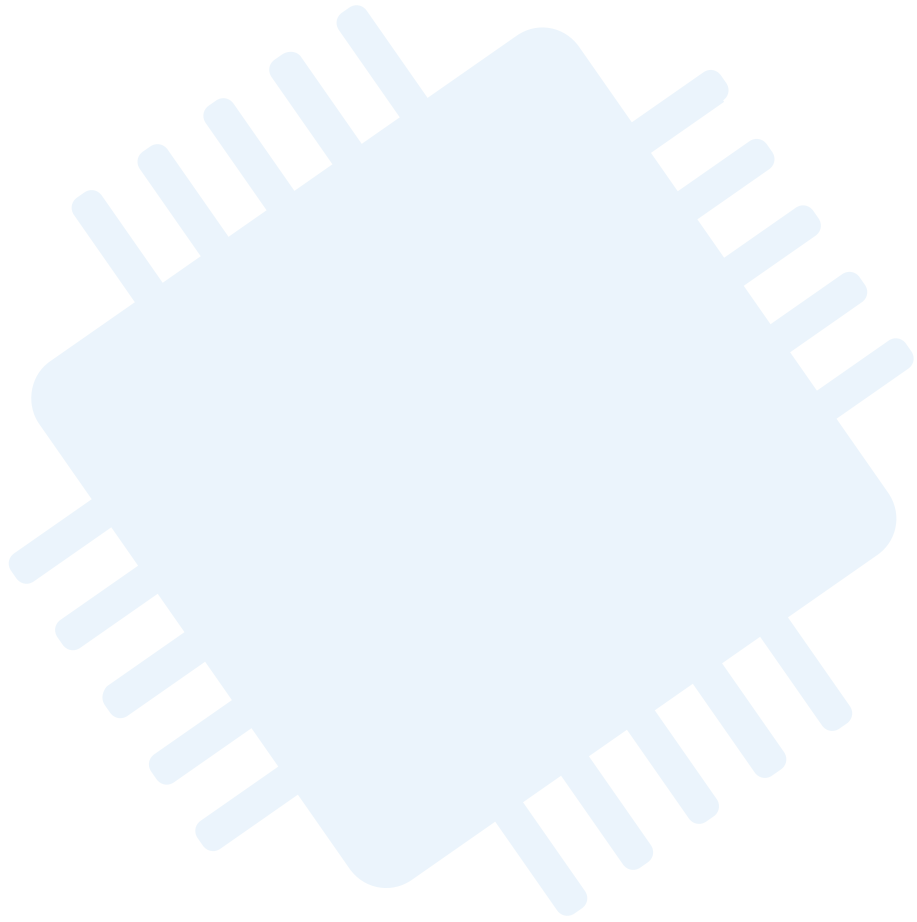# ADVANCED MACHINE LEARNING & TEXT MINING

Session 1

# Data Pre-processing

- **Definitions**

  - Text mining, also referred to as text data mining, roughly equivalent to text analytics, refers to the process of deriving high-quality information from text." –Wikipedia

  - "Another way to view text data mining is as a process of exploratory data analysis that leads to heretofore unknown information, or to answers for questions for which the answer is not currently known." -Hearst, 1999

  - "Text mining applies techniques such as categorization, entity extraction, sentiment analysis and natural language processing to transform text into data that can be used for further analysis." – expert system

- ## Why Text Mining?

  - ### 1.8 Zettabytes unstructured data

  - ### Unlike search engines, which surface documents based on keywords, text mining tools analyze documents to identify entities and extract relationships between them, unlocking hidden information to help:

    - identify and develop new hypotheses attain knowledge

    - attain knowledge

    - improve understanding.

- **Why Text Mining?**

  - Text mining ensures the use of all of available information to make better informed decisions, automate information-intensive processes, gather business critical insights and mitigate operational risks.

  - Text mining can solve high-value knowledge discovery problems in many different areas of application.

- Types of text mining

  - **Search and information retrieval**

    - Storage and retrieval of text documents, including search engines, and keyword search

  - **Document Clustering**

    - Grouping and categorizing terms snippets, paragraph or documents using clustering algorithms

  - **Document classification**

    - Grouping and categorizing terms snippets, paragraph or documents using classification methods, based on models trained on labels example

- # Types of text mining

  - ## Web mining

    - Data and text mining on the Internet with specific focus on scale and interconnectedness

  - ## Information Extraction

    - Identification and extraction relevant facts and relationships from unstructured text; the process of making structured data from unstructured and semi-structured data

- Types of text mining

  - **NLP**

    - Low level language processing and understanding tasks often used synonymously with computational linguistics

  - **Concept extraction**

    - Grouping of words and phrases into semantically similar groups

# Text mining around us

## Sentiment Analysis

- Text mining around us

Document Summarization

## ▪ Text mining around us

### Document Summarization

■ Text mining around us

## Movie Recommendation

## Text mining around us

### Text mining in healthcare

- ## Text mining around us

### Text mining in Finance

- ## Text Mining Process Flow

(Source: Claudia Peersman)

- ## Text Mining Process Flow

(Source: Miner and Eder et al)



```
Phase 1: Determine
the purpose of
study
        ↓
Phase 2: Explore
Availability and
nature of data
        →
Phase 3: Prepare
Data
        ↑
Phase 4: Develop
and asses the model
        →
Phase 5: Evaluate
the result
        ↓
Phase 6: Deploy the
model
```

- ## Text Mining Process Flow

(Source: Data Flair)

- **Text Mining is not easy**
  - Language is ambiguous

- **Homonomy**: same word, different meaning by accident of history
  - Bank
    - a. Mary walked along the <u>bank</u> of the river.
    - b. HarborBank is the richest <u>bank</u> in the city.

**Synonymy**: synonyms, different words, similar or same meaning; can substitute one word for the other without changing the meaning of the sentence substantively.

Synonyms can have differing connotations...
  - a. Miss Nelson became a kind of <u>big</u> sister to Benjamin.
  - b. Miss Nelson became a kind of <u>large</u> sister to Benjamin.

**Polysemy**: same word or form, but different, albeit related meaning

Bank
  - a. The <u>bank</u> raised its interest rates yesterday.
  - b. The store is next to the newly constructed <u>bank</u>.
  - c. The <u>bank</u> appeared first in Italy in the Renaissance.

**Hyponymy**: concept hierarchy or subclass (subordinates)

Animal (noun)
  - a. dog
  - b. cat

Injury
  - a. Broken leg, contusion...

- Concepts and word extraction usually results in a huge dimension

  - Thousands of new fields

  - Each field typically has low information

- Misspelling abbreviations, spelling variants

# TEXT MINING

| Label | Challenges |
|---|---|
| Words and morphemes | Word segmentation: dividing text into words. Fairly easy for English and other languages that use whitespace; much harder for languages like Chinese and Japanese.<br>– Assigning part of speech.<br>– Identifying synonyms; synonyms are useful for searching.<br>– Stemming: the process of shortening a word to its base or root form. For example, a simple stemming of *words* is *word*.<br>– Abbreviations, acronyms, and spelling also play important roles in understanding words. |

| | |
|---|---|
| Multiword and sentence | Phrase detection: *quick red fox*, *hockey legend Bobby Orr*, and *big brown shoe* are all examples of phrases.<br>– Parsing: breaking sentences down into subject-verb and other relationships often yields useful information about words and their relationships to each other.<br>– Sentence boundary detection is a well-understood problem in English, but is still not perfect.<br>– Coreference resolution: "Jason likes dogs, but he would never buy one." In this example, *he* is a coreference to Jason. The need for coreference resolution can also span sentences. |

| | |
|---|---|
| Multiword and sentence | – Words often have multiple meanings; using the context of a sentence or more may help choose the correct word. This process is called *word sense disambiguation* and is difficult to do well. – Combining the definitions of words and their relationships to each other to determine the meaning of a sentence |
| Multisentence and paragraph | At this level, processing becomes more difficult in an effort to find deeper understanding of an author's intent. Algorithms for summarization often require being able to identify which sentences are more important than others. |

# Introduction to Machine Learning

# LEARNING – HOW IT IS DONE?

**1**

- Mr. Micky found unknown food
- He takes the first bite

**Scenario 1: Learning**

**2**

- Mr. Micky experience poisonous features in food

**2**

- Mr. Micky experience no poisonous features in food

- Mr. Micky eats known food without second thought
- He has knowledge (through experience) of the features of the food

**Scenario 2: How to**

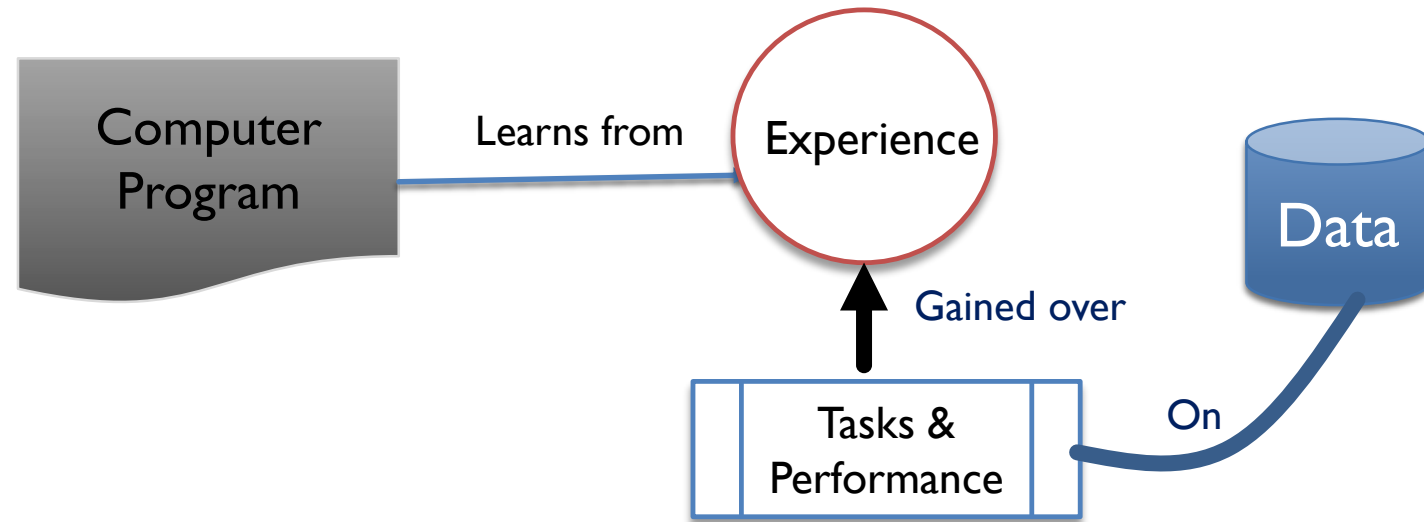- Machine learns through **computer programs** (the language of machine) which learns of from **experience**.



- Machine learns from **changes** in computer program or **data** (input); e.g., samples of speech from different person improve speech recognition.

- Definitions of Machine Learning

  - Machine learning is a method that gives computer the ability to learn **without** being explicitly programmed (Arthur, S.)

  - Machine learning is a set of **tools** that, broadly speaking, allow to "**teach**" computers how to perform tasks by providing examples of how they should be done (Mitchel, T.)

  - Machine learning is a method of data analysis that **automates** analytical model building (SAS)

▪ Learning methods of machines

▪ Inductive learning: It involves the process of learning by example – where a system tries to induce a model (general rule/function) from training data.



Training Data

Model

**General Rule**
Apple = {size = big, shape rounded, color = red}

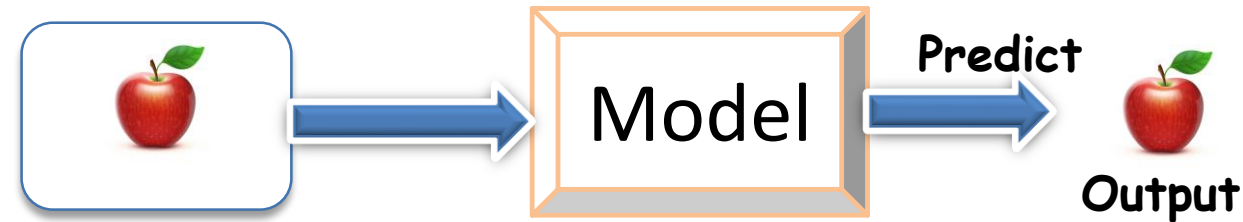- **Deductive learning**: It applies the learning model to predict the outcome.



- **Instance-based learning**
  - The most trivial form of learning
  - The system learns the examples by heart, then generalizes to new cases by using a similarity measure to compare them to the learned examples (or a subset of them)



SPAM Email

New Email

**Measure the similarity**

- Find common words between two mails
- If number of common words is significant then flag the new email as Spam

- **Model-based learning**
  - Generalize from a set of examples is to build a model of these examples and then use that model to make predictions

- There are different approaches for machine learning systems
  - **Semi Automatic**
  - **Fully Automatic**



**Semi-automatic system approach**



**Automatic system approach**

- Machine learning helps human to learn

▪ **Supervised Learning**: In supervised learning, the training set you feed to the algorithm includes the desired solutions, called **labels**
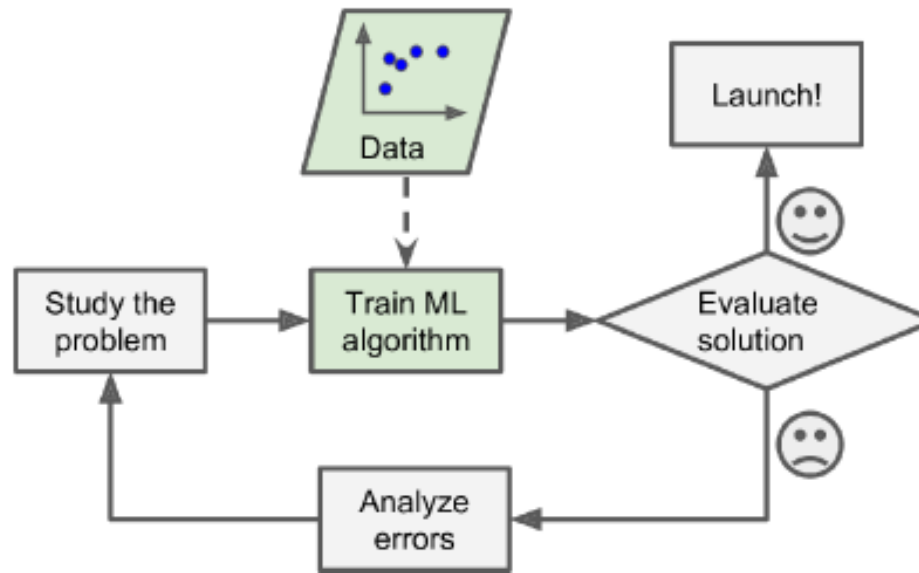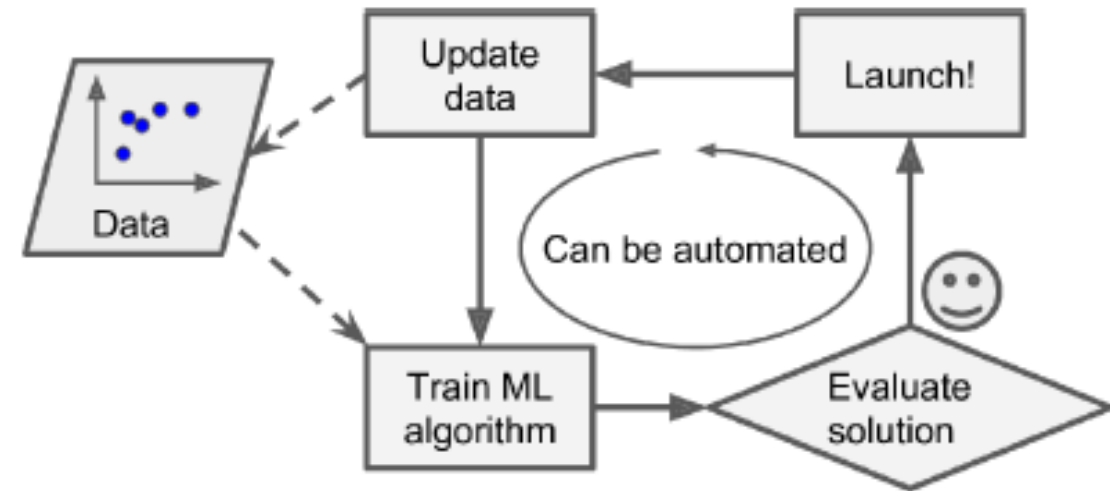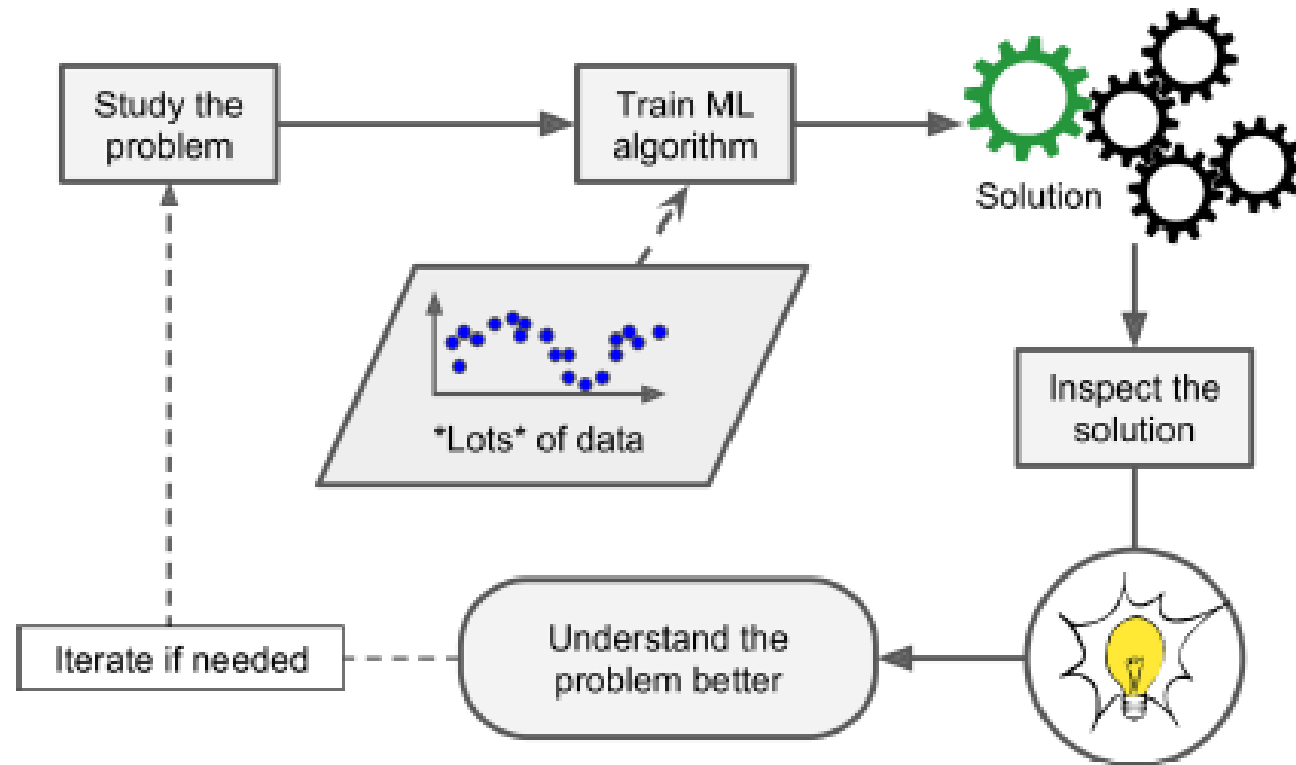
**Data:** $D = \{d_1, d_2, ..., d_n\}$   **a set of $n$ examples**

$d_i = <\mathbf{x}_i, y_i>$

$\mathbf{x}_i$ is input vector, and $y$ is desired output

**Objective:** learn the mapping $f : X \rightarrow Y$

s.t.  $y_i \approx f(x_i)$   for all $i = 1, ..., n$

**Two types of problems:**

- **Regression:** X discrete or continuous $\rightarrow$

   Y is **continuous**

- **Classification:** X discrete or continuous $\rightarrow$

   Y is **discrete**

Ex: Price of a car, given a set of *features*

Ex: Classifying emails along with their class spam or Not Spam

- **Note:** Some regression algorithms can be used for classification as well, and vice versa.
  - For example, Logistic Regression is commonly used for classification, as it can output a value that corresponds to the probability of belonging to a given class

- The most widely used Supervised Learning Algorithms

  - k-Nearest Neighbors
  - Linear Regression
  - Logistic Regression
  - Support Vector Machines (SVMs)
  - Decision Trees and Random Forests
  - Neural networks

- **Unsupervised Learning**

  - It is a type of learning where algorithm used to draw inferences from datasets consisting of input data (x) without labeled responses.
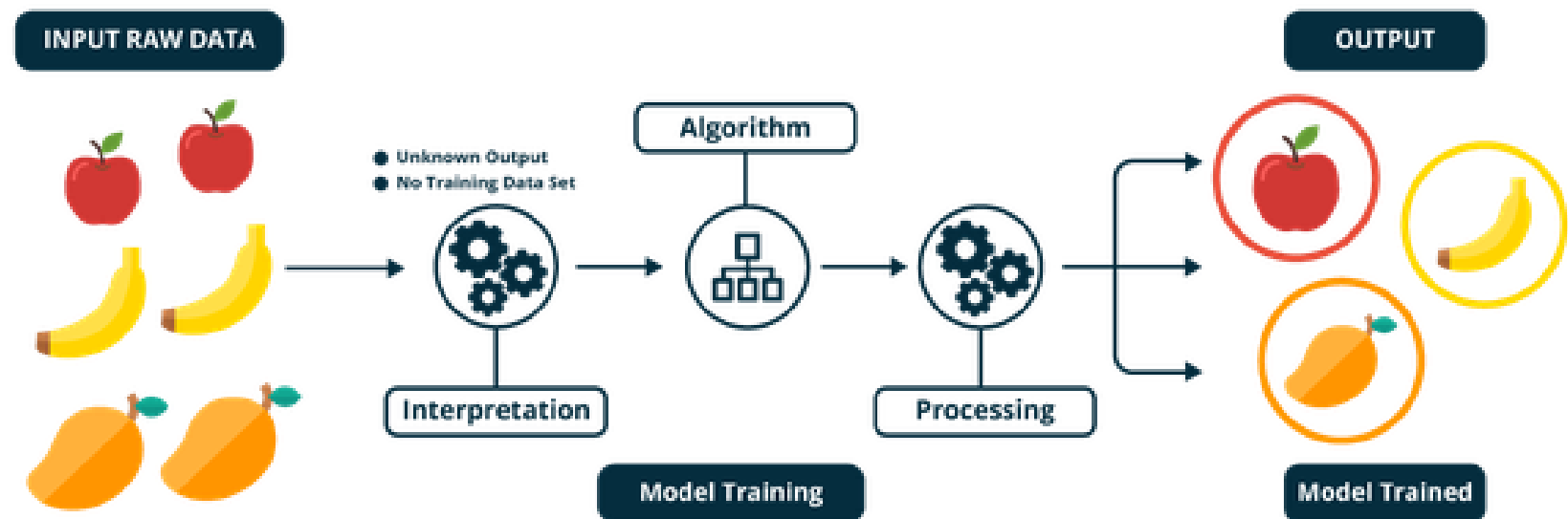
  - Goal:

    - In some pattern recognition problems, the training data consists of a set of input vectors x without any corresponding target values.

    - The goal in such unsupervised learning problems may be to discover groups of similar examples within the data.

- **Unsupervised Learning**
  - Training data provides "example", but we have no specific outcomes.
  - In simple word **there is no label associated with this learning.**
  - In unsupervised learning the machine tries to find interesting patterns in the data.
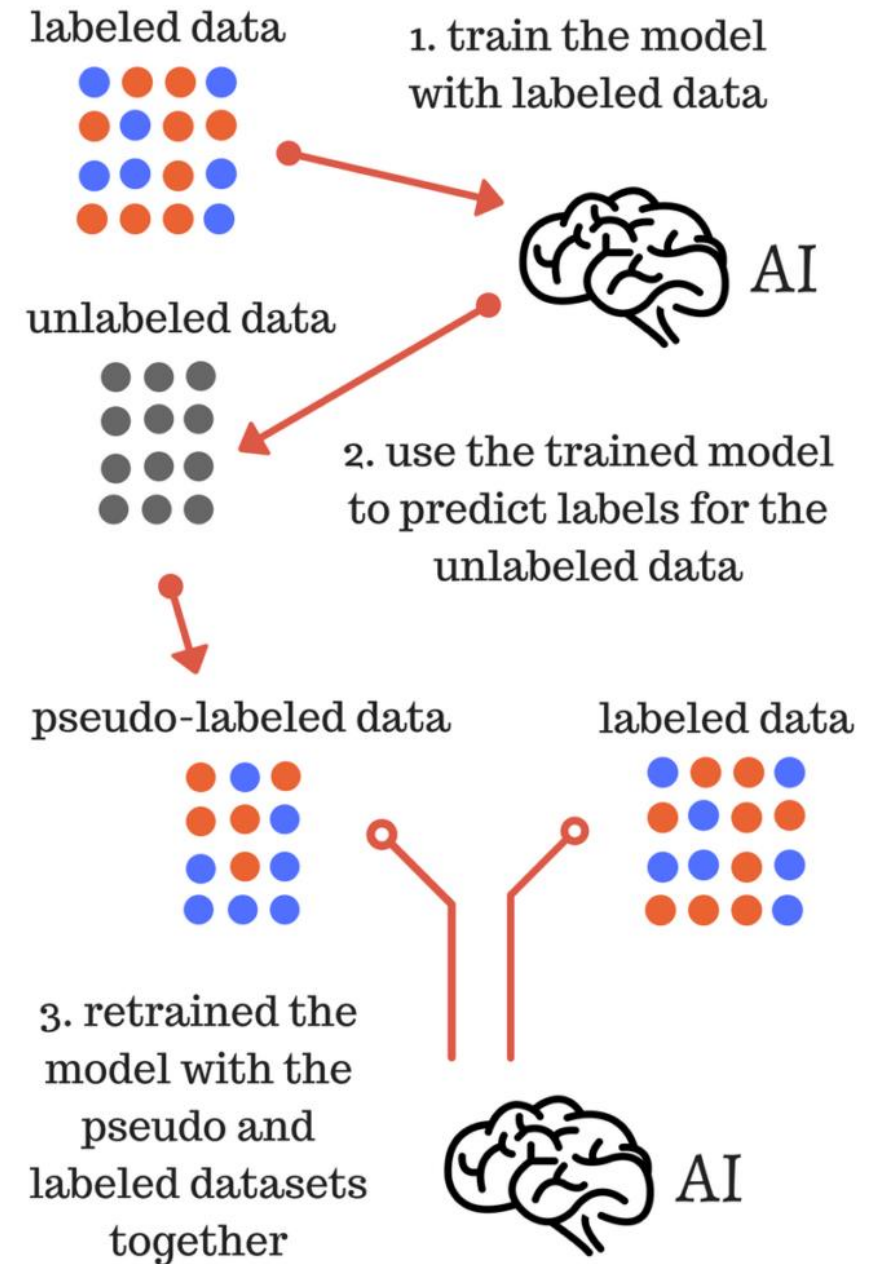
▪ **Unsupervised Algorithms**

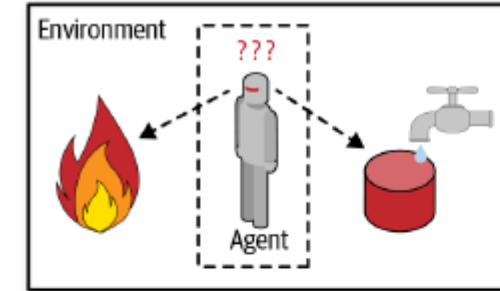| Unsupervised Algorithms | |
|---|---|
| Clustering | —K-Means<br>—DBSCAN<br>—Hierarchical Cluster Analysis (HCA) |
| Anomaly detection and novelty detection | —One-class SVM<br>—Isolation Forest |
| Visualization and dimensionality reduction | —Principal Component Analysis (PCA)<br>—Kernel PCA<br>—Locally Linear Embedding (LLE)<br>—t-Distributed Stochastic Neighbor Embedding (t-SNE) |
| Association rule learning | —Apriori<br>—Eclat |

- **Semi-supervised Learning**
  - Semi-supervised learning (SSL) is halfway between supervised and unsupervised learning
  - In addition to unlabelled data, the algorithm is provided with some super-vision information – but not necessarily for all examples

labeled data

1. train the model with labeled data

unlabeled data

2. use the trained model to predict labels for the unlabeled data

pseudo-labeled data

labeled data

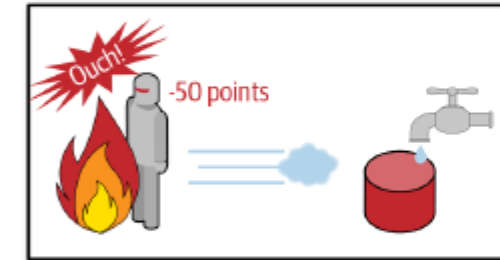3. retrained the model with the pseudo and labeled datasets together

## ▪ Reinforcement Learning

  - ▪ The learning system, called an agent in this context, can observe the environment, select and perform actions, and get rewards in return (or penalties in the form of negative rewards.

  - ▪ It must then learn by itself what is the best strategy, called a policy, to get the most reward over time

  - ▪ A policy defines what action the agent should choose when it is in a given situation.



Environment
??? 
Agent

**1** Observe

**2** Select action using policy

Ouch! -50 points

**3** Action!
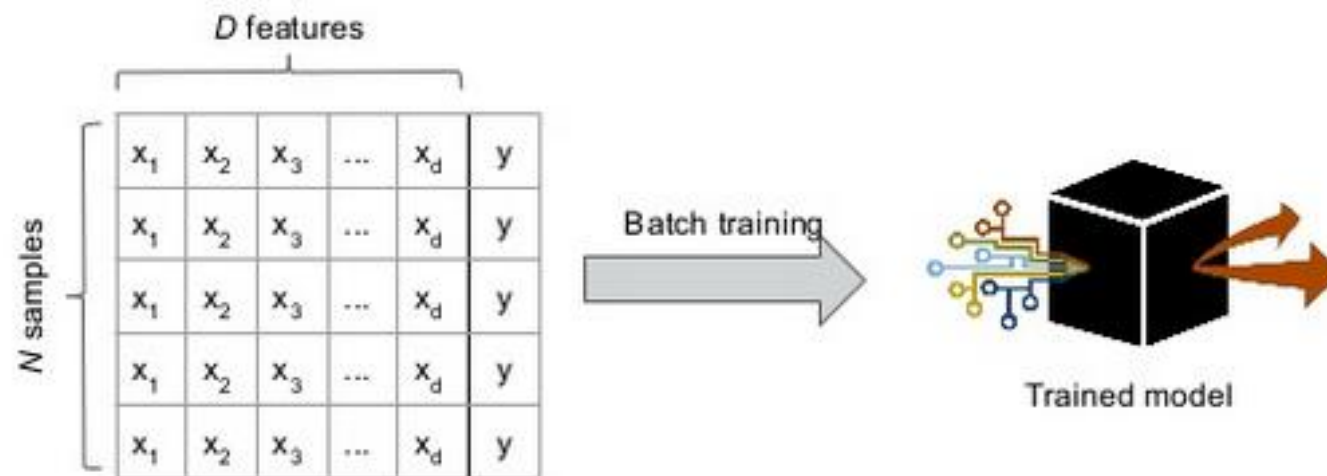
**4** Get reward or penalty

🔥 = bad!
Next time avoid it.

**5** Update policy (learning step)

**6** Iterate until an optimal policy is found
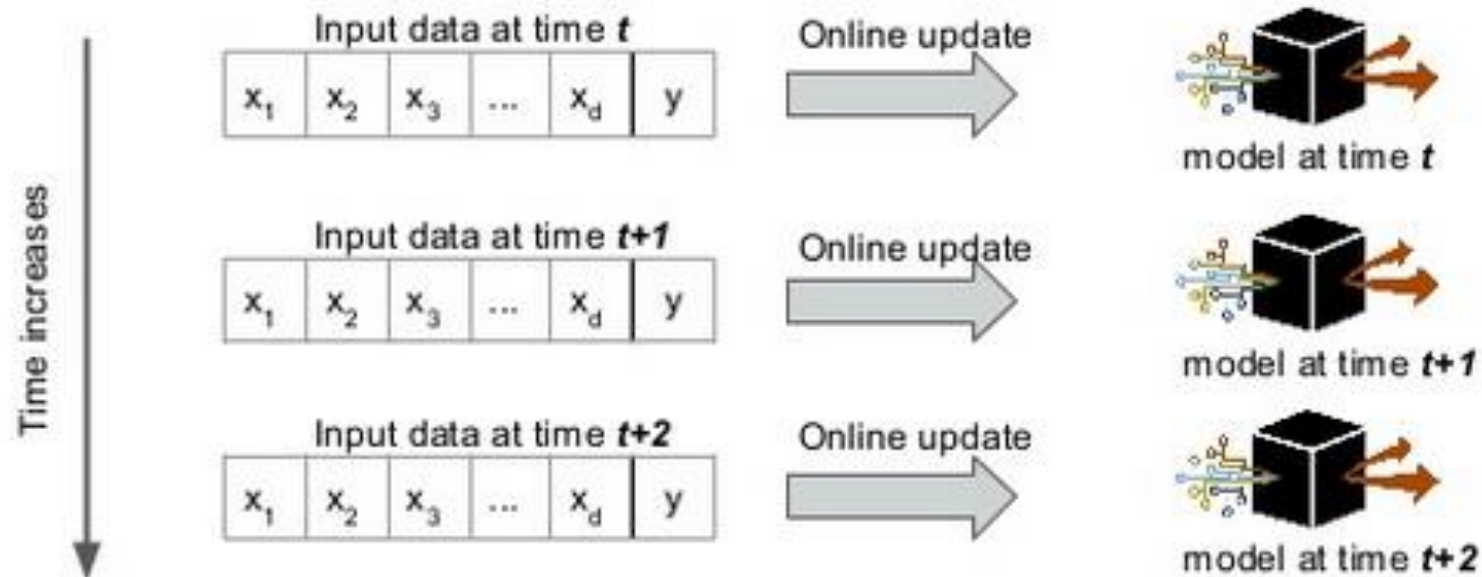
39

- **Batch Learning**
  - The system must be trained using all available data
  - The system is incapable of learning incrementally
  - For new data (such as a new type of spam),
    - Train a new version of the system from scratch on the full dataset (not just the new data, but also the old data)
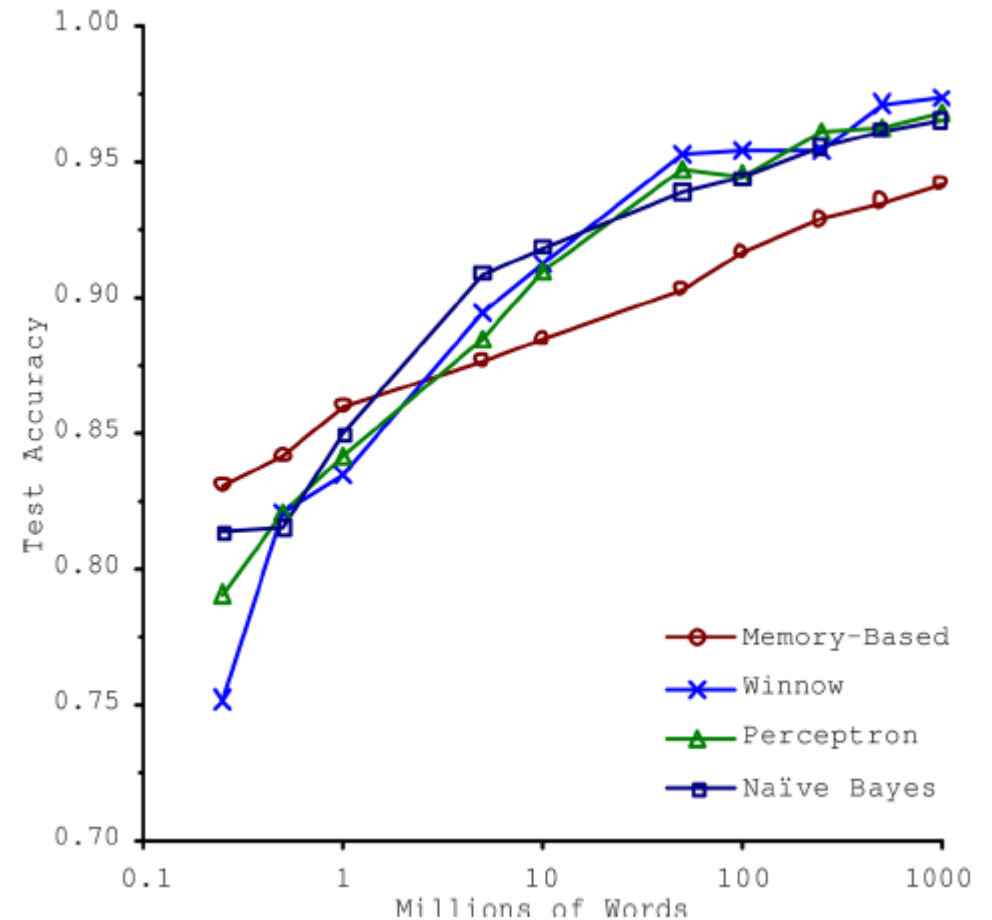    - Then stop the old system and replace it with the new one.

## Online Learning

- Train the system incrementally by feeding it data instances sequentially, either individually or in small groups called mini-batches.

- Each learning step is fast and cheap, so the system can learn about new data on the fly, as it arrives
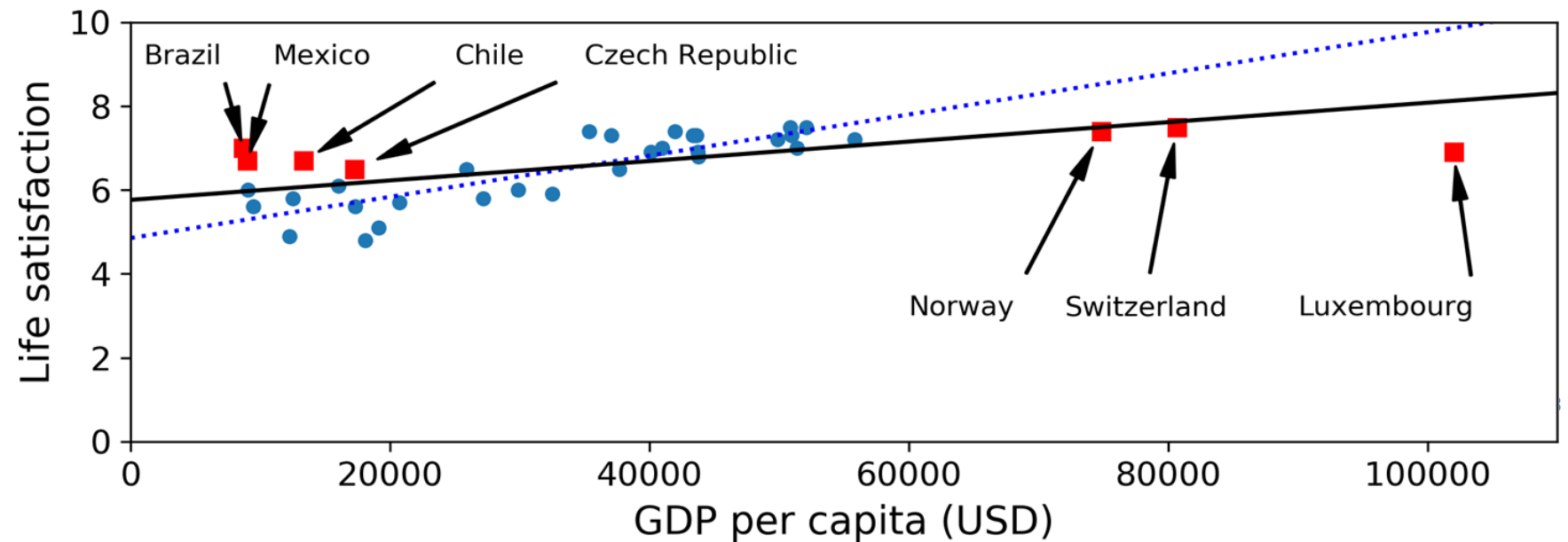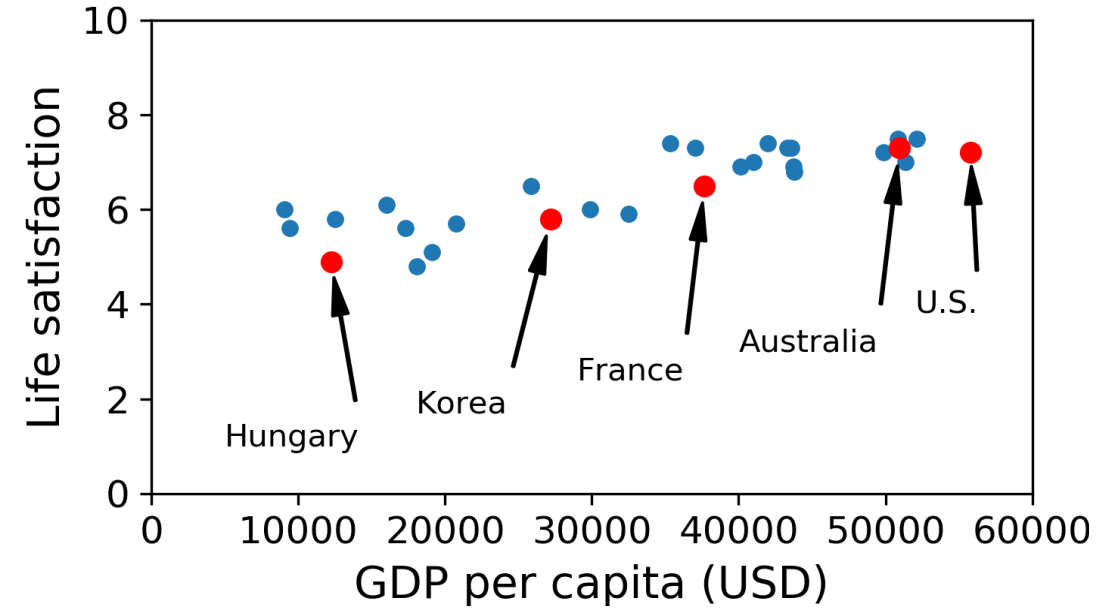
■ **Insufficient Quantity of Data**

   ■ For very simple problems you typically need thousands of examples

   ■ For complex problems such as image or speech recognition you may need millions of examples

## Non-Representative Training Data

- It is crucial that your training data be representative of the new cases you want to generalize to

## Poor Quality Data

- If training data is full of errors, outliers, and noise it will make it harder for the system to detect the underlying patterns, so your system is less likely to perform well

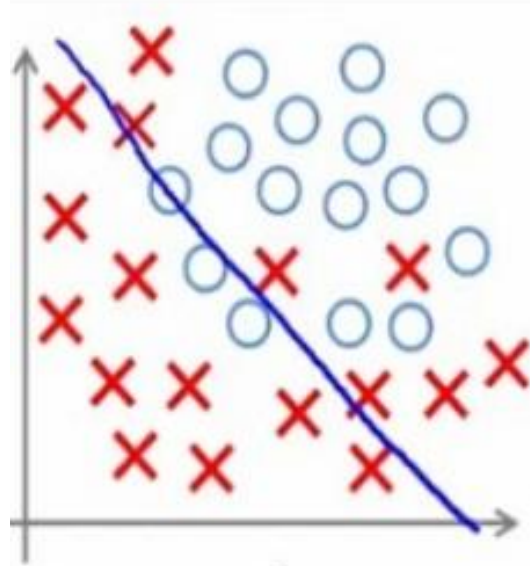## Overfitting training data

- The model performs well on the training data, but it does not generalize well.

## Underfitting

- It occurs when your model is too simple to learn the underlying structure of the data
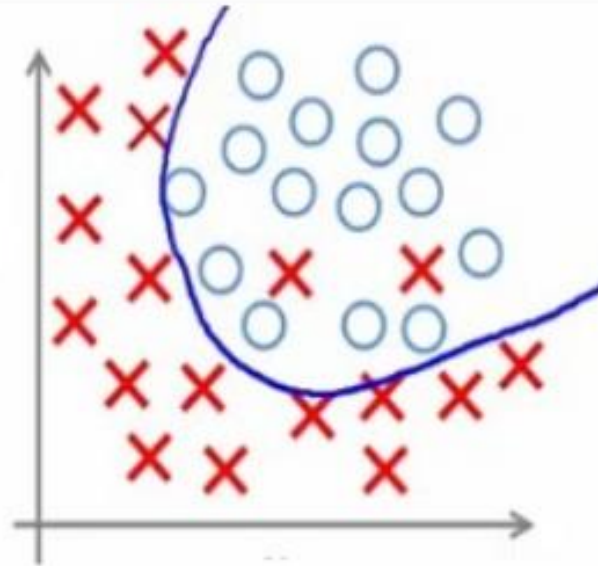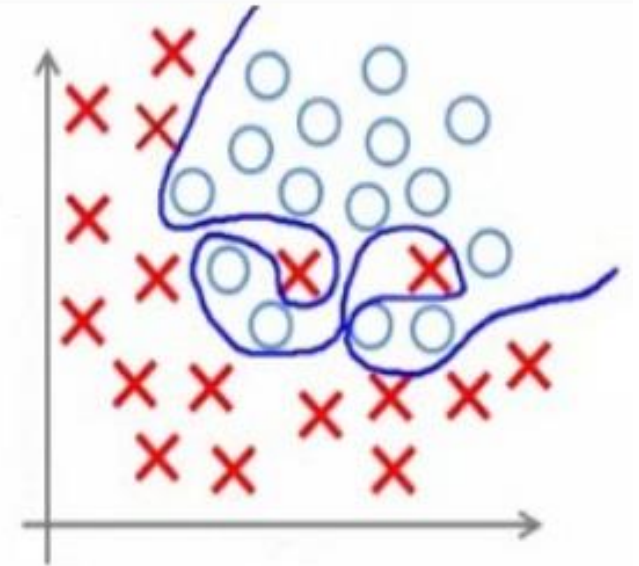
**Under-fitting**

(too simple to explain the variance)

**Appropriate-fitting**

**Over-fitting**

(forcefitting -- too good to be true)

- **Data Pre-processing**

Feature extraction

- Tokenization

- Removing unnecessary punctuation, tags

- Removing stop words

- Stemming

- Lemmatization/POS tagging

Feature selection

- Term frequency

- Chi-square

- Expected cross entropy

- Odds Ratio

- Algorithms
  - Support Vector Machines
  - Naive Bayes
  - Maximum Entropy
  - K-nearest neighbour
  - Decision Tree
  - TF-IDF
  - Conditional Random Field (CRF)
  - Latent Dirichlet Allocation (LDA)
  - Artificial Neural Networks
  - LSTM
  - Bi-directional LSTM
  - Recurrent neural networks
  - Models of word embedding:
    - ·Word2Vec
    - ·Glove

- Common NLP techniques

    - Named Entity Recognition

    - Sentiment Analysis

    - Language detection

    - Topic Modeling

    - Text Summarization

    - Machine Translation