# Slot 01 - Introduction to Text Mining

Presenter:

Dr. LE Thanh Tung
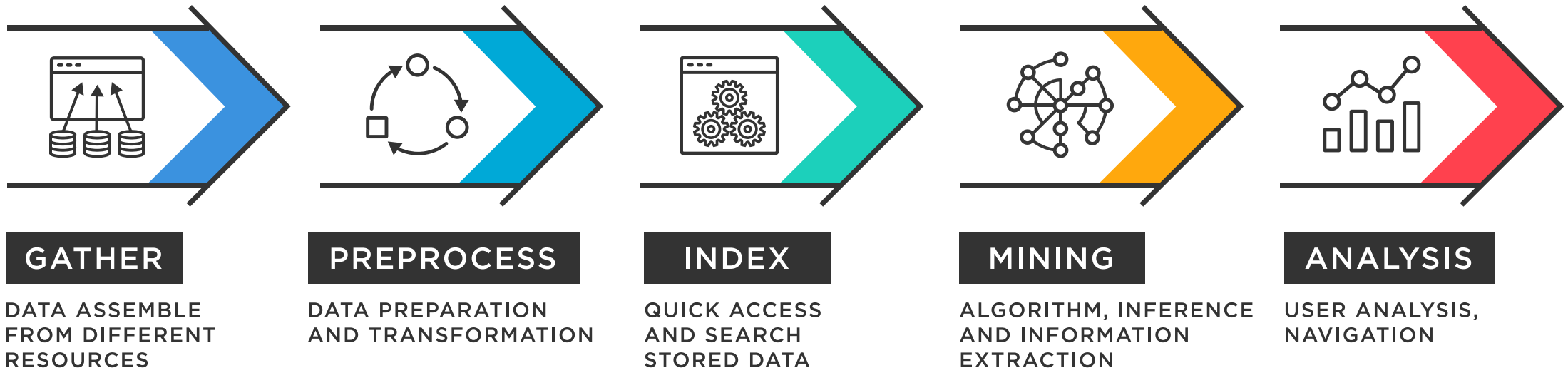
# Content

**1** Text Mining

**2** Machine Learning

**3** Evaluation Metrics
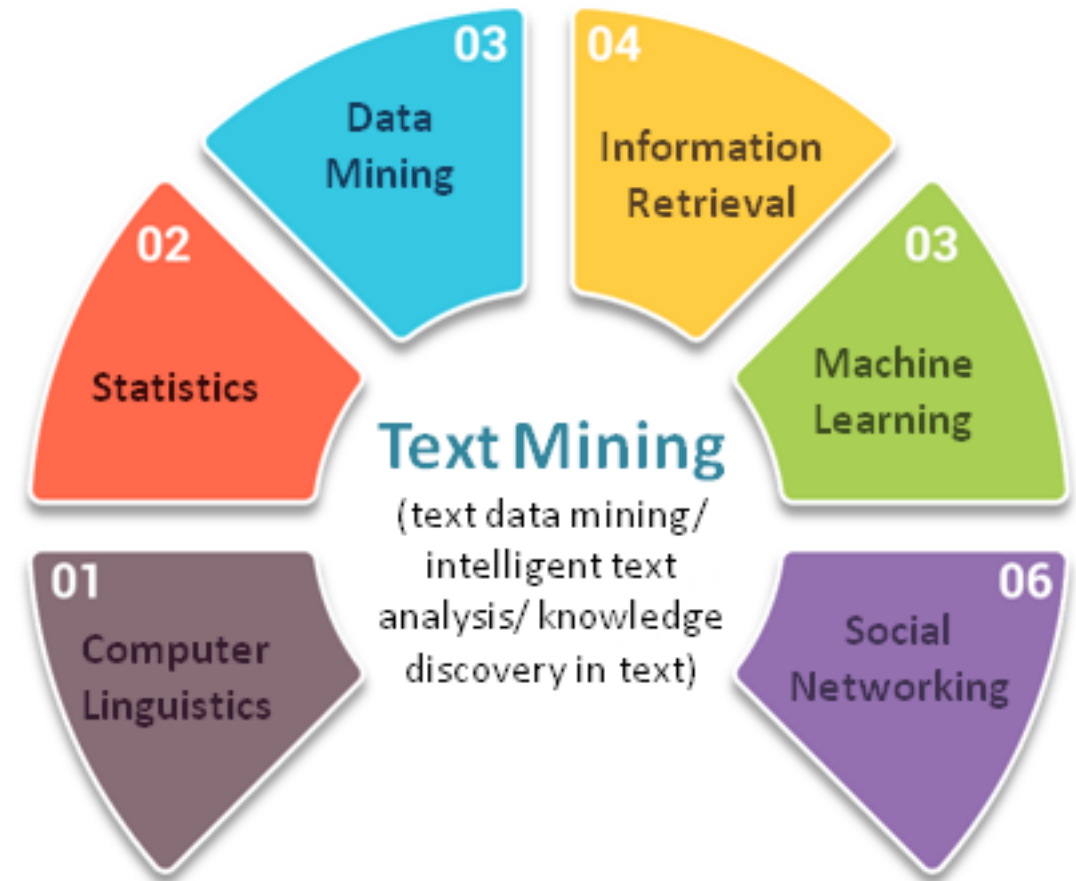
**4** Programming Language

# Introduction

- Text mining is the process of transforming unstructured **text** into a structured format to identify meaningful patterns and new insights

TEXT MINING INVOLVES A SERIES OF ACTIVITIES TO BE PERFORMED IN ORDER TO EFFICIENTLY MINE THE INFORMATION. THESE ACTIVITIES ARE:
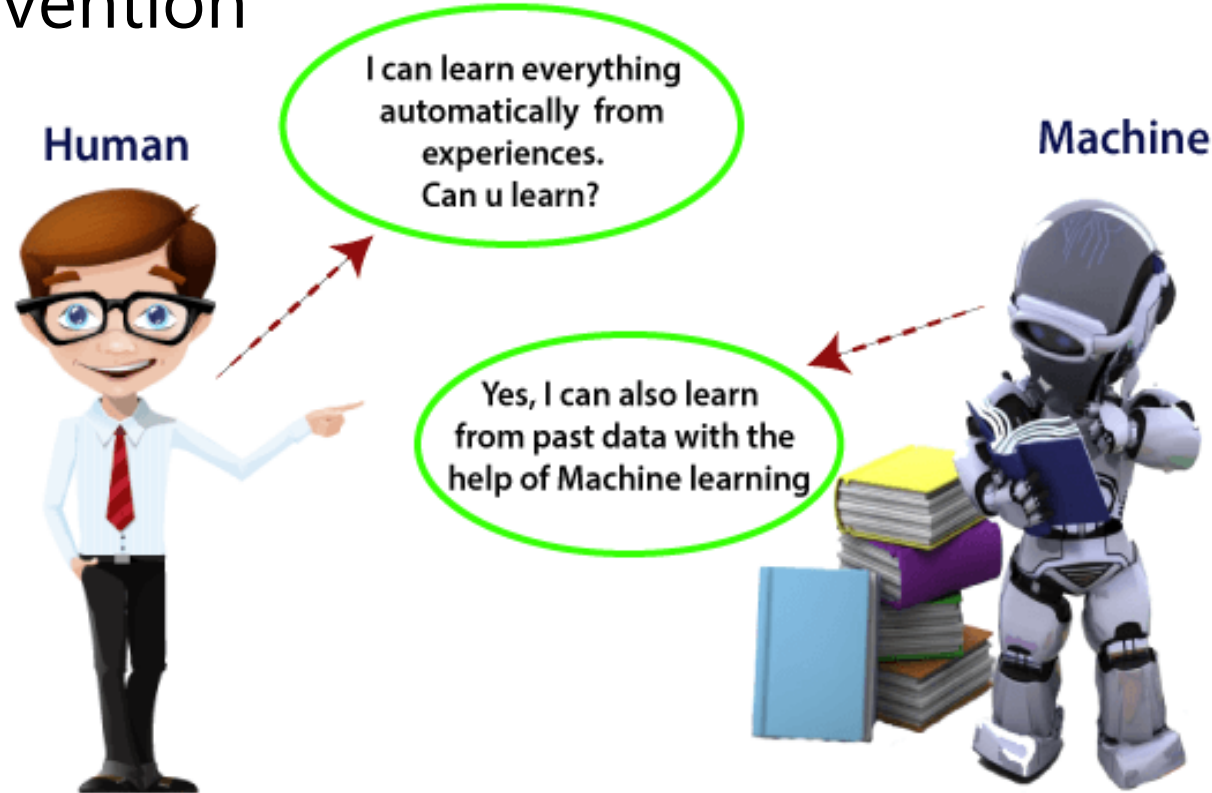
**GATHER**

DATA ASSEMBLE FROM DIFFERENT RESOURCES

**PREPROCESS**

DATA PREPARATION AND TRANSFORMATION

**INDEX**

QUICK ACCESS AND SEARCH STORED DATA

**MINING**

ALGORITHM, INFERENCE AND INFORMATION EXTRACTION

**ANALYSIS**

USER ANALYSIS, NAVIGATION

# Introduction

■ Gupta & Lehal (2009) have regarded text mining as new interdisciplinary area which is an amalgamation of data mining, information retrieval, machine learning, computer linguistic and statistics
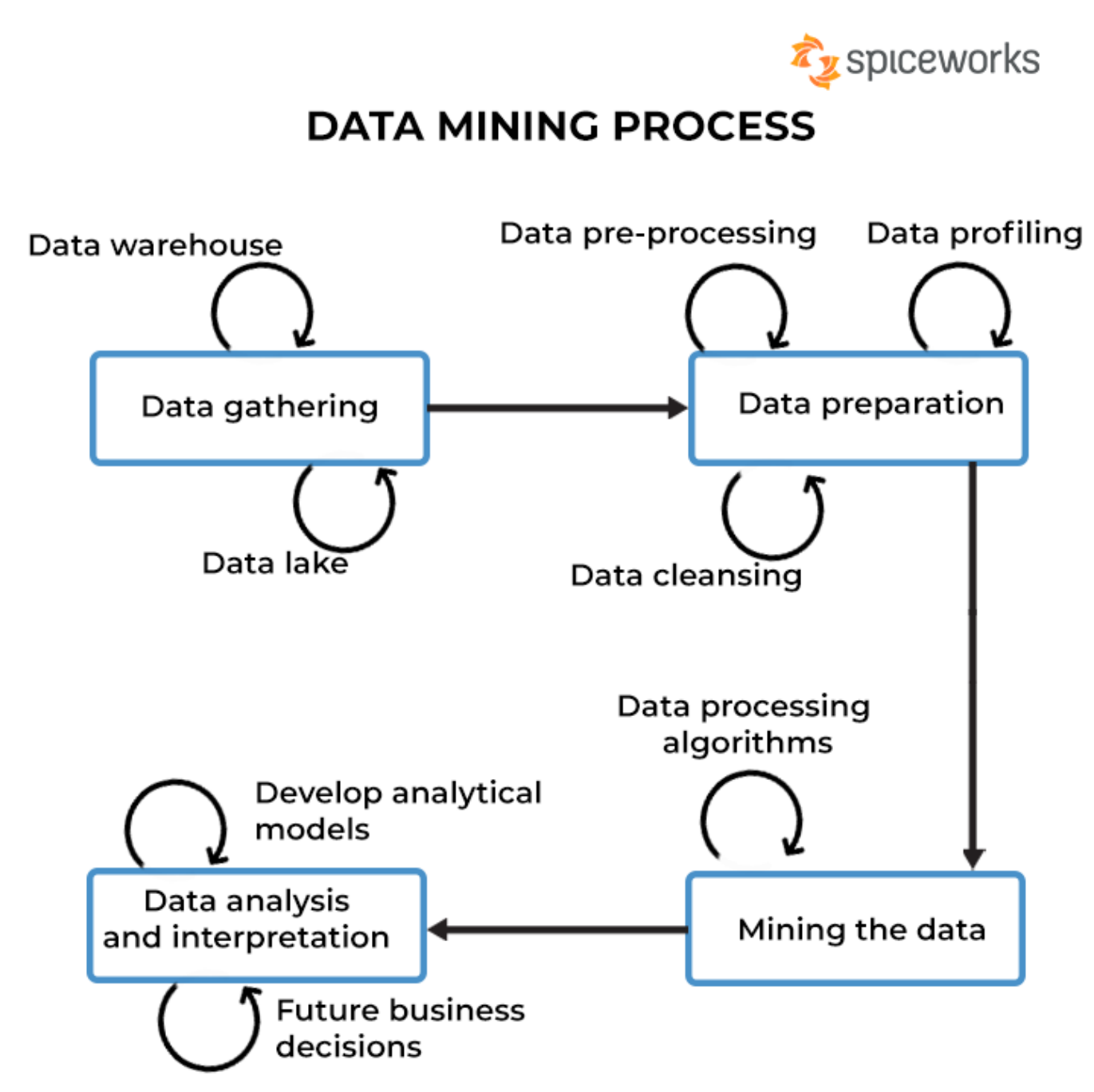
- Machine Learning is a discipline of artificial intelligence (AI) that provides machines with the ability to automatically **learn from data and past experiences** while identifying patterns to make predictions with minimal human intervention
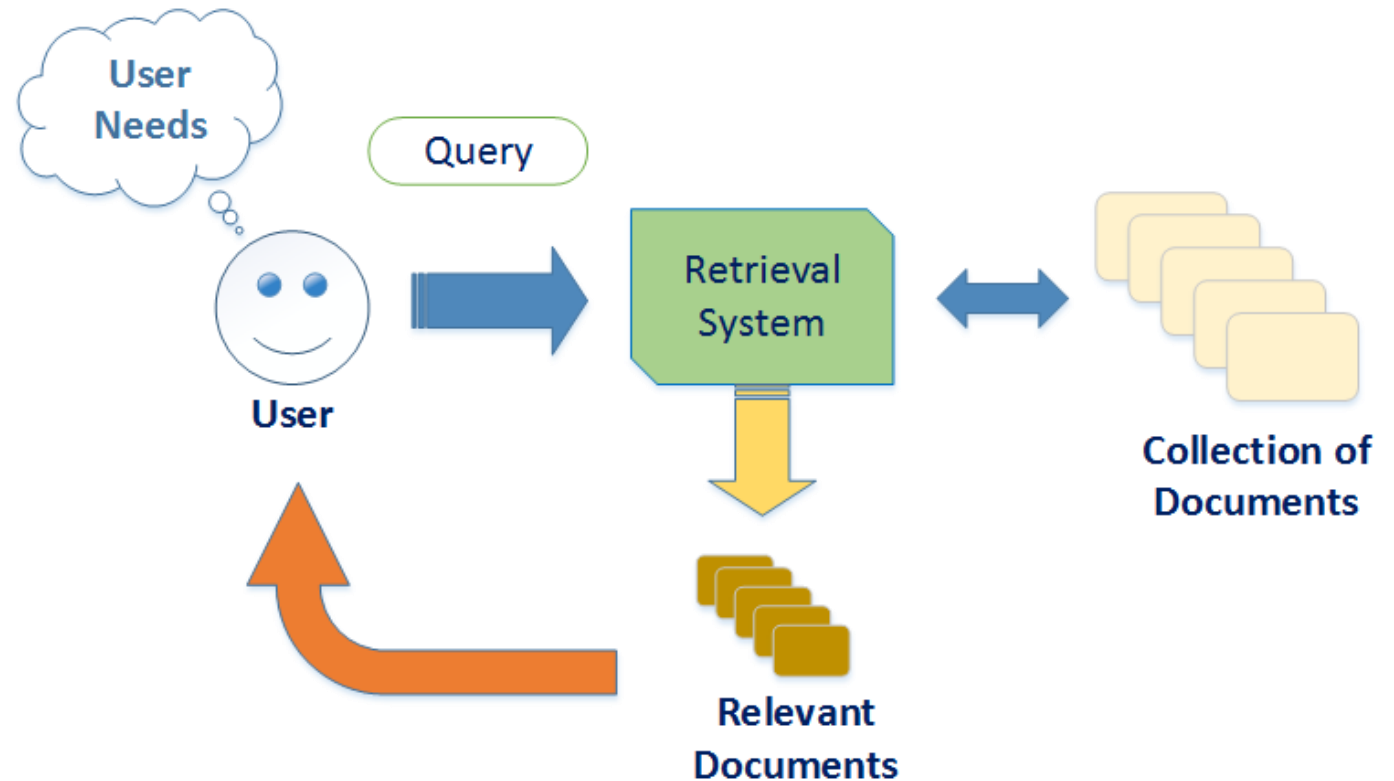
# Introduction

- Data mining is the process of analyzing a large batch of information to **discern trends and patterns**

## DATA MINING PROCESS

# Introduction

- Information Retrieval is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers)

- Two main approaches are matching words in the query against the database index (keyword searching) and traversing the database using hypertext or hypermedia links

DATA MINING **V/S** TEXT MINING

| DATA MINING | TEXT MINING |
|---|---|
| Technique of processing raw data in a structured form | Technique of processing of text from documents |
| Data is stored in structured format | Data is stored in unstructured format |
| Processing of data is done directly | Processing of data is done linguistically |
| Easy to retrieve data as it is homogeneous | Not so easy to retrieve data as it is heterogeneous |
| Areas of uses – fraud detection, medicine, healthcare etc. | Areas of uses – online reviews, customer surveys etc. |

SKILL<S/ASH>

# Text Analytics Application



**Manufacturers**
- Identify root causes of product issue quicker
- Identify trends in market segments
- Understand competitors products

**Government**
- Identify fraud
- Understand public sentiments about unmet needs
- Find emerging concerns that can shape policy

**Financial Institutions**
- Use contact center transcriptions
- Understand customers
- Identify money laundering or other fraudulent situation

**Retail**
- Identify profitable customers and understand the reasons for their loyalty
- Manage the brand on social media

**Legal**
- Identify topics and keywords in discovery documents
- Find patterns in defendant's communications

**Healthcare**
- Find similar patterns in doctor's reports
- Use social media to detect outbreaks earlier
- Identify patterns in patient claims data

**Telecommunications**
- Prevent customer churn
- Suggest up-sell/cross-sell opportunities by understanding customer comments

**Life Sciences**
- Identify adverse events in medicines or vaccines
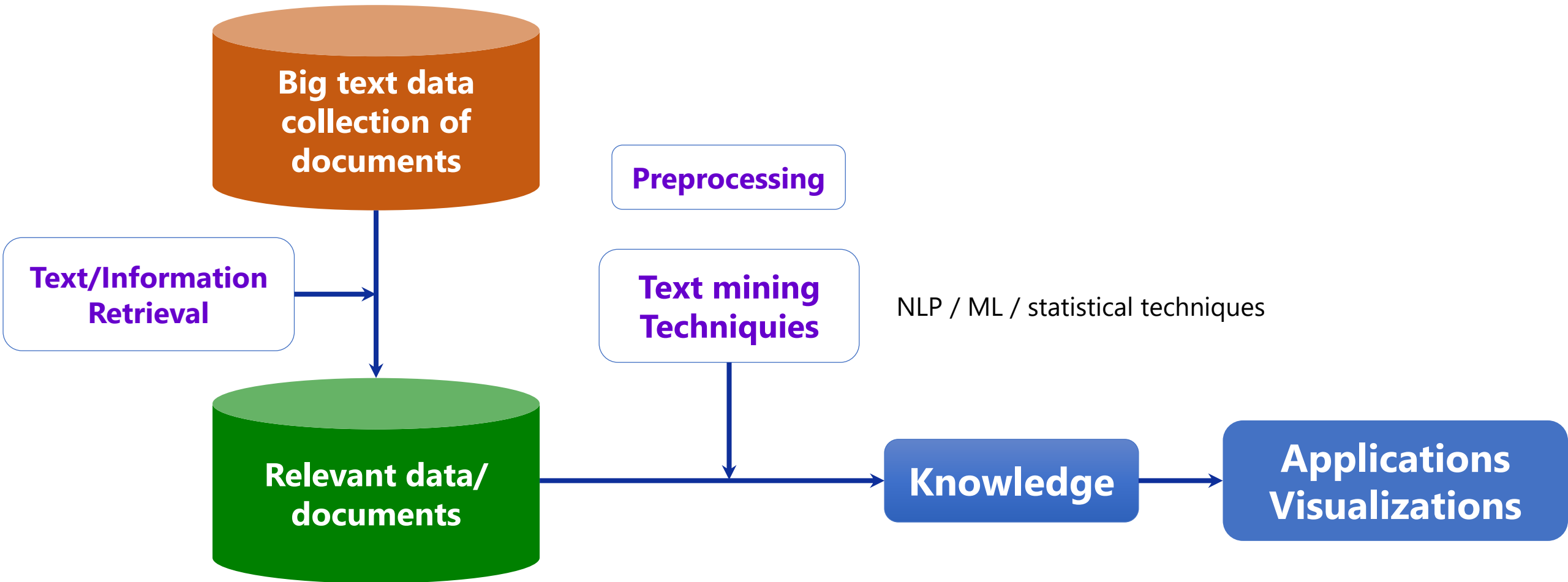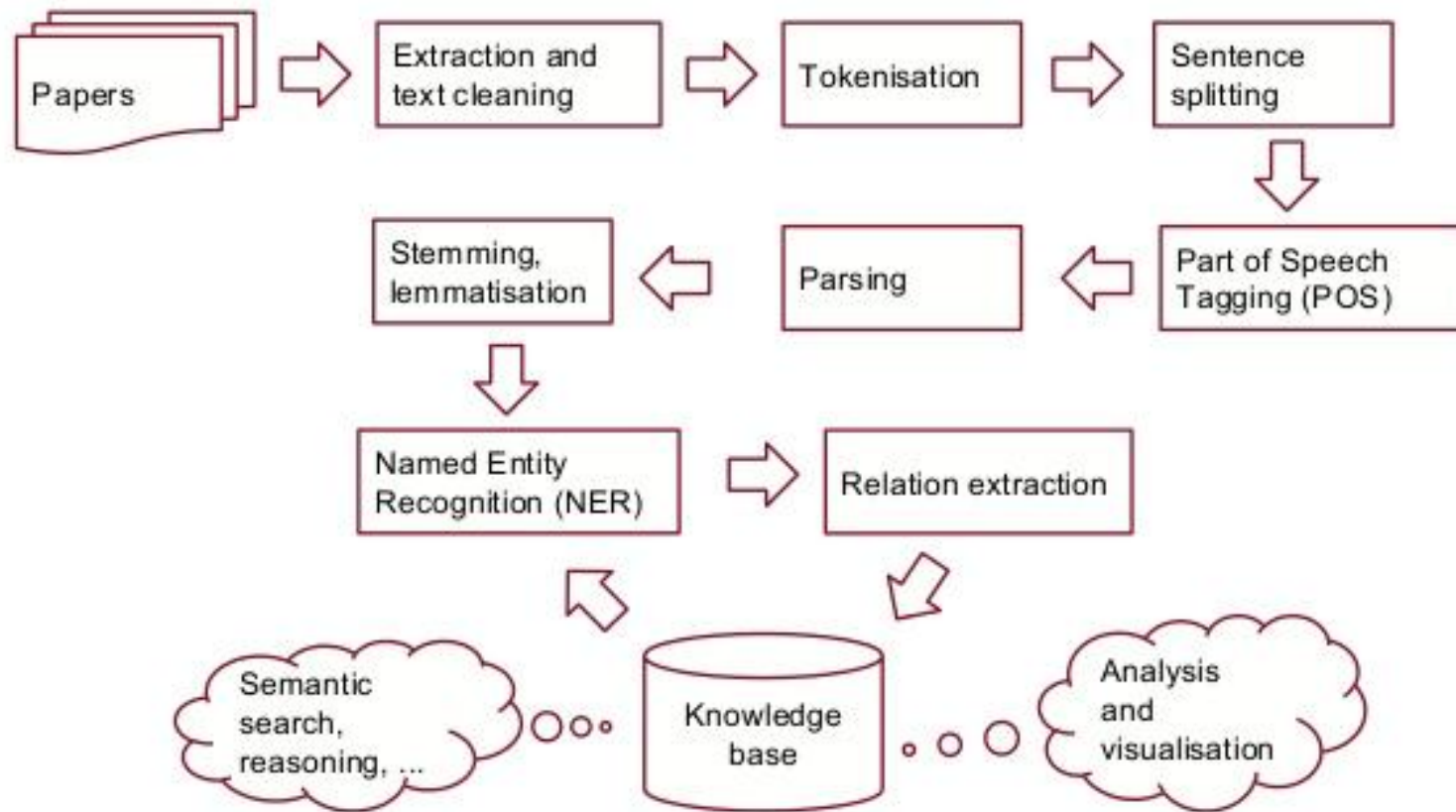- Recommend appropriate research materials

**Insurance**
- Identify fraudulent claims
- Track competitive intelligence
- Manage the brand on social media

zencos

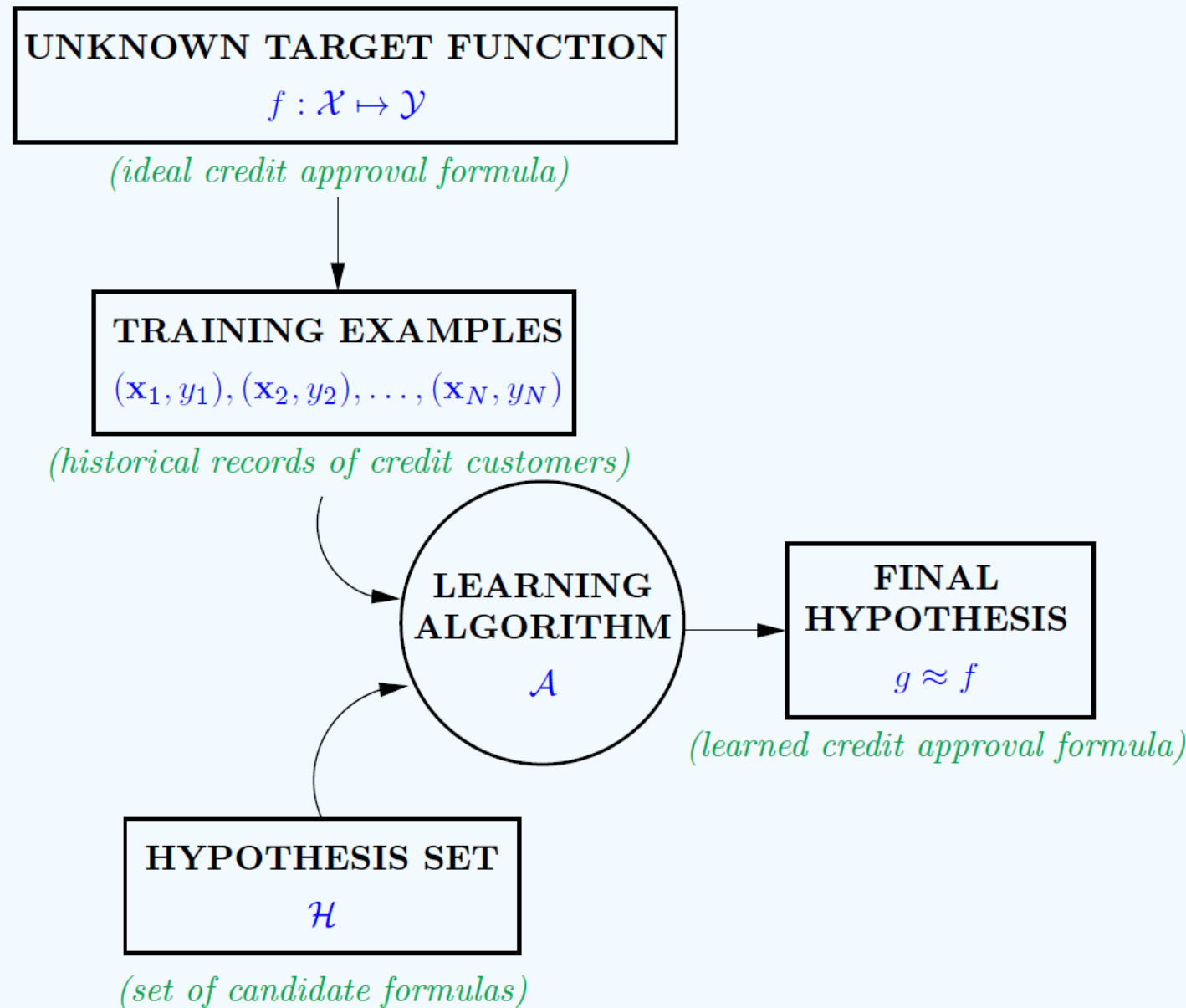# General Architecture

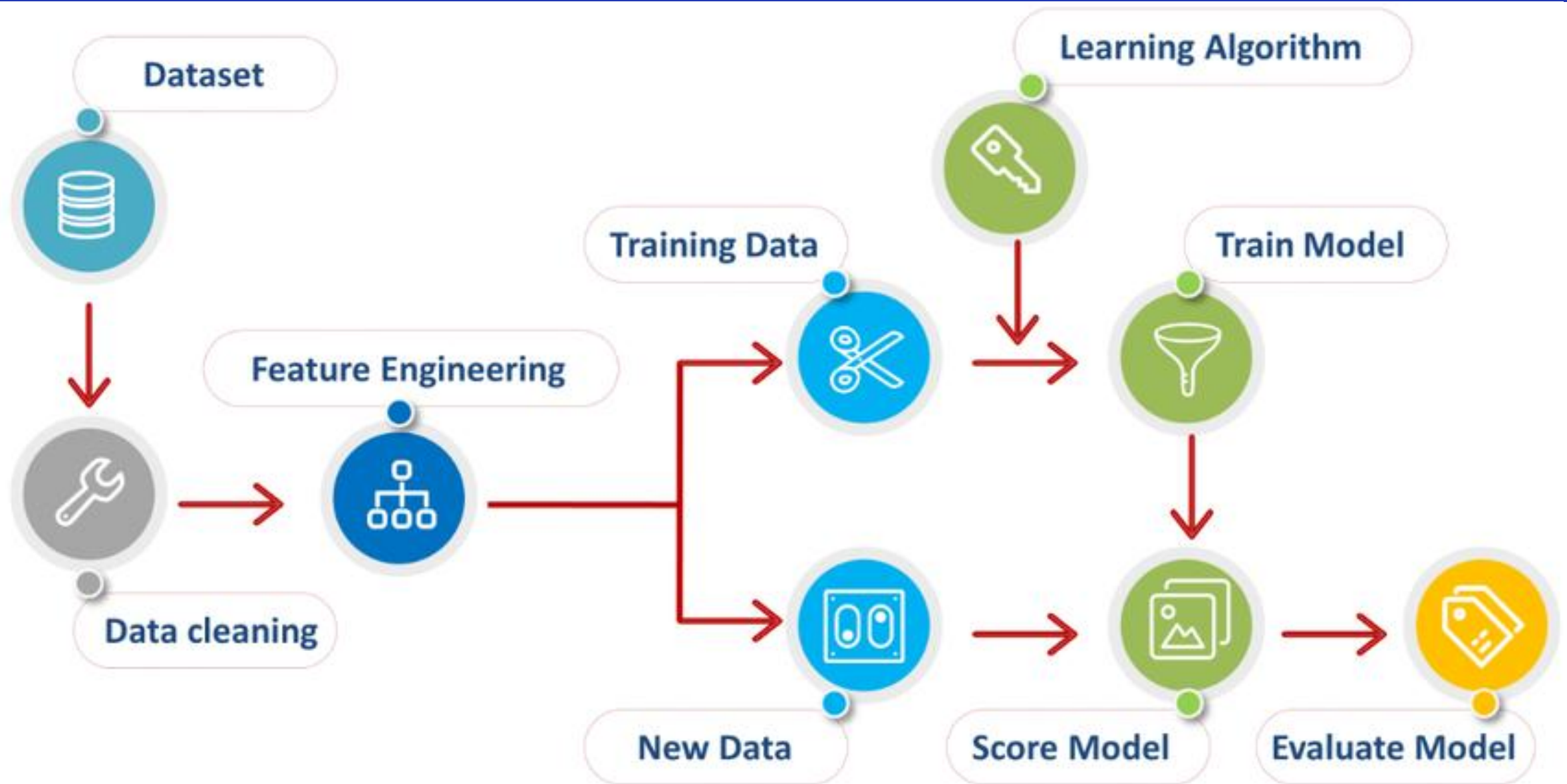Generic text mining workflow

- Machine learning tasks
  - Build a **model** from some **data**
    - choose how to map raw data to feature vectors
    - choose a model form
    - choose parameter values in the model
  - **Test** or **validate** the model
    - Evaluate the model on unseen data to assess its performance

# Learning Algorithms

- The data is formed by a specific distribution.

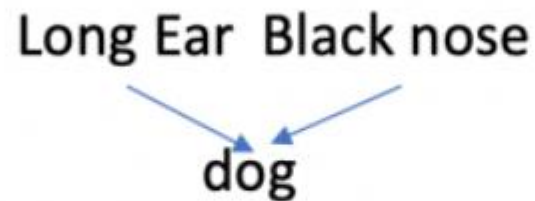- The goal of learning is to find the data distribution which is generalized in this kind of data



UNKNOWN TARGET FUNCTION
$f : \mathcal{X} \mapsto \mathcal{Y}$

*(ideal credit approval formula)*

TRAINING EXAMPLES
$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_N, y_N)$

*(historical records of credit customers)*

LEARNING ALGORITHM
$\mathcal{A}$

FINAL HYPOTHESIS
$g \approx f$

*(learned credit approval formula)*

HYPOTHESIS SET
$\mathcal{H}$

*(set of candidate formulas)*

- Feature

  - A measurable property or characteristic of phenomenon being observed.

- Feature

  - A measurable property or characteristic of phenomenon being observed.



Machine Learning:

Sample → Label

dog     cat     horse

- Feature
  - A measurable property or characteristic of phenomenon being observed.

- Feature

  - A measurable property or characteristic of phenomenon being observed.
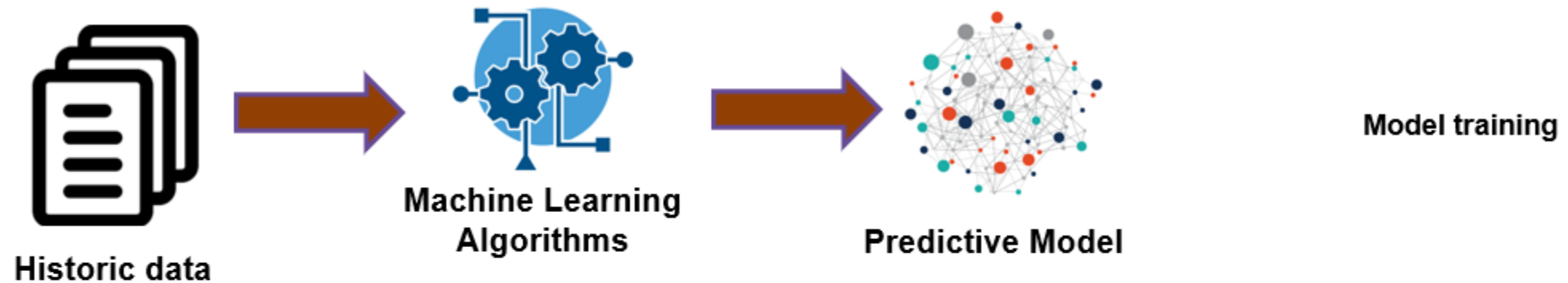
| | Color | Weight | my Rank | ... |
|---|---|---|---|---|
|  | Red | 200g | 1st | |
|  | Yellow | 300g | 3rd | |

# Predictive Tasks

- Build a Machine Learning model:

  - Training

  - Evaluating

  - Testing

# Learning Algorithms

- The goal of learning is to find the hidden pattern/distribution from historic data

# Type of Machine Learning

# Type of Machine Learning

**Supervised**

**Unsupervised**

**Reinforcement**

**Input:**

- Features

- Gold labels

**Tasks:** task-oriented

- Classification

- Regression

**Input:**

- Features

**Tasks:** data-oriented

- Clustering

- Dimension Reduction

**Input:**

- Features

- Feed-back function

**Tasks:** environment

- Real-time Decision

- Learning tasks

# Type of Machine Learning



**Supervised**

**Unsupervised**

**Reinforcement**

- Learn from **labeled** data
  - Classification methods: predict a label
  - Regression methods: predict a quantity



**Regression**
What is the temperature going to be tomorrow?

**Classification**
Will it be Cold or Hot tomorrow?

- Handwriting recognition:

  - Data: a set of handwriting images and their labels.

  - Goal: build a model to predict a digit for a given image.

- Medical Costs Analysis:

    - Data: containing medical records

    - Goal: Build a model to predict medical costs for a given individual based on their demographic and health-related information

|   | age | sex | bmi | children | smoker | region | charges |
|---|-----|-----|-----|----------|--------|--------|---------|
| 0 | 19 | female | 27.900 | 0 | yes | southwest | 16884.92400 |
| 1 | 18 | male | 33.770 | 1 | no | southeast | 1725.55230 |
| 2 | 28 | male | 33.000 | 3 | no | southeast | 4449.46200 |
| 3 | 33 | male | 22.705 | 0 | no | northwest | 21984.47061 |
| 4 | 32 | male | 28.880 | 0 | no | northwest | 3866.85520 |

- Learn from unlabeled data
    - Dimension reduction
    - Clustering
    - Association



**Association Rule Learning**

*"93% of people who purchased item A also purchased item B"*

- Learn from mistakes

# Feature Extraction

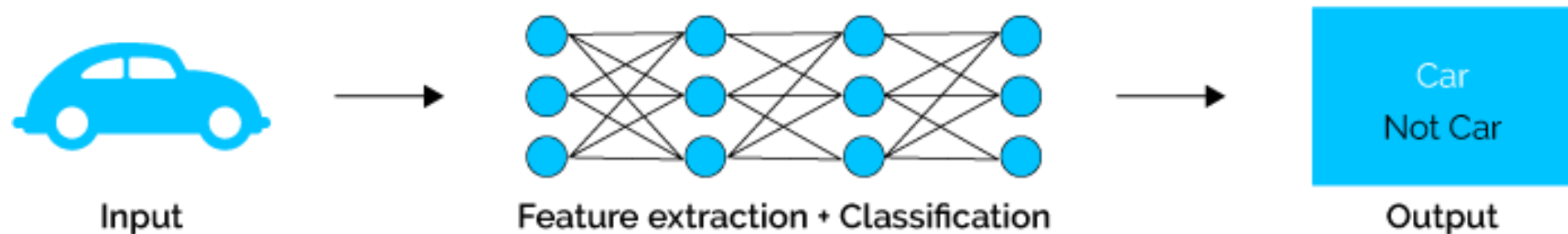- Represent the data into the "learnable" object in computer



- Image: pixels, rules, ...

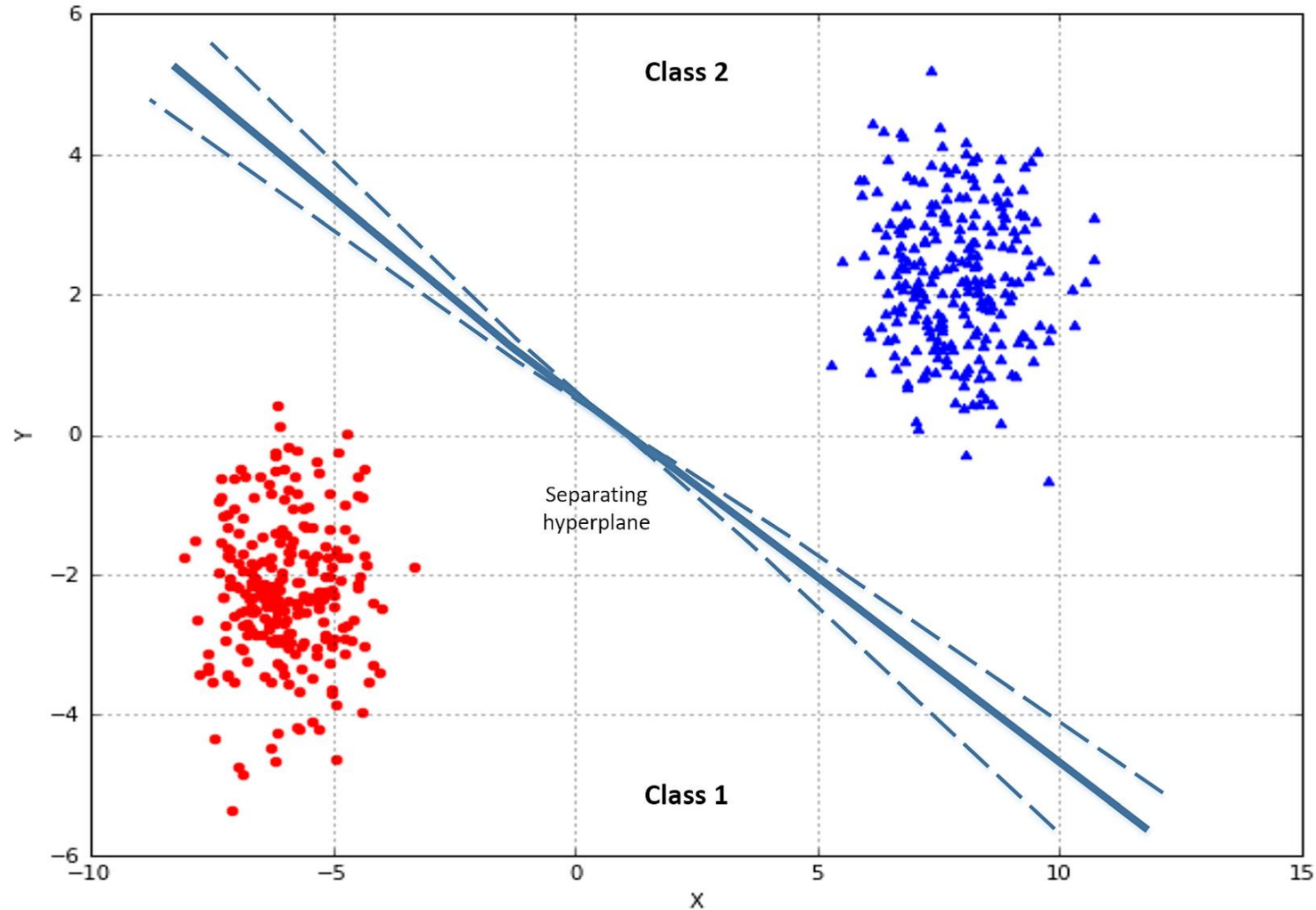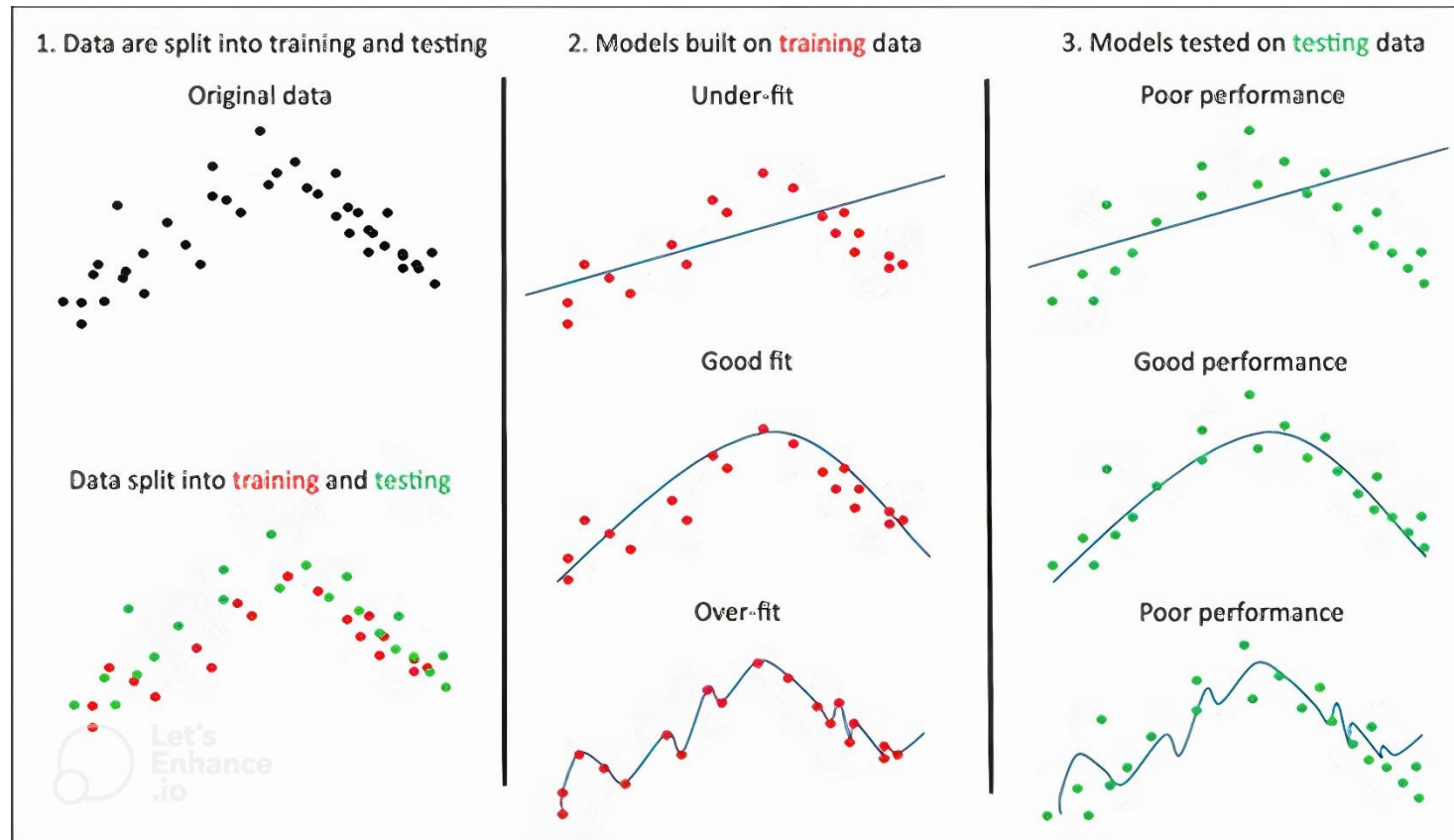- Audio: sampling, quantization

- Text: ????

# Feature Extraction

- How to choose the best "suitable" model

# Evaluation

- Generalization refers to your model's ability to adapt properly to new, previously **unseen data**, drawn from the same distribution as the one used to create the mode

# Evaluation

- Based on some metrics that compares the model's predicted labels with the known true labels

| Regression | Classification | Recommender System |
|---|---|---|
| • Mean Absolute Error (MAE)<br>• Root Mean Squared Error (RMSE)<br>• R-Squared and Adjusted R-Squared | • Recall<br>• Precision<br>• F1-Score<br>• Accuracy<br>• Area Under the Curve (AUC) | • Mean Reciprocal Rank<br>• Root Mean Squared Error (RMSE) |

# Evaluation: Classification

- A confusion matrix is a performance measurement technique

- Accuracy is a metric that generally describes how the model performs across all classes.

$$Accuracy = \frac{Correct\ prediction}{Total\ cases} * 100\%$$

$$Accuracy = \frac{(TP + TN)}{(TP + TN + FP + FN)} * 100\%$$

Actual Values

| | Positive | Negative |
|---|---|---|
| **Positive** | True Positive | False Positive |
| **Negative** | False Negative | True Negative |

Predicted Values

# Evaluation: Classification

- Precision:

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

$$= \frac{True\ Positive}{Total\ Predicted\ Positive}$$

- Recall:

$$\text{Recall} = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

$$= \frac{True\ Positive}{Total\ Actual\ Positive}$$

Actual Values

|  | Positive | Negative |
|---|---|---|
| **Positive** | True Positive | False Positive |
| **Negative** | False Negative | True Negative |

Predicted Values

# Data Splitting Techniques

- Kinds of data:

  - Train: learn from the historical patterns to reveal the hidden models

  - Validating: choose the best/good model, tune the model's hyperparameters and configurations

  - Testing: evaluate the practical performance

# Data Splitting Techniques

- Random sampling

  - Simple

  - Stratified: random from each group (group division by label)

  - Cluster: group division by features

  - Multistage: divide by many kinds of random sampling

# Data Splitting Techniques

- Cross-validation:
    - Cross over training data
    - Cross over all data

- This alleviates any bias occurring as selecting data in the training and validation sets.

- The goal of this technique is validating and testing, not training model.

- If an algorithm performs well on a certain class of problems, then it necessarily pays for that with degraded performance on the set of all remaining problems. (Wolpert and Macready)
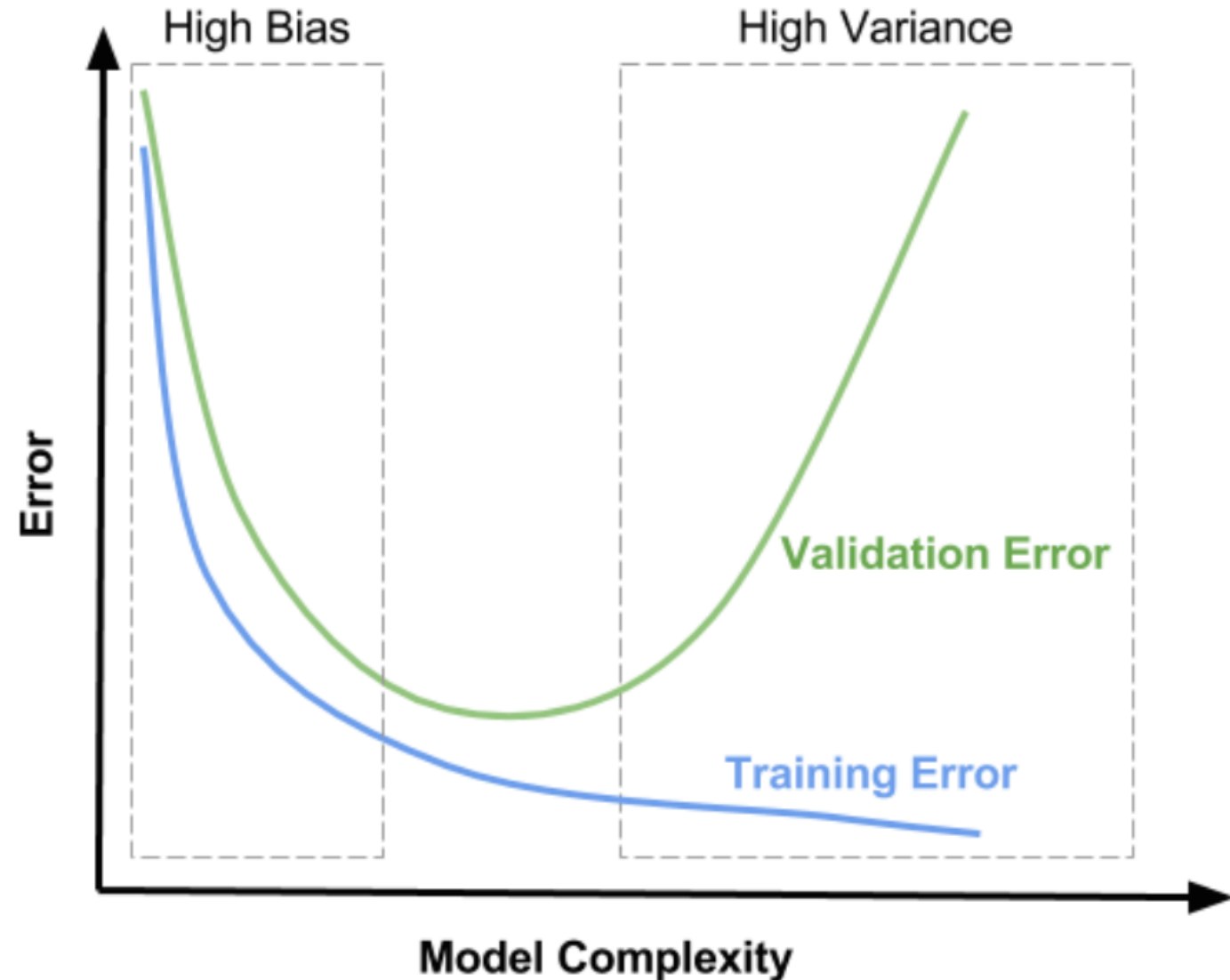
- There is **no** privileged or **best feature** representation, and that even the notion of similarity between patterns depends implicitly on assumption that may or may not be correct. (Watanabe)

- Various phenomena that arise when analyzing data in high-dimensional spaces that do not occur in low-dimensional settings



a) 1D - 4 regions    b) 2D - 16 regions    c) 3D - 64 regions

- Estimation Error

– Bias: difference between expected prediction and correct value

– Variance: variability of prediction for a given data point

– Trade-off

# Machine Learning Algorithms

- Read more at:
  - [https://www.javatpoint.com/machine-learning-algorithms](https://www.javatpoint.com/machine-learning-algorithms)
  - [https://machinelearningcoban.com/](https://machinelearningcoban.com/)

# Basic Statistics

# Reviews

|  | Population | Sample |
|---|---|---|
| # of subjects | $N$ | $n$ |
| Mean | $\mu = \dfrac{\sum_{i=1}^{N} x_i}{N}$ | $\bar{x} = \dfrac{\sum_{i=1}^{n} x_i}{n}$ |
| Variance | $\sigma^2 = \dfrac{\sum_{i=1}^{N} (x_i - \mu)^2}{N}$ | $S^2 = \dfrac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n-1}$ |

Note: $S^2$ is the formula for unbiased sample variance, since we're dividing by $n - 1$.

| | | |
|---|---|---|
| Standard deviation | $\sigma = \sqrt{\dfrac{\sum_{i=1}^{N} (x_i - \mu)^2}{N}}$ | $S = \sqrt{\dfrac{\sum_{i=1}^{n} (x_i - \bar{x})^2}{n-1}}$ |

Note: Finding $S$ by taking $\sqrt{S^2}$ reintroduces bias.

# Exercise

Given a sample: X = [ 10, 7, -22, 4, -9 ]

1. Calculate the mean, variance and standard deviation

2. Normalize the data points into the range [0, 1]

3. Calculate the mean, variance and standard deviation after normalizing

4. Standardize X s.t. the mean of X is 0 and standard deviation is 1

**Note**: round the value up with 2 digits after the floating point (e.g., 10.23)

$$\text{Normalization: } X_{new} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

$$\text{Standardization: } X_{new} = \frac{X - \bar{X}}{s}$$

# Exercise: Solution

```
1.
Sample: [10, 7, -22, 4, -9]
Mean = -2
Variance = 177.5
Standard deviation = 13.32
Max value = 10
Min value = -22


2.
[10, 7, -22, 4, -9] after normalization: [1.0, 0.91, 0.0, 0.81, 0.41]


3.
Sample: [1.0, 0.90625, 0.0, 0.8125, 0.40625]
Mean = 0.62
Variance = 0.17
Standard deviation = 0.42
Max value = 1.0
Min value = 0.0


4.
[10, 7, -22, 4, -9] after standardization: [0.9, 0.68, -1.5, 0.45, -0.53]
```

# Data Standardization



Source: https://scikit-learn.org/stable/auto_examples/preprocessing/plot_all_scaling.html

# Data Standardization



- Needed before statistic-based algorithm (e.g. clustering, PCA, SVM, …)
- No need in regression and tree-based algorithms

# Python: Introduction

# Introduction



Table of Contents

Link:

https://colab.research.google.com/github/jakevdp/PythonDataScienceHandbook/blob/master/notebooks/Index.ipynb

Keras

Pytorch

# Google Colab

# THANK YOU
## for YOUR ATTENTION