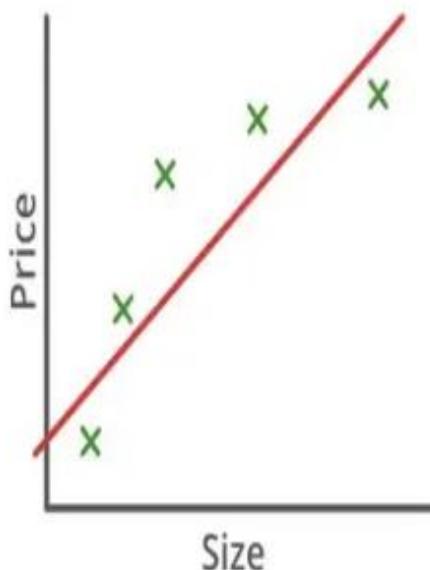
---

- We can get the feature importance of each feature of your dataset by using the feature importance property of the model.
- Feature importance gives us a score for each feature of our data, the higher the score more important or relevant is the feature towards your output variable.
- Feature importance is an inbuilt class that comes with Tree Based Classifiers, we will be using Extra Tree Classifier for extracting the top 10 features for the dataset.

# Correlation Matrix with Heatmap

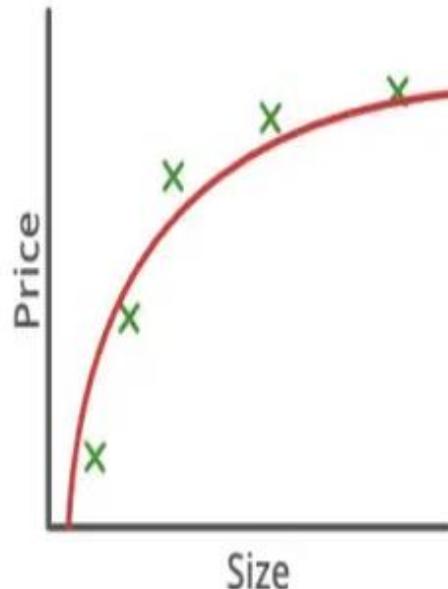
- Correlation states how the features are related to each other or the target variable.
- Correlation can be positive (increase in one value of feature increases the value of the target variable) or negative (increase in one value of feature decreases the value of the target variable)
- Heatmap makes it easy to identify which features are most related to the target variable, we will plot heatmap of correlated features using the seaborn library.

# Under Fitting and Over Fitting



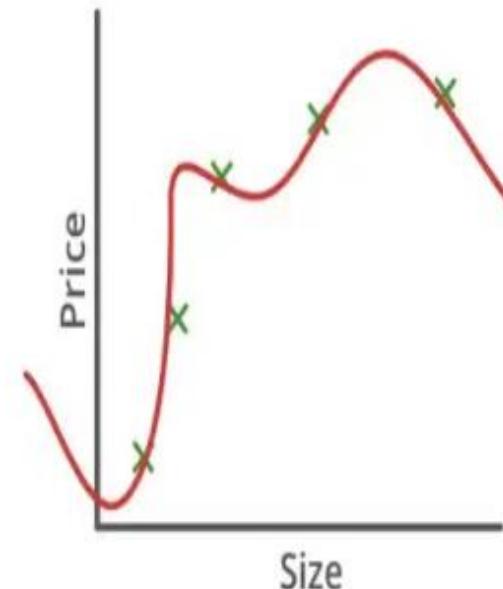
$$\theta_0 + \theta_1 x$$

**High Bias**  
(Underfitting)



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

**Low Bias, Low Variance**  
(Goodfitting)



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

**High Variance**  
(Overfitting)

# Under fitting

- A statistical model or a machine learning algorithm is said to have underfitting when a model is too simple to capture data complexities.
- It represents the inability of the model to learn the training data effectively result in poor performance both on the training and testing data.
- In simple terms, an underfit model's are inaccurate, especially when applied to new, unseen examples.
- It mainly happens when we uses very simple model with overly simplified assumptions.
- To address underfitting problem of the model, we need to use more complex models, with enhanced feature representation, and less regularization.

## Reasons For Under fitting

---

- The model is too simple, So it may be not capable to represent the complexities in the data.
- The input features which is used to train the model is not the adequate representations of underlying factors influencing the target variable.
- The size of the training dataset used is not enough.
- Excessive regularization are used to prevent the overfitting, which constraint the model to capture the data well.
- Features are not scaled.

# Techniques to Reduce Underfitting

---

- Increase model complexity.
- Increase the number of features, performing feature engineering.
- Remove noise from the data.
- Increase the number of epochs or increase the duration of training to get better results.

# Overfitting in Machine Learning

- A statistical model is said to be overfitted when the model does not make accurate predictions on testing data.
- When a model gets trained with so much data, it starts learning from the noise and inaccurate data entries in our data set.
- when testing with test data results in High variance.
- Then the model does not categorize the data correctly, because of too many details and noise.
- The causes of overfitting are the non-parametric and non-linear methods because these types of machine learning algorithms have more freedom in building the model based on the dataset and therefore they can really build unrealistic models.
- A solution to avoid overfitting is using a linear algorithm if we have linear data or using the parameters like the maximal depth if we are using decision trees.

## Reasons for Overfitting:

---

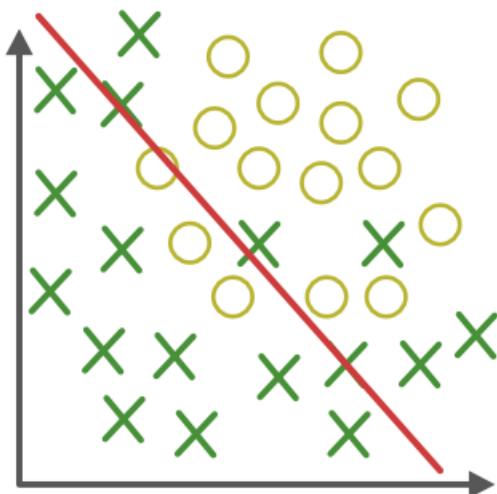
- High variance and low bias.
- The model is too complex.
- The size of the training data.

# Techniques to Reduce Overfitting

---

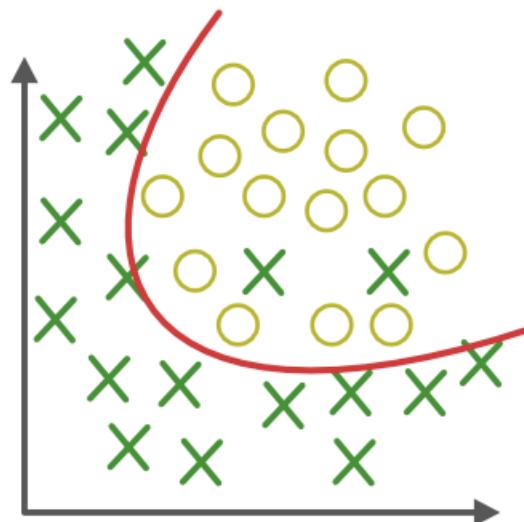
- Increase training data.
- Reduce model complexity.
- Early stopping during the training phase (have an eye over the loss over the training period as soon as loss begins to increase stop training).
- Ridge Regularization and Lasso Regularization.
- Use dropout for neural networks to tackle overfitting.

# Under fitting and Over fitting

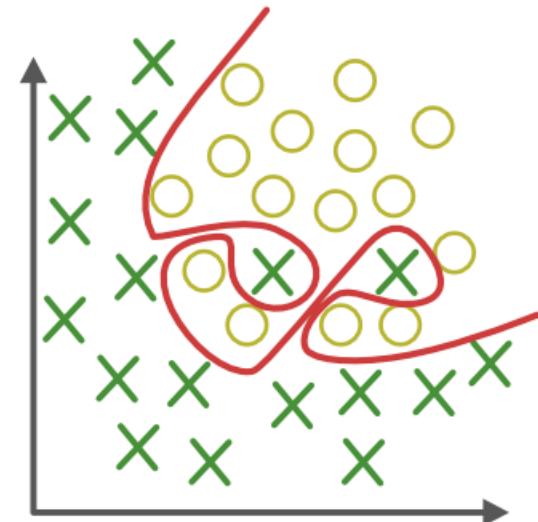


**Under-fitting**

(too simple to explain the variance)



**Appropriate-fitting**



**Over-fitting**

(forcefitting--too good to be true)

# Good Fit in a Statistical Model

- Ideally, the case when the model makes the predictions with 0 error, is said to have a good fit on the data.
- This situation is achievable at a spot between overfitting and underfitting.
- With the passage of time, our model will keep on learning, and thus the error for the model on the training and testing data will keep on decreasing.
- If it will learn for too long, the model will become more prone to overfitting due to the presence of noise and less useful details.
- Hence the performance of our model will decrease.
- To get a good fit, we will stop at a point just before where the error starts increasing.
- At this point, the model is said to have good skills in training datasets as well as our unseen testing dataset.

# Bias-Variance Trade-off

---

- Bias is the difference between the average prediction of our model and the correct value which we are trying to predict.
- Model with high bias pays very little attention to the training data and oversimplifies the model.
- It always leads to high error on training and test data.

# Bias-Variance Trade-off

- Variance is the variability of model prediction for a given data point or a value which tells us spread of our data.
- Model with high variance pays a lot of attention to training data and does not generalize on the data which it hasn't seen before.
- As a result, such models perform very well on training data but has high error rates on test data.

## Mathematically

---

- Let the variable we are trying to predict as Y and other covariates as X. We assume there is a relationship between the two such that
- $Y = f(X) + e$
- Where e is the error term and it's normally distributed with a mean of 0.
- We will make a model  $\hat{f}(X)$  of  $f(X)$  using linear regression or any other modeling technique.

# Mathematically

So the expected squared error at a point x is

$$Err(x) = E \left[ (Y - \hat{f}(x))^2 \right]$$

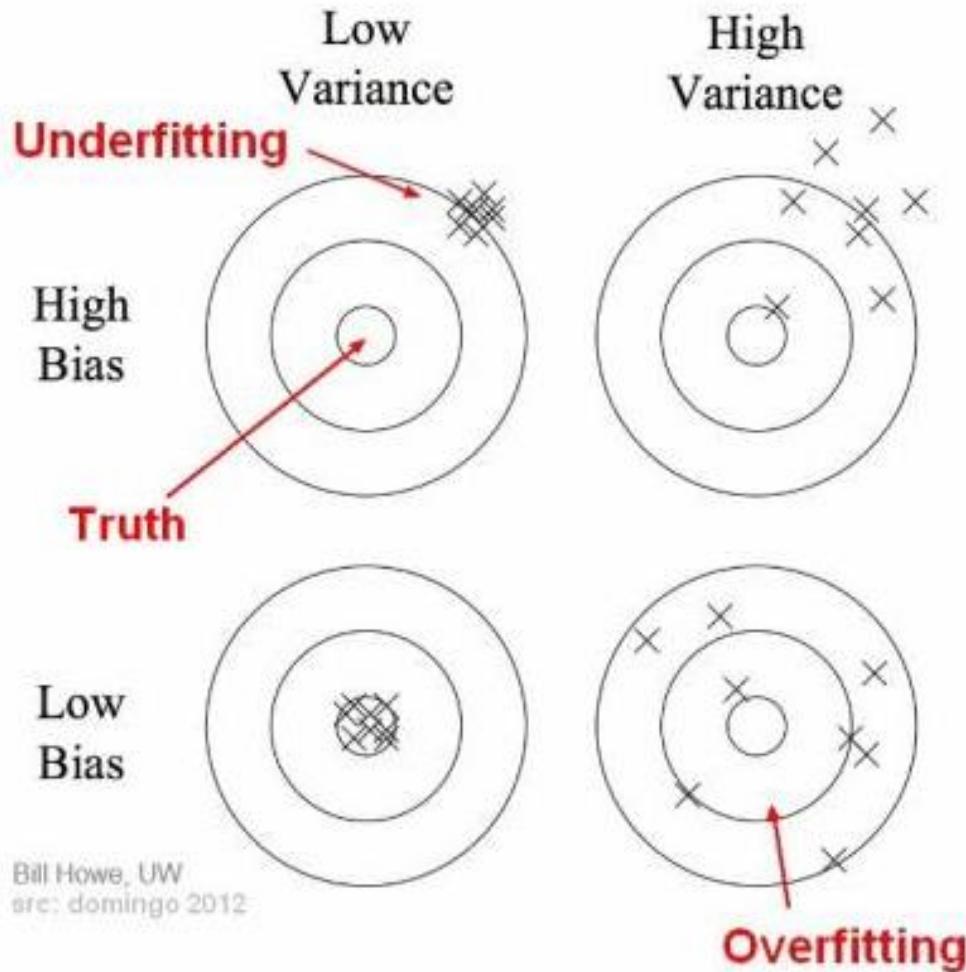
The Err(x) can be further decomposed as

$$Err(x) = \left( E[\hat{f}(x)] - f(x) \right)^2 + E \left[ \left( \hat{f}(x) - E[\hat{f}(x)] \right)^2 \right] + \sigma_e^2$$

$$Err(x) = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

Err(x) is the sum of Bias<sup>2</sup>, variance and the irreducible error.

# Mathematically



# Why is Bias Variance Tradeoff?

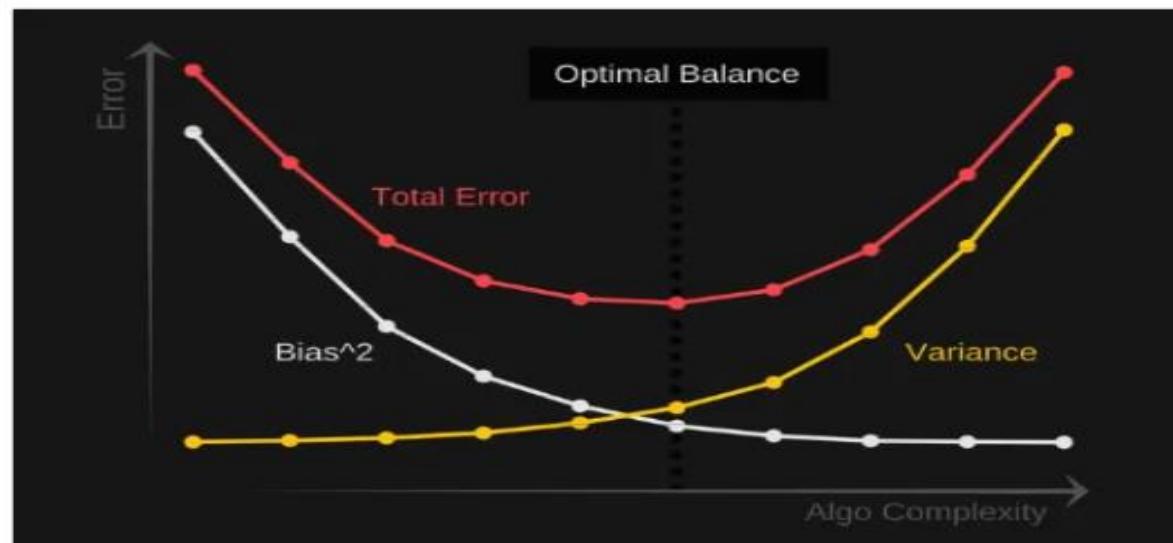
- If our model is too simple and has very few parameters, then it may have high bias and low variance.
- On the other hand, if our model has large number of parameters then it's going to have high variance and low bias.
- So, we need to find the right/good balance without overfitting and underfitting the data.
- This tradeoff in complexity is why there is a tradeoff between bias and variance.
- An algorithm can't be more complex and less complex at the same time.

# Total Error

## Total Error

To build a good model, we need to find a good balance between bias and variance such that it minimizes the total error.

$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$



An optimal balance of bias and variance would never overfit or underfit the model.



# Evaluation Metrics

# Hypothesis Testing

---

- Hypothesis testing refers to a process used by analysts to assess the plausibility of a hypothesis by using sample data.
- In hypothesis testing, statisticians formulate two hypotheses: the null hypothesis and the alternative hypothesis.
- A null hypothesis determines there is no difference between two groups or conditions, while the alternative hypothesis determines that there is a difference.
- Researchers evaluate the statistical significance of the test based on the probability that the null hypothesis is true.

# Hypothesis Testing



Lead content  $\leq 2.5$  PPM

so they extracted out a small sample out  
of a very book population



Sample = 10 Thousand

# Hypothesis Testing

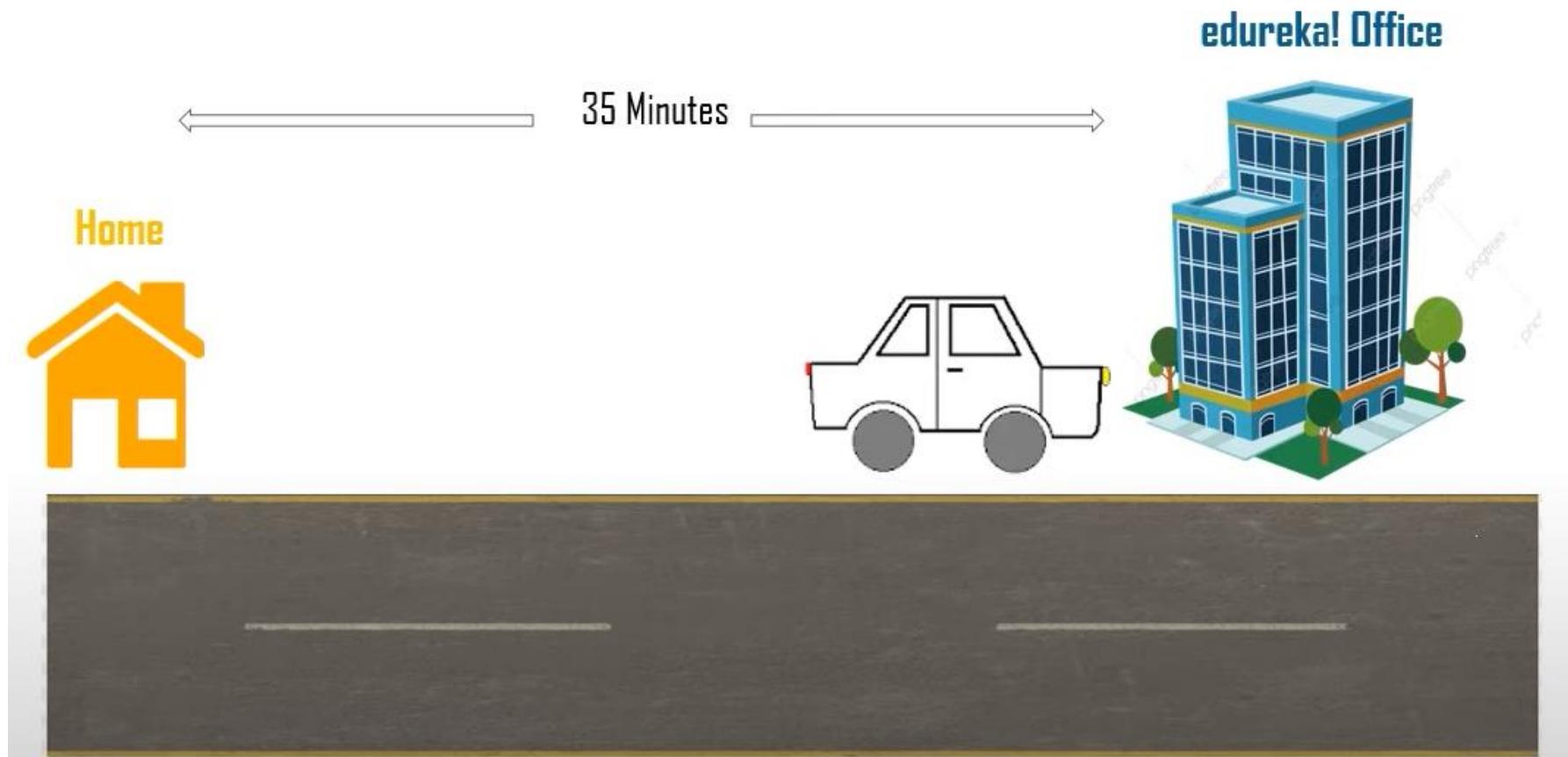
---

- Hypothesis testing is a statistical method used to determine if there is enough evidence in a sample data to draw conclusions about a population.
- Let's discuss few examples of statistical hypothesis from real-life -
- A teacher assumes that 60% of his college's students come from lower-middle-class families.
- A doctor believes that 3D (Diet, Dose, and Discipline) is 90% effective for diabetic patients.



DATA SCIENCE

# Hypothesis Testing



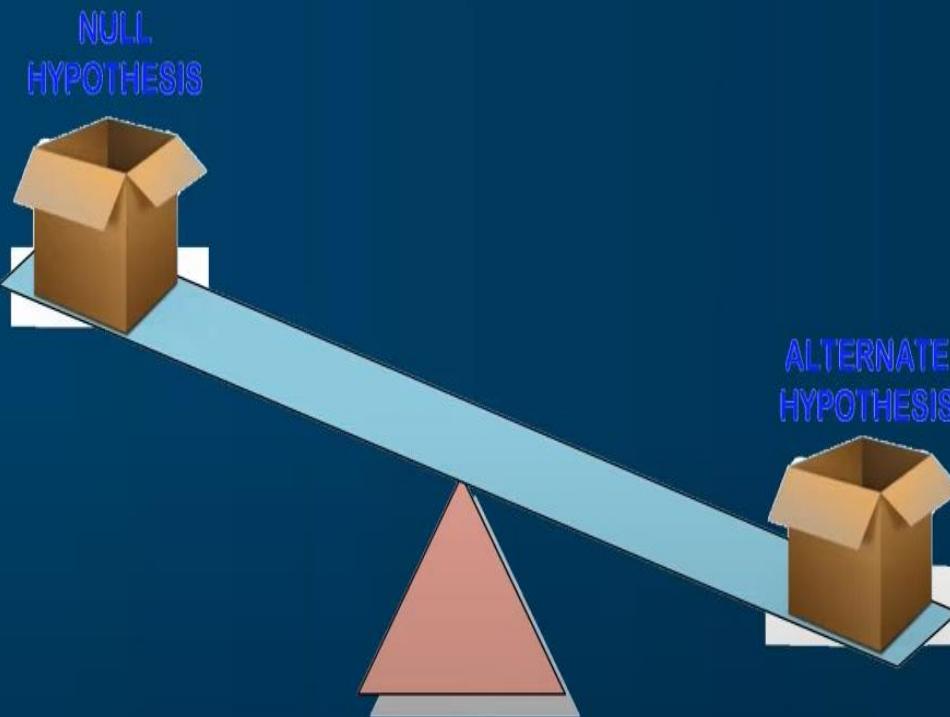


DATA SCIENCE

# Null & Alternate Hypothesis



# Null & Alternate Hypothesis



## (H<sub>0</sub>) Null Hypothesis:

- Prevailing belief about the population
- Assumes that status quo is true

## (H<sub>1</sub>) Alternate Hypothesis

- Claim that opposes the null hypothesis

# Null & Alternate Hypothesis

		The Person is	
		Innocent	Guilty
The Judge Says	Innocent		
	Guilty		
		No Error	Type 2 error
		Type 1 error	No Error

# Null & Alternate Hypothesis

---

- First Case: The person is innocent, and the judge identifies the person as innocent
- Second Case: The person is innocent, and the judge identifies the person as guilty
- Third Case: The person is guilty, and the judge identifies the person as innocent
- Fourth Case: The person is guilty, and the judge identifies the person as guilty

# Null & Alternate Hypothesis

## Formulation of Null & Alternate Hypothesis:

If your claim statement has words like "at least", "at most", "less than", or "greater than", you cannot formulate the null hypothesis just from the claim statement

Rule to formulate the null and alternate hypotheses:

- Null hypothesis signs: =  $\bar{O}R \leq \bar{O}R \geq$
- Alternate hypothesis signs:  $\neq \bar{O}R > \bar{O}R <$

Situation 1:



Situation 2:

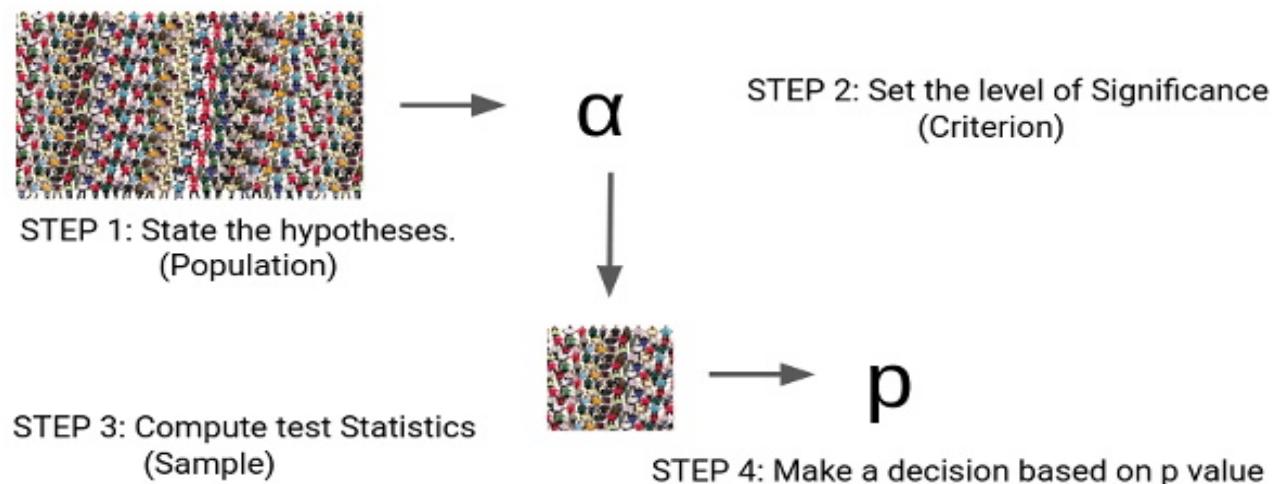


# Steps to Perform Hypothesis Testing

## Steps to Perform Hypothesis Testing

There are four steps to performing Hypothesis Testing:

1. Set the Null and Alternate Hypotheses
2. Set the Significance Level, Criteria for a decision
3. Compute the test statistic
4. Make a decision



# Steps to Perform Hypothesis Testing

- It must be noted that z-Test & t-Tests are Parametric Tests, which means that the Null Hypothesis is about a population parameter, which is less than, greater than, or equal to some value.
- Steps 1 to 3 are quite self-explanatory but on what basis can we make decision in step 4? What does this p-value indicate?
- We can understand this p-value as the measurement of the Defense Attorney's argument.
- If the p-value is less than  $\alpha$ , we reject the Null Hypothesis, and if the p-value is greater than  $\alpha$ , we fail to reject the Null Hypothesis.

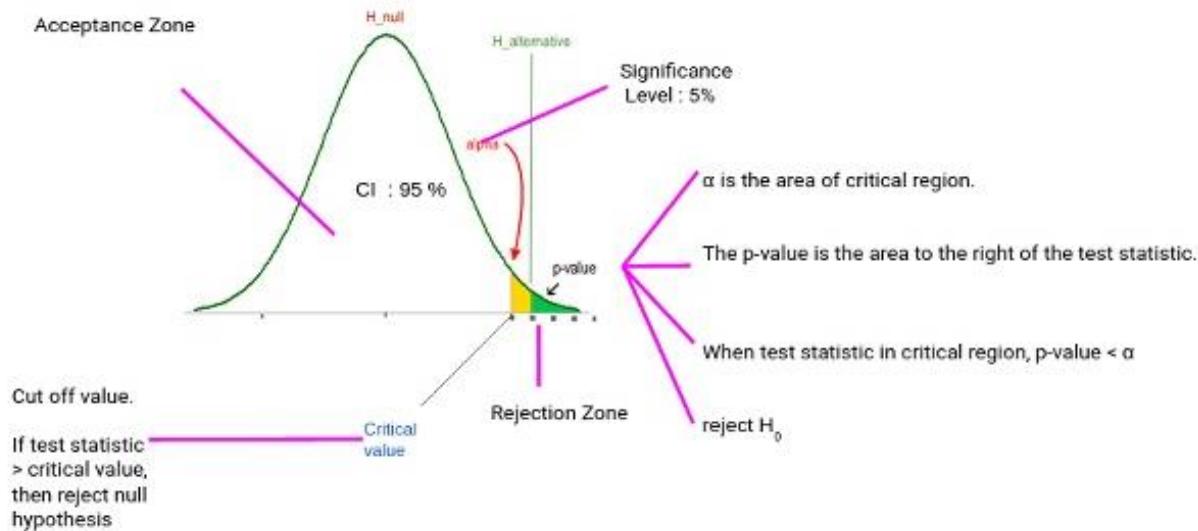
## Critical P-Value

---

- Critical Value is the cut off value between Acceptance Zone and Rejection Zone.
- We compare our test score to the critical value and if the test score is greater than the critical value, that means our test score lies in the Rejection Zone and we reject the Null Hypothesis.
- On the other hand, if the test score is less than the Critical Value, that means the test score lies in the Acceptance Zone and we fail to reject the null Hypothesis.

# Critical P-Value

Let's understand the logic of Hypothesis Testing with the graphical representation for [Normal Distribution](#).



The above visualization helps to understand the z-value and its relation to the critical value. Typically, we set the Significance level at 10%, 5%, or 1%. If our test score lies in the Acceptance Zone, we fail to reject the Null Hypothesis. If our test score lies in the Critical Zone, we reject the Null Hypothesis and accept the Alternate Hypothesis.

# Z-Test

- z tests are a statistical way of testing a Null Hypothesis when either:
  - We know the population variance, or
  - We do not know the population variance, but our sample size is large  $n \geq 30$
- If we have a sample size of less than 30 and do not know the population variance, we must use a t-test.
- This is how we judge when to use the z-test vs the t-test.
- Further, it is assumed that the z-statistic follows a standard normal distribution.
- In contrast, the t-statistics follows the t-distribution with a degree of freedom equal to  $n-1$ , where  $n$  is the sample size.
- It must be noted that the samples used for z-test or t-test must be independent sample and must have a distribution identical to the population distribution.
- This makes sure that the sample is not “biased” to/against the Null Hypothesis which we want to validate/invalidate.

# Types of Hypothesis Testing

---

- Z Test
- T Test
- Chi Square Test

# Finding P-Value

---

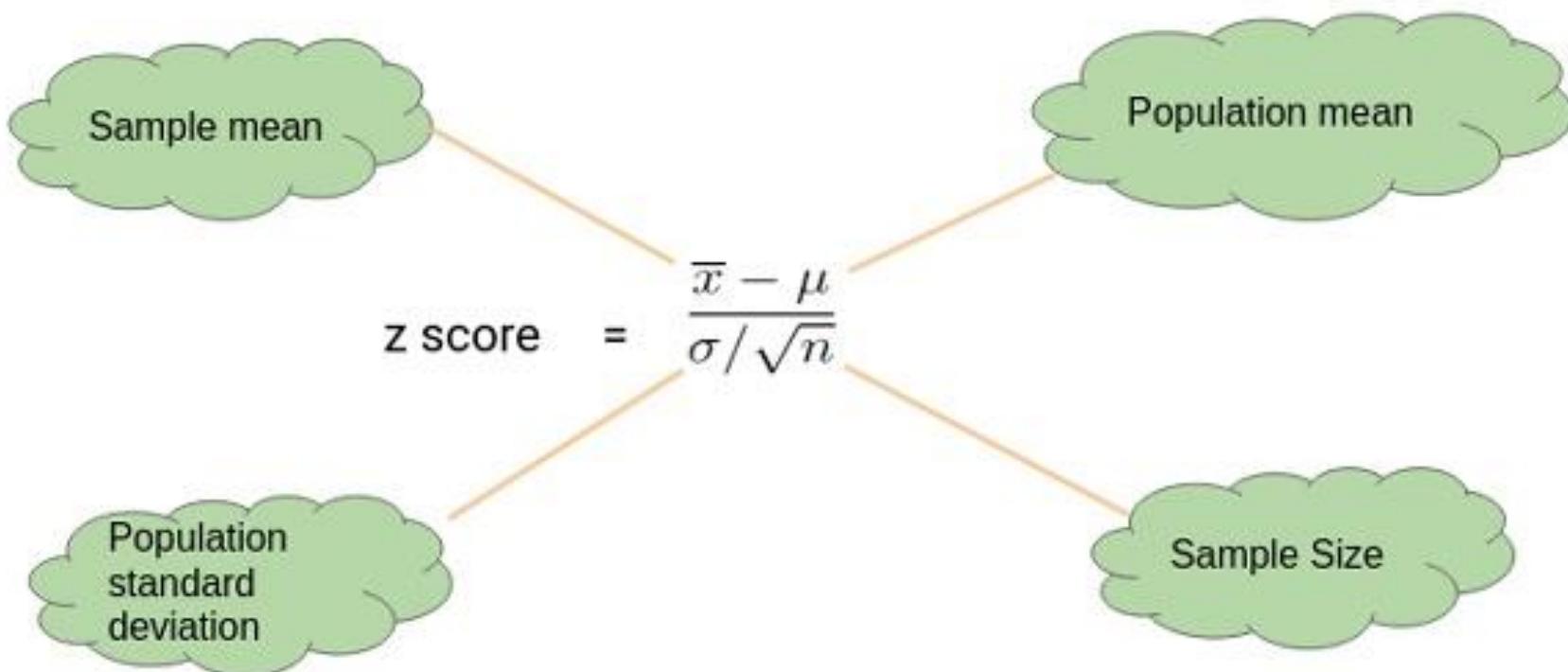
- P stands for probability here.
- To calculate the p-value, the chi-square test is used in statistics.
- The different values of p indicates the different hypothesis interpretation, are given below:
- $P \leq 0.05$ ; Hypothesis rejected
- $P > .05$ ; Hypothesis Accepted
- Probability is all about chance or risk or uncertainty.
- It is the possibility of the outcome of the sample or the occurrence of an event.

# Finding P-Value

---

- But when we talk about statistics, it is more about how we handle various data using different techniques.
- It helps to represent complicated data or bulk data in a very easy and understandable way.
- It describes the collection, analysis, interpretation, presentation, and organization of data.
- The concept of both probability and statistics is related to the chi-squared test.

# Z-Test



# Z-Test

Let's say we need to determine if girls on average score higher than 600 in the exam. We have the information that the standard deviation for girls' scores is 100. So, we collect the data of 20 girls by using random samples and record their marks. Finally, we also set our  $\alpha$  value (significance level) to be 0.05.



Score
650
730
510
670
480
800
690
530
590
620
710
670
640
780
650
490
800
600
510
700

# Z-Test

In this example:

- Mean Score for Girls is 641
- The number of data points in the sample is 20
- The population mean is 600
- Standard Deviation for Population is 100

$$\begin{aligned}
 z \text{ score} &= \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \\
 &= \frac{641 - 600}{100 / \sqrt{20}} \\
 &= 1.8336
 \end{aligned}$$

$$p \text{ value} = .033357.$$

Critical Value = 1.645

Z score > Critical Value

P value < 0.05

# P-Value : Probability of getting less than a Z-score  
 $p\_value = 1 - \text{stats.norm.cdf}(z\_score)$



$$H_0: \mu \leq 600$$

$$H_1: \mu > 600$$



# Z-Test

## P Value from Z Score Calculator

This is very easy: just stick your Z score in the box marked Z score, select your significance level and whether you're testing a one or two-tailed hypothesis (if you're not sure, go with the defaults), then press the button!

If you need to derive a Z score from raw data, [you can find a Z test calculator here](#).

Z score:

1.8335

Significance Level:

- 0.01
- 0.05
- 0.10

One-tailed or two-tailed hypothesis?:

- One-tailed
- Two-tailed

The P-Value is .033364.

The result is significant at  $p < .05$ .

Data  
volume vs  
speed.  
You can  
have both

# Types of Hypothesis Testing

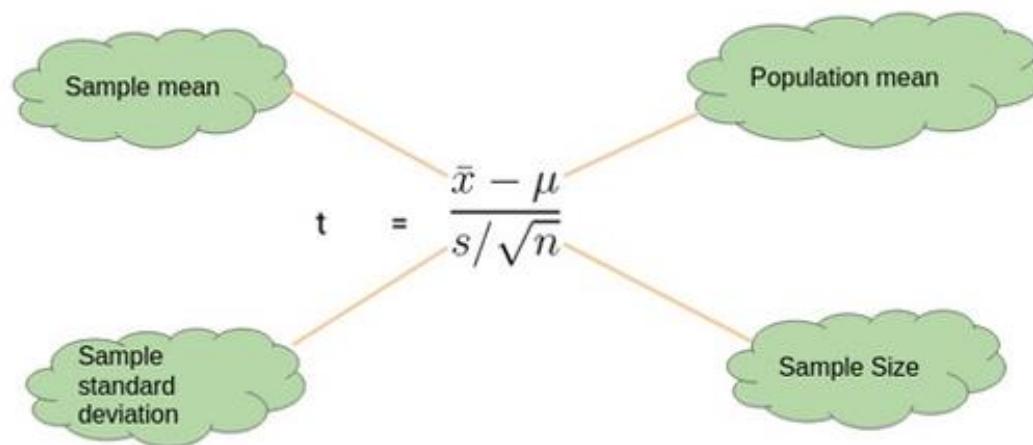
---

- T Test
  - T-tests are a statistical way of testing a hypothesis when:
  - We do not know the population variance
  - Our sample size is small,  $n < 30$ .

# Types of Hypothesis Testing

## One-Sample T-Test

We perform a One-Sample t-test when we want to **compare a sample mean with the population mean**. The difference from the z-Test is that we do **not have the information on Population Variance** here. We use the **sample standard deviation** instead of population standard deviation in this case.



# Types of Hypothesis Testing

## Here's an Example to Understand a One Sample T-Test

Let's say we want to determine if on average girls score more than 600 in the exam. We do not have the information related to variance (or standard deviation) for girls' scores. To perform t-test, we randomly collect the data of 10 girls with their marks and choose our  $\alpha$  value (significance level) to be 0.05 for Hypothesis Testing.



Girls_Score
587
602
627
610
619
622
605
608
596
592

# Types of Hypothesis Testing

In this example:

- Mean Score for Girls is 606.8
- The size of the sample is 10
- The population mean is 600
- Standard Deviation for the sample is 13.14

$$\begin{aligned}
 t &= \frac{\bar{x} - \mu}{s/\sqrt{n}} \\
 &= \frac{606.8 - 600}{13.14/\sqrt{10}} \\
 &= 1.64
 \end{aligned}$$

Critical Value = 1.833

t score < Critical Value

P value = 0.0678

P value > 0.05



$$H_0 : \mu \leq 600$$

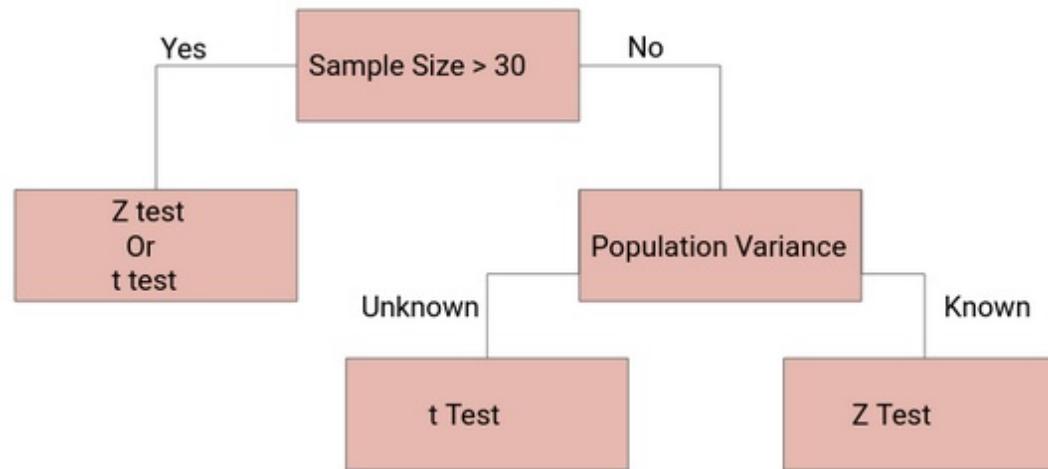
$$H_1 : \mu > 600$$



Our p-value is greater than 0.05 thus we fail to reject the null hypothesis and don't have enough evidence to support the hypothesis that on average, girls score more than 600 in the exam.

# Deciding Between Z-Test and T-Test

So when should we perform the z-test, and when should we perform the t-Test? It's a critical question we need to answer if we want to master statistics.



If the sample size is large enough, then the z-Test and t-Test will conclude with the same results. For a **large sample size**, **Sample Variance will be a better estimate** of Population variance, so even if population variance is unknown we can **use the z-test using sample variance**.

Similarly, for a **Large Sample**, we have a high degree of freedom. And since t-distribution approaches the **normal distribution**, the difference between the z score and t score is negligible.

# Deciding Between Z-Test and T-Test

	Z Test	T Test
Assumption	Population standard deviation is known	Population standard deviation is unknown
Sample Size	Large sample size ( $n > 30$ )	Small sample size ( $n < 30$ )
Distribution	Z-distribution	T-distribution
Test Statistic	$(\text{Sample mean} - \text{Population mean}) / (\text{Population SD} / \sqrt{n})$	$(\text{Sample mean} - \text{Population mean}) / (\text{Sample SD} / \sqrt{n})$
Hypothesis Testing	Test for a population mean or proportion	Test for a population mean
Degrees of Freedom	Not applicable	$n - 1$
Application	Used when the population standard deviation is known and the sample size is large	Used when the population standard deviation is unknown or the sample size is small
Example	Testing whether the average height of male adults is significantly different from a known value	Testing whether a new teaching method improves student test scores compared to the old method

# Pearson's Chi Square Test

- A chi-squared test (symbolically represented as  $\chi^2$ ) is basically a data analysis on the basis of observations of a random set of variables.
- Usually, it is a comparison of two statistical data sets.
- This test was introduced by Karl Pearson in 1900 for categorical data analysis and distribution.
- So, it was mentioned as Pearson's chi-squared test.
- The chi-square test is used to estimate how likely the observations that are made would be, by considering the assumption of the null hypothesis as true.

# Chi-Square Test of Independence

---

- The chi-square test of independence also known as the chi-square test of association which is used to determine the association between the categorical variables.
- It is considered as a non-parametric test.
- It is mostly used to test statistical independence.

# Chi-Square Test of Independence

---

- The chi-square test of independence is not appropriate when the categorical variables represent the pre-test and post-test observations.
- For this test, the data must meet the following requirements:
  - Two categorical variables
  - Relatively large sample size
  - Categories of variables (two or more)
  - Independence of observations.

# Chi-Square Test of Independence

- Let us take an example of a categorical data where there is a society of 1000 residents with four neighbourhoods, P, Q, R and S.
- A random sample of 650 residents of the society is taken whose occupations are doctors, engineers and teachers.
- The null hypothesis is that each person's neighbourhood of residency is independent of the person's professional division.

# Finding P-Value

## Formula

The chi-squared test is done to check if there is any difference between the observed value and expected value. The formula for chi-square can be written as;

$$\chi^2 = \sum \frac{(\text{Observed value} - \text{Expected value})^2}{\text{Expected value}}$$

or

$$\chi^2 = \sum (O_i - E_i)^2 / E_i$$

where  $O_i$  is the observed value and  $E_i$  is the expected value.

# Chi-Square Test of Independence

---

Categories	P	Q	R	S	Total
Doctors	90	60	104	95	349
Engineers	30	50	51	20	151
Teachers	30	40	45	35	150
Total	150	150	200	150	650

# Chi-Square Test of Independence

- Assume the sample living in neighborhood P, 150, to estimate what proportion of the whole 1,000 people live in neighborhood P.
- In the same way, we take 349/650 to calculate what ratio of the 1,000 are doctors.
- By the supposition of independence under the hypothesis, we should “expect” the number of doctors in neighbourhood P is;
- $150 \times 349/650 \approx 80.54$

# Chi-Square Test of Independence

- So, by the chi-square test formula for that particular cell in the table, we get;
- $(\text{Observed} - \text{Expected})^2 / \text{Expected Value} = (90 - 80.54)^2 / 80.54 \approx 1.11$
- Some of the exciting facts about the Chi-square test are given below:
- The Chi-square statistic can only be used on numbers.
- We cannot use them for data in terms of percentages, proportions, means or similar statistical contents.
- Suppose, if we have 20% of 400 people, we need to convert it to a number, i.e. 80, before running a test statistic.
- A chi-square test will give us a p-value. The p-value will tell us whether our test results are significant or not.

# Chi-Square Test of Independence

- However, to perform a chi-square test and get the p-value, we require two pieces of information:
- (1) Degrees of freedom. That's just the number of categories minus 1.
- (2) The alpha level( $\alpha$ ). You or the researcher chooses this. The usual alpha level is 0.05 (5%), but you could also have other levels like 0.01 or 0.10.
- In elementary statistics, we usually get questions along with the degrees of freedom(DF) and the alpha level.
- Thus, we don't usually have to figure out what they are. To get the degrees of freedom, count the categories and subtract 1.

# Confusion Matrix

- A confusion matrix is a matrix that summarizes the performance of a machine learning model on a set of test data.
- It is often used to measure the performance of classification models, which aim to predict a categorical label for each input instance.
- The matrix displays the number of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN) produced by the model on the test data.

# Confusion Matrix

---

- A confusion matrix is a matrix that summarizes the performance of a machine learning model on a set of test data.
- It is often used to measure the performance of classification models, which aim to predict a categorical label for each input instance.
- The matrix displays the number of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN) produced by the model on the test data.

# Confusion Matrix

		Actual	
		Dog	Not Dog
Predicted	Dog	True Positive (TP)	False Positive (FP)
	Not Dog	False Negative (FN)	True Negative (TN)

- **True Positive (TP):** It is the total counts having both predicted and actual values are Dog.
- **True Negative (TN):** It is the total counts having both predicted and actual values are Not Dog.
- **False Positive (FP):** It is the total counts having prediction is Dog while actually Not Dog.
- **False Negative (FN):** It is the total counts having prediction is Not Dog while actually, it is Dog.

# Confusion Matrix

## Example

Index	1	2	3	4	5	6	7	8	9	10
Actual	Dog	Dog	Dog	Not Dog	Dog	Not Dog	Dog	Dog	Not Dog	Not Dog
Predicted	Dog	Not Dog	Dog	Not Dog	Dog	Dog	Dog	Dog	Not Dog	Not Dog
Result	TP	FN	TP	TN	TP	FP	TP	TP	TN	TN

- Actual Dog Counts = 6
- Actual Not Dog Counts = 4
- True Positive Counts = 5
- False Positive Counts = 1
- True Negative Counts = 3
- False Negative Counts = 1

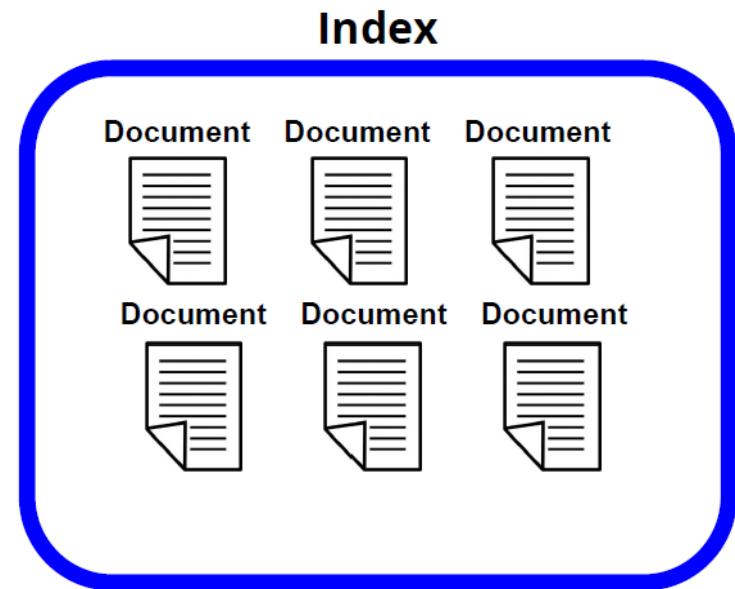
# Relevance



## Elasticsearch

Store | Search | Analyze

I store data as documents!

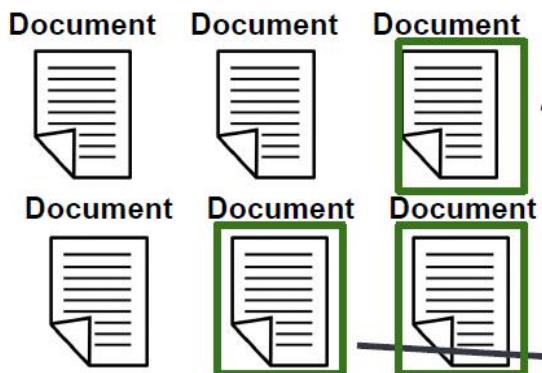
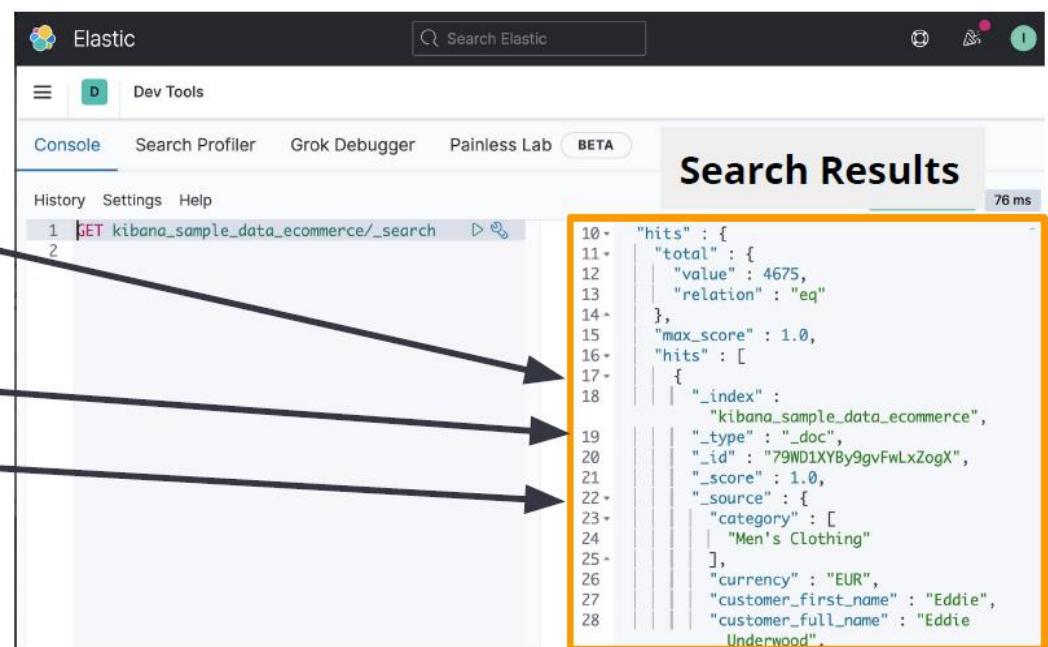


**Documents with similar traits are grouped into an index!**

# Relevance

**When search query is sent, Elasticsearch retrieves relevant documents and presents the documents as search results.**

## Index

The screenshot shows the Elasticsearch Dev Tools interface with the "Search Results" tab selected. The search query is `GET kibana_sample_data_ecommerce/_search`. The results are displayed as a JSON object:

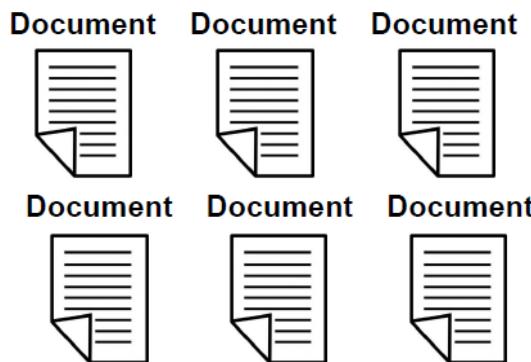
```

10. "hits" : {
11.   "total" : {
12.     "value" : 4675,
13.     "relation" : "eq"
14.   },
15.   "max_score" : 1.0,
16.   "hits" : [
17.     {
18.       "_index" :
19.         "kibana_sample_data_ecommerce",
20.       "_type" : "_doc",
21.       "_id" : "79WD1XYBy9gvFwLxZogX",
22.       "_score" : 1.0,
23.       "_source" : {
24.         "category" : [
25.           "Men's Clothing"
26.         ],
27.         "currency" : "EUR",
28.         "customer_first_name" : "Eddie",
29.         "customer_full_name" : "Eddie Underwood".
30.       }
31.     }
32.   ]
33. }
34. 
```

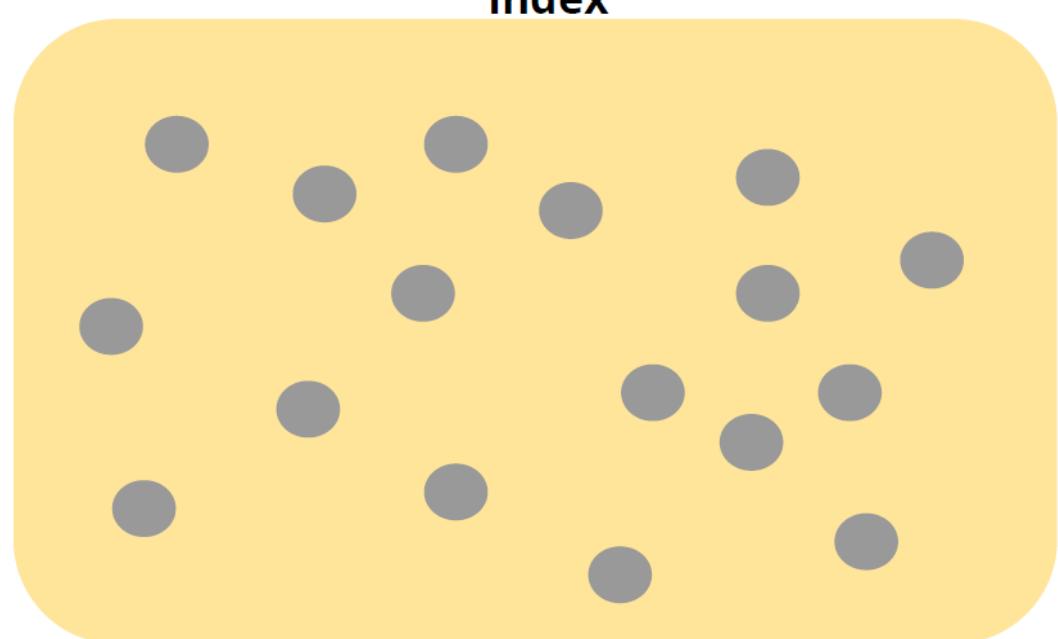
# Relevance

These two diagrams depict the same thing!

Index



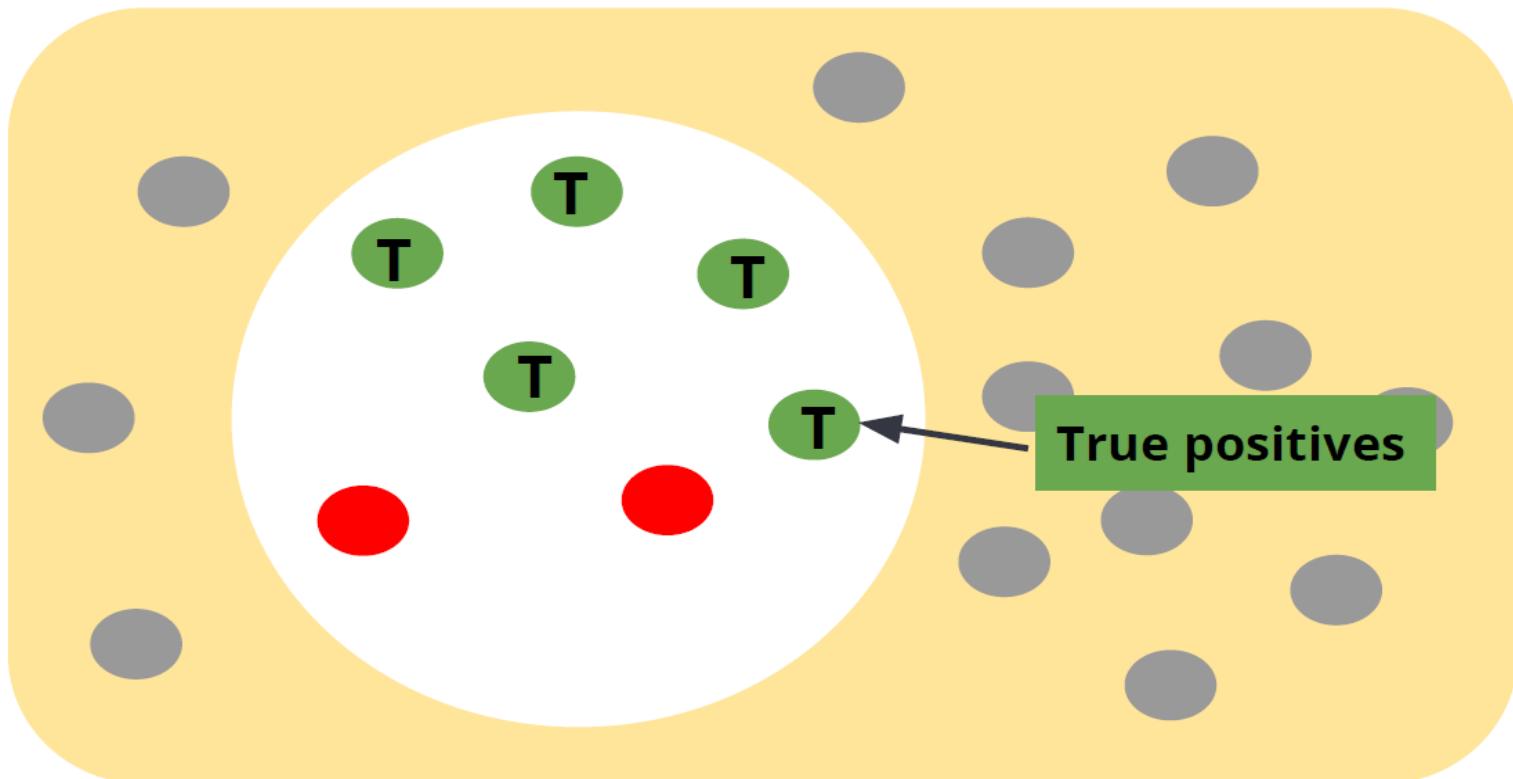
Index



# Relevance

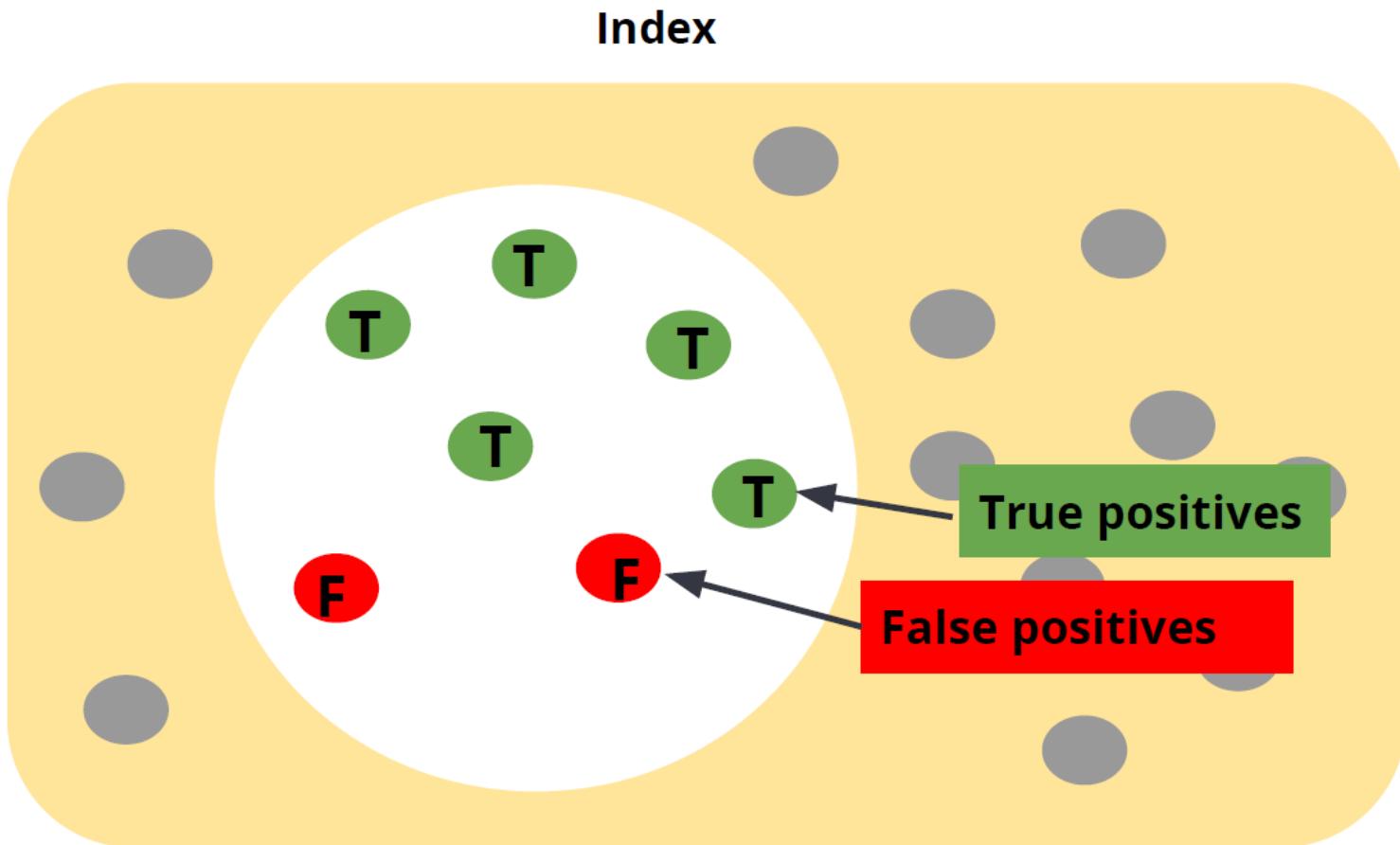
**True positives** are relevant documents that are returned to the user.

Index



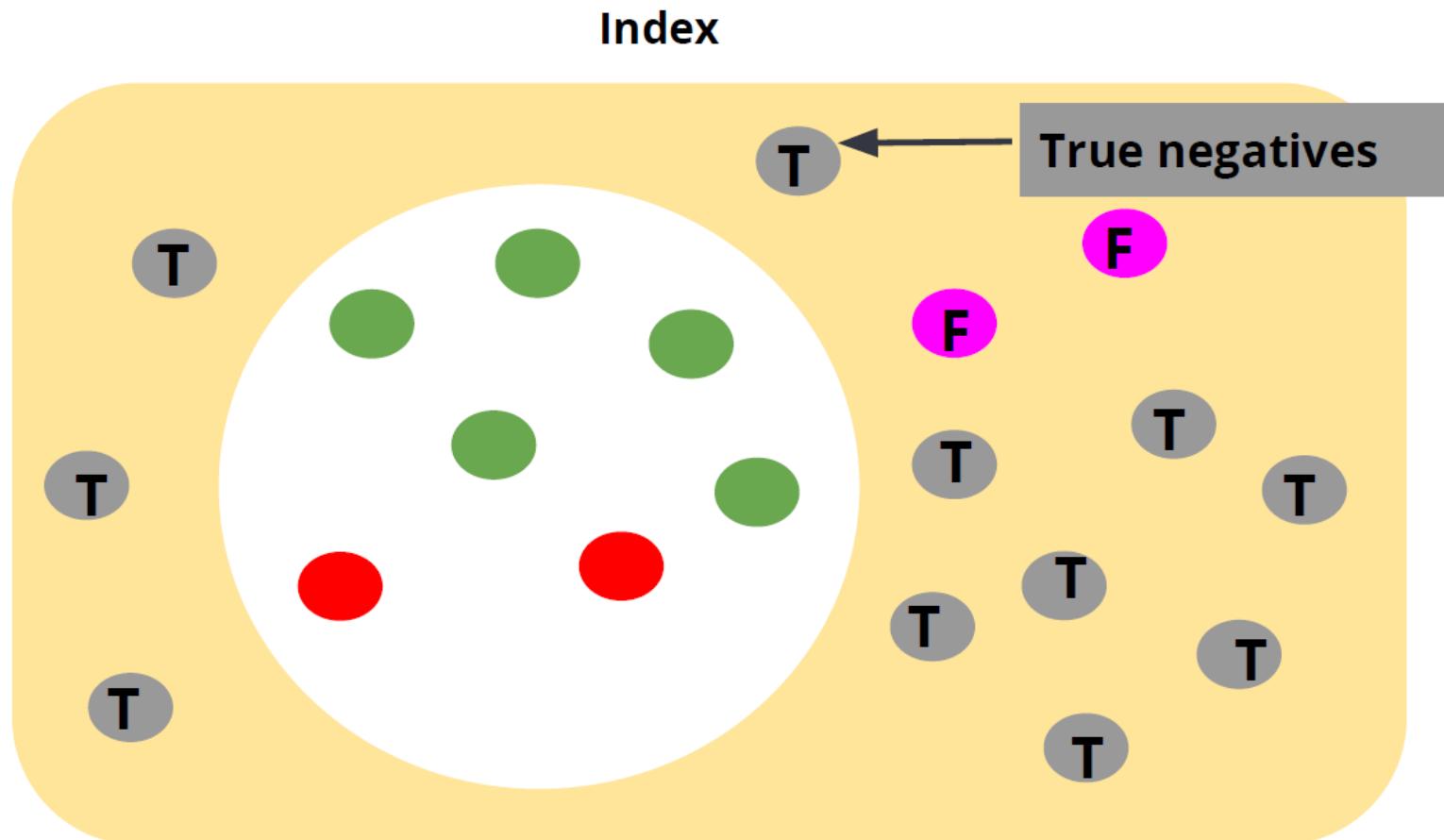
# Relevance

**False positives** are irrelevant documents that are returned to the user.



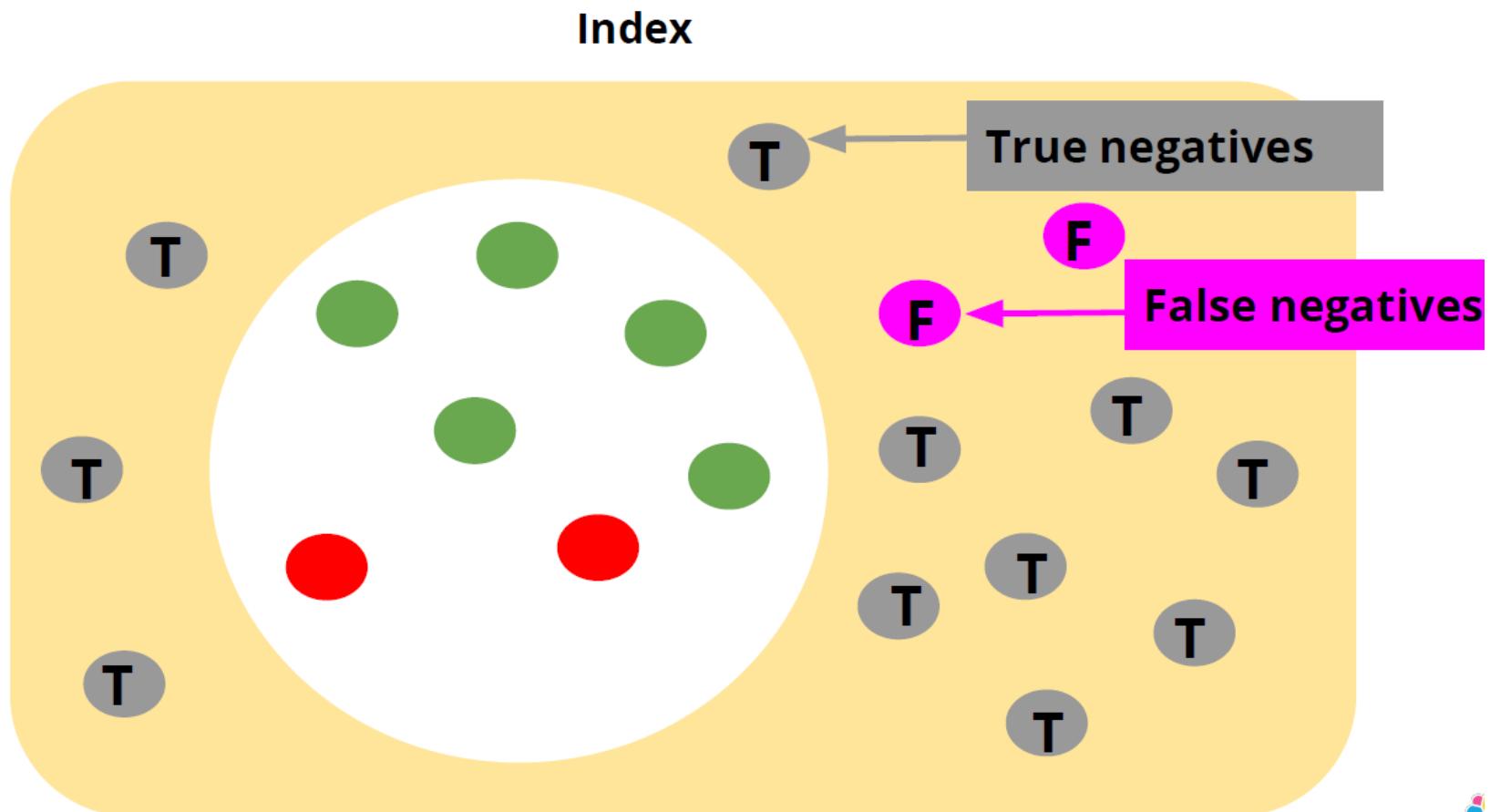
# Relevance

True negatives are irrelevant documents that are not returned to the user.



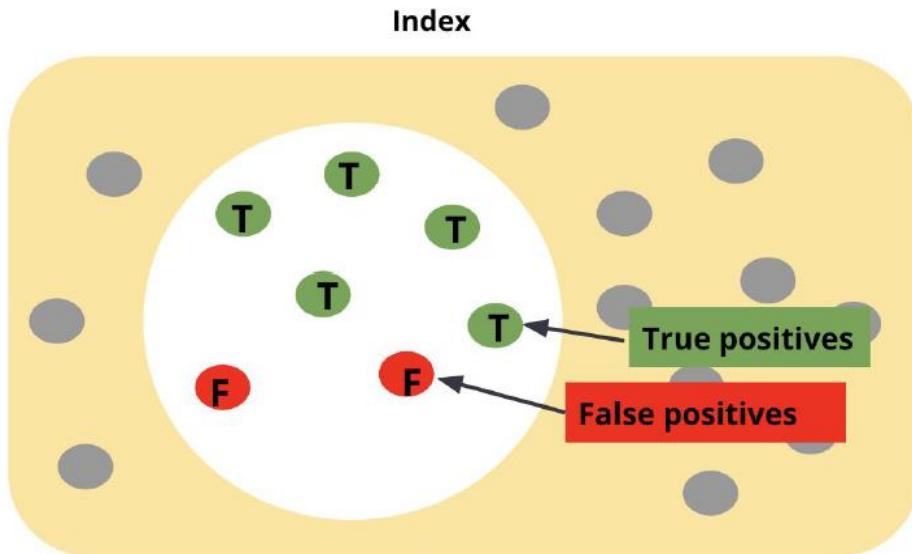
# Relevance

**False negatives** are relevant documents that were not returned to the user.



# Relevance

## What is precision?

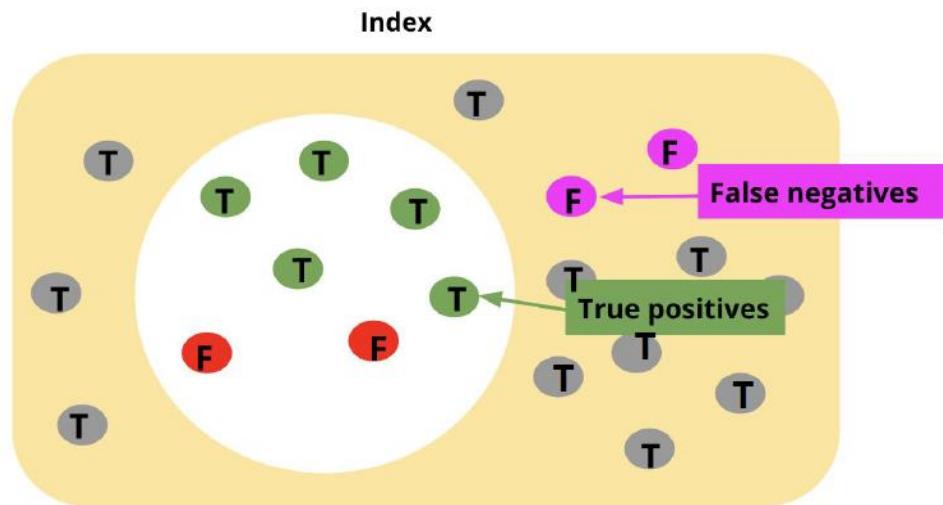


$$\text{Precision} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

What portion of the retrieved data is actually relevant to the search query?

# Relevance

## What is recall?



Recall =

$$\frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

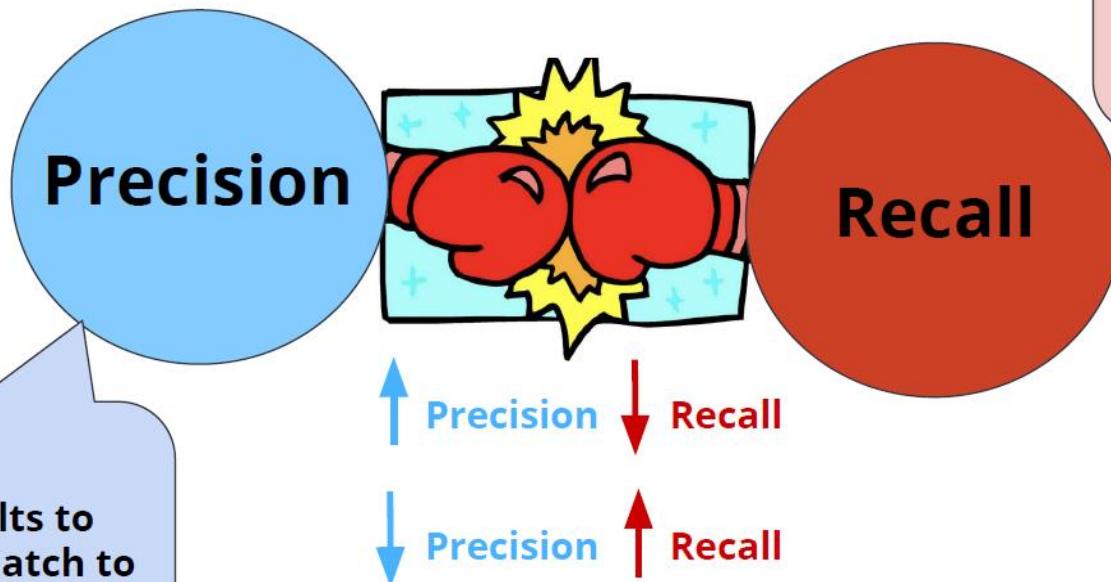
What portion of relevant data is being returned as search results?



DATA SCIENCE

# Relevance

## Precision and Recall are inversely related

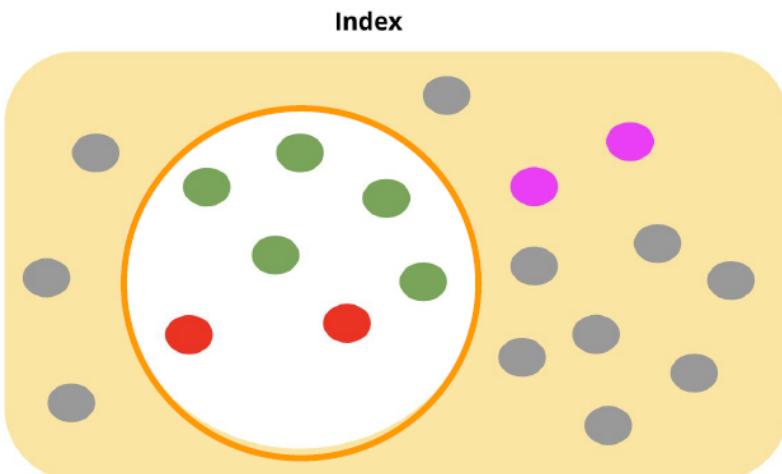


I want all the retrieved results to be a perfect match to the query, even if it means returning less documents.

I want to retrieve more results even if documents may not be a perfect match to the query.

# Relevance

**Precision and recall determine which documents are included in the search results.**



Elastic

Search Elastic

Dev Tools

Console Search Profiler Grok Debugger Painless Lab BETA

History Settings Help

```
1 GET kibana_sample_data_ecommerce/_search
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
```

200 - success 76 ms

```

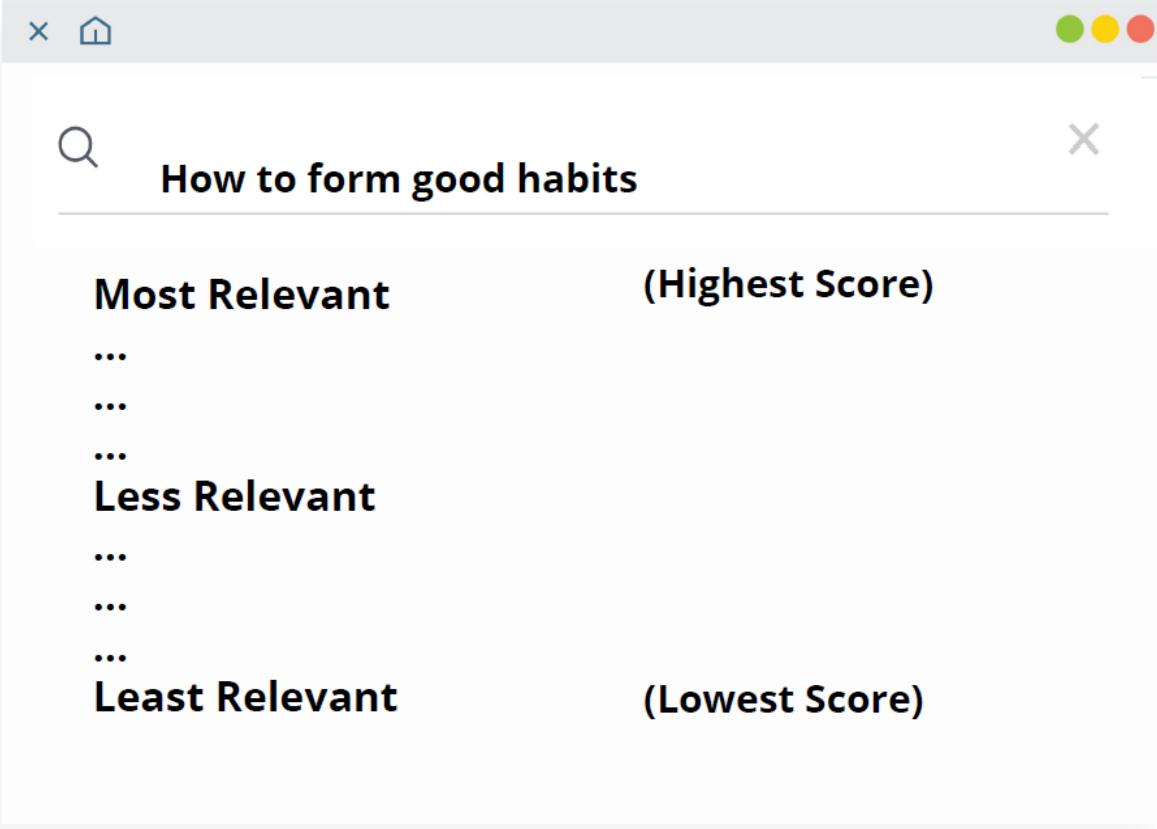
"hits" : {
  "total" : {
    "value" : 4675,
    "relation" : "eq"
  },
  "max_score" : 1.0,
  "hits" : [
    {
      "_index" :
        "kibana_sample_data_ecommerce",
      "_type" : "_doc",
      "_id" : "79WD1XYBy9gvFwLxZogX",
      "_score" : 1.0,
      "_source" : {
        "category" : [
          "Men's Clothing"
        ],
        "currency" : "EUR",
        "customer_first_name" : "Eddie",
        "customer_full_name" : "Eddie Underwood".
      }
    }
  ]
}

```

**Precision and recall do not determine which of the returned documents are more relevant compared to the other!**

# Relevance

**Ranking refers to ordering of the results (from most relevant results at the top, to least relevant at the bottom).**



A screenshot of a search interface with a light gray header bar containing a close button (X), a home icon, and three colored dots (green, yellow, red). The search bar below contains a magnifying glass icon and the text "How to form good habits". To the right of the search bar is another close button (X). The main content area is divided into two columns by a vertical line. The left column is labeled "Most Relevant" at the top and "Less Relevant" in the middle, with three ellipsis ("...") lines between them. The right column is labeled "(Highest Score)" at the top and "(Lowest Score)" in the middle, also with three ellipsis ("...") lines between them. This visualizes how search results are ordered from highest to lowest relevance.

# Score

---

The score is a value that represents how relevant a document is to that specific query

A score is computed for each document that is a hit

# Specificity

---

- Specificity of a classifier is the ratio between how much were correctly classified as negative to how much was negative.
- $\text{Specificity} = \text{TN}/(\text{FP} + \text{TN})$
- Where is it used?
- Places where classification of negatives are high priority.
- Eg: Diagnosing for a health condition before treatment.

# Sensitivity

---

- Sensitivity of a classifier is the ratio between how much were correctly identified as positive to how much were positive.
- $\text{Sensitivity} = \text{TP} / (\text{FN} + \text{TP})$
- Where is it used ?
- Places where classification of positives are high priority.
- Eg: Security checks in airports.

# Specificity

---

- Specificity of a classifier is the ratio between how much were correctly classified as negative to how much was negative.
- $\text{Specificity} = \text{TN}/(\text{FP} + \text{TN})$
- Where is it used?
- Places where classification of negatives are high priority.
- Eg: Diagnosing for a health condition before treatment.

# Precision

---

- How much were correctly classified as positive out of all positives.
- Precision =  $TP/TP+FP$

# Recall

---

- Recall and sensitivity are one and the same.

$$\text{Recall} = \text{TP} / (\text{FN} + \text{TP})$$

# Where does precision and recall are used ?

- The harmonic mean of precision and recall gives a score call f1 score which is a measure of performance of the model's classification ability.
- $F1 \text{ score} = 2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$
- F1 score is considered a better indicator of the classifier's performance than the regular accuracy measure.

## AUC-ROC curve

---

- ROC stands for Receiver Operating Characteristics, and the ROC curve is the graphical representation of the effectiveness of the binary classification model.
- It plots the true positive rate (TPR) vs the false positive rate (FPR) at different classification thresholds.

# AUC-ROC curve

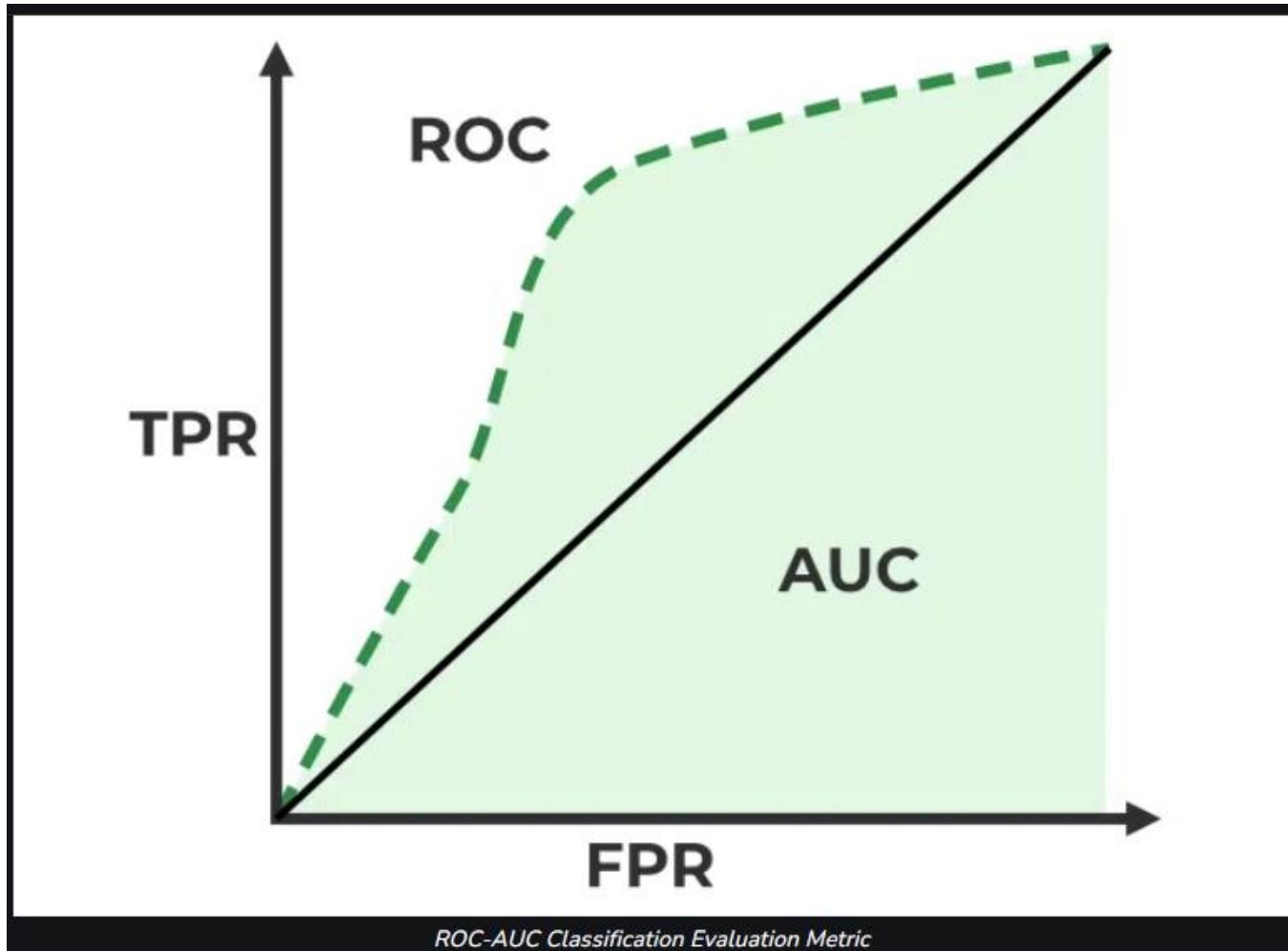
- AUC stands for Area Under the Curve, and the AUC curve represents the area under the ROC curve.
- It measures the overall performance of the binary classification model.
- As both TPR and FPR range between 0 to 1, So, the area will always lie between 0 and 1, and A greater value of AUC denotes better model performance.
- Our main goal is to maximize this area to have the highest TPR and lowest FPR at the given threshold.
- The AUC measures the probability that the model will assign a randomly chosen positive instance a higher predicted probability compared to a randomly chosen negative instance.

# AUC-ROC curve

---

- TPR and FPR
- This is the most common definition that you would have encountered when you would Google AUC-ROC
- Basically, the ROC curve is a graph that shows the performance of a classification model at all possible thresholds( threshold is a particular value beyond which you say a point belongs to a particular class).
- The curve is plotted between two parameters
  - TPR – True Positive Rate
  - FPR – False Positive Rate

# AUC-ROC curve



# Absolute Error Definition

- Absolute error is the difference between measured or inferred value and the actual value of a quantity.
- The absolute error is inadequate since it does not give any details regarding the importance of the error.
- While measuring distances between cities kilometers apart, an error of a few centimeters is negligible and is irrelevant.
- Consider another case where an error of centimeters when measuring small machine parts is a very significant error.
- Both the errors are in the order of centimeters, but the second error is more severe than the first one.

# Absolute Error Definition

- For example, 24.13 is the actual value of a quantity and 25.09 is the measure or inferred value, then the absolute error will be:
- Absolute Error =  $25.09 - 24.13$
- = 0.86
- Most of the time it is sufficient to record only two decimal digits of the absolute error.
- Thus, it is sufficient to state that the absolute error of the approximation 4.55 to the correct value 4.538395 is 0.012.

# Relative Error

- The relative error is defined as the ratio of the absolute error of the measurement to the actual measurement.
- Using this method, we can determine the magnitude of the absolute error in terms of the actual size of the measurement.
- If the true measurement of the object is not known, then the relative error can be found using the measured value.
- The relative error gives an indication of how good measurement is relative to the size of the object being measured.

# Relative Error

- If  $x$  is the actual value of a quantity,  $x_0$  is the measured value of the quantity and  $\Delta x$  is the absolute error, then the relative error can be measured using the below formula.
- Relative error =  $(x_0-x)/x = (\Delta x)/x$
- An important note that relative errors are dimensionless.
- When writing relative errors, it is usual to multiply the fractional error by 100 and express it as a percentage.

# Relative Error

---

- Example 1:
- Find the absolute and relative errors of the approximation 125.67 to the value 119.66.
- Solution:
- Absolute error =  $|125.67 - 119.66| = 6.01$
- Relative error =  $|125.67 - 119.66| / 119.66 = 0.05022$

# Mean Absolute Error

- The mean absolute error is the average of all absolute errors of the data collected.
- It is abbreviated as MAE (Mean Absolute Error).
- It is obtained by dividing the sum of all the absolute errors with the number of errors.
- The formula for MAE is:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

Here,

$|x_i - \bar{x}|$  = absolute errors

n = number of errors