# Machine Learning Introduction

# *What is Machine Learning?*

Machine learning provides systems the ability to automatically learn and improve from experience without being explicitly programmed. Machine learning focuses on the development of computer programs that can access data and use it learn for themselves.

# *Why "Learn"?*

- Machine learning is programming computers to optimize a performance criterion using example data or past experience.

# *What We Talk About When We Talk About "Learning"*

- Learning general models from a data of particular examples

- Data is cheap and abundant (data warehouses, data marts); knowledge is expensive and scarce.

- Example in retail: Customer transactions to consumer behavior:

    *People who bought "Da Vinci Code" also bought "The Five People You Meet in Heaven" (www.amazon.com)*

- Build a model that is *a good and useful approximation* to the data.

# *Data Mining/KDD*

Definition := *"KDD(Knowledge Discovery in Database) is the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data"* (Fayyad)

Applications:

- Retail: Market basket analysis, Customer relationship management (CRM)
- Finance: Credit scoring, fraud detection
- Manufacturing: Optimization, troubleshooting
- Medicine: Medical diagnosis
- Telecommunications: Quality of service optimization
- Bioinformatics: Motifs, alignment
- Web mining: Search engines
- ...

# *What is Machine Learning?*

- Machine Learning
  - Study of algorithms that
  - improve their performance
  - at some task
  - with experience
- Optimize a performance criterion using example data or past experience.
- Role of Statistics: Inference from a sample
- Role of Computer science: Efficient algorithms to
  - Solve the optimization problem
  - Representing and evaluating the model for inference

# *Growth of Machine Learning*

- Machine learning is preferred approach to
  - Speech recognition, Natural language processing
  - Computer vision
  - Medical outcomes analysis
  - Robot control
  - Computational biology
- This trend is accelerating
  - Improved machine learning algorithms
  - Improved data capture, networking, faster computers
  - Software too complex to write by hand
  - New sensors / IO devices
  - Demand for self-customization to user, environment

# *Applications*

- Supervised Learning
  - Classification
  - Regression/Prediction
- Unsupervised Learning
- Reinforcement Learning

# *Supervised Learning*

The majority of practical machine learning uses supervised learning.

Supervised learning is where you have input variables (x) and an output variable (Y) and you use an algorithm to learn the mapping function from the input to the output.

$$Y = f(X)$$

The goal is to approximate the mapping function so well that when you have new input data (x) that you can predict the output variables (Y) for that data.

# *Supervised Learning*

It is called supervised learning because the process of an algorithm learning from the training dataset can be thought of as a teacher supervising the learning process.

We know the correct answers, the algorithm iteratively makes predictions on the training data and is corrected by the teacher. Learning stops when the algorithm achieves an acceptable level of performance.

# *Supervised Learning*

Supervised learning problems can be further grouped into regression and classification problems.

- **Classification**: A classification problem is when the output variable is a category, such as "red" or "blue" or "disease" and "no disease".
- **Regression**: A regression problem is when the output variable is a real value, such as "dollars" or "weight".

Some common types of problems built on top of classification and regression include recommendation and time series prediction respectively.

# *Unsupervised Machine Learning*

Unsupervised learning is where you only have input data (X) and no corresponding output variables.

These are called unsupervised learning because unlike supervised learning above there is no correct answers and there is no teacher. Algorithms are left to their own devices to discover and present the interesting structure in the data.

Some popular examples of unsupervised learning algorithms are:

- K-means for clustering problems.

# *Unsupervised Learning*

Unsupervised learning problems can be further grouped into clustering and association problems.

- **Clustering**: A clustering problem is where you want to discover the inherent groupings in the data, such as grouping customers by purchasing behavior.
- **Association**:  An association rule learning problem is where you want to discover rules that describe large portions of your data, such as people that buy X also tend to buy Y.

# *Learning Associations*

- Basket analysis:

  $P(Y|X)$ probability that somebody who buys $X$ also buys $Y$ where $X$ and $Y$ are products/services.

  Example: $P(\text{chips} \mid \text{beer}) = 0.7$

Market-Basket transactions

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

# *Classification*

Classification is a process of finding a function which helps in dividing the dataset into classes based on different parameters.

In Classification, a computer program is trained on the training dataset and based on that training, it categorizes the data into different classes.

The task of the classification algorithm is to find the mapping function to map the input(x) to the discrete output(y).

# *Classification*

**Example:** The best example to understand the Classification problem is Email Spam Detection. The model is trained on the basis of millions of emails on different parameters, and whenever it receives a new email, it identifies whether the email is spam or not. If the email is spam, then it is moved to the Spam folder.
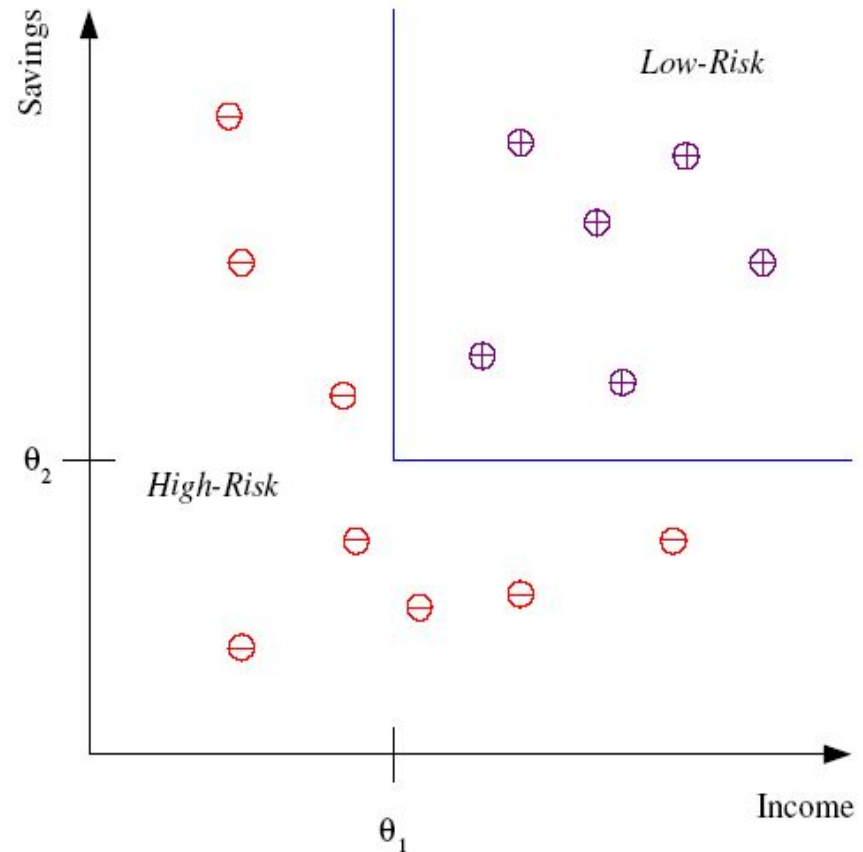
# *ML Classification Algorithms*

- Logistic Regression

- K-Nearest Neighbours

- Support Vector Machines

- Naïve Bayes

- Decision Tree Classification

- Random Forest Classification

# *Classification*

- Example: Credit scoring
- Differentiating between low-risk and high-risk customers from their *income* and *savings*



Discriminant: IF *income* > $\theta_1$ AND *savings* > $\theta_2$
THEN low-risk ELSE high-risk

Model

# *Classification: Applications*

- Aka Pattern recognition
- Face recognition: Pose, lighting, occlusion (glasses, beard), make-up, hair style
- Character recognition: Different handwriting styles.
- Speech recognition: Temporal dependency.
  - Use of a dictionary or the syntax of the language.
  - Sensor fusion: Combine multiple modalities; eg, visual (lip image) and acoustic for speech
- Medical diagnosis: From symptoms to illnesses
- Web Advertizing: Predict if a user clicks on an ad on the Internet.

# *Face Recognition*

Training examples of a person



Test images



AT&T Laboratories, Cambridge UK
http://www.uk.research.att.com/facedatabase.html

# *Regression*

Regression is a process of finding the correlations between dependent and independent variables. It helps in predicting the continuous variables such as prediction of Market Trends, prediction of House prices, etc.

The task of the Regression algorithm is to find the mapping function to map the input variable(x) to the continuous output variable(y).

**Example:** Suppose we want to do weather forecasting, so for this, we will use the Regression algorithm. In weather prediction, the model is trained on the past data, and once the training is completed, it can easily predict the weather for future days.

# *Regression Algorithms*

- Simple Linear Regression

- Polynomial Regression

- Support Vector Regression

- Decision Tree Regression

- Random Forest Regression

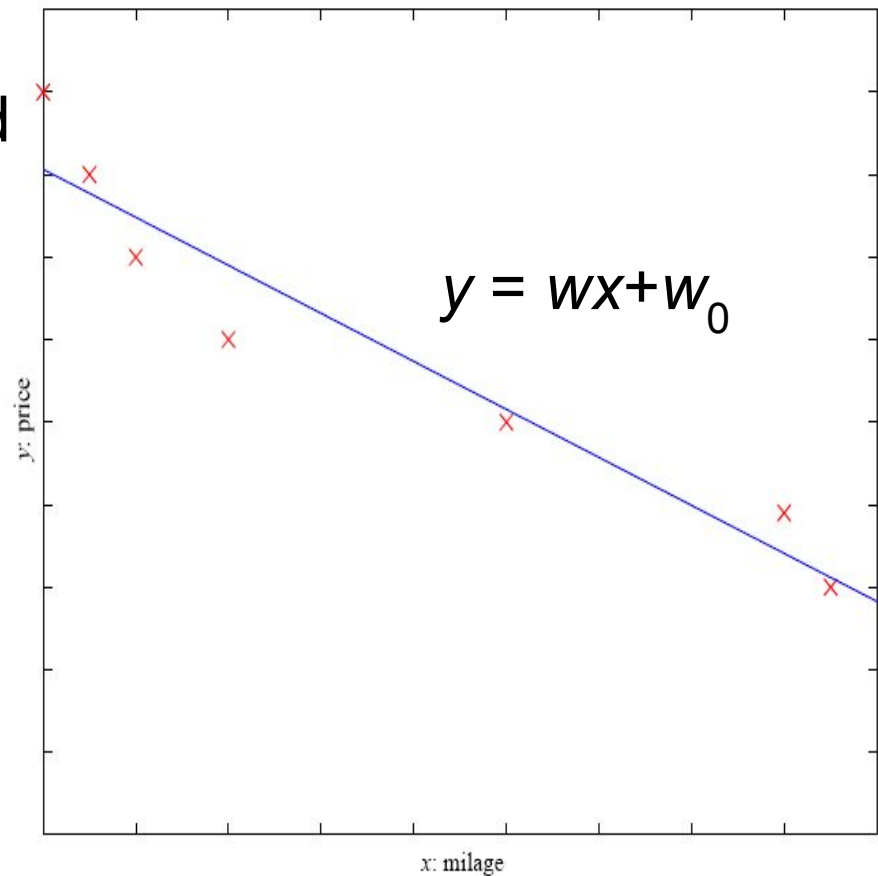| Regression Algorithm | Classification Algorithm |
|---|---|
| In Regression, the output variable must be of continuous nature or real value. | In Classification, the output variable must be a discrete value. |
| The task of the regression algorithm is to map the input value (x) with the continuous output variable(y). | The task of the classification algorithm is to map the input value(x) with the discrete output variable(y). |
| Regression Algorithms are used with continuous data. | Classification Algorithms are used with discrete data. |
| In Regression, we try to find the best fit line, which can predict the output more accurately. | In Classification, we try to find the decision boundary, which can divide the dataset into different classes. |
| Regression algorithms can be used to solve the regression problems such as Weather Prediction, House price prediction, etc. | Classification Algorithms can be used to solve classification problems such as Identification of spam emails, Speech Recognition, Identification of cancer cells, etc. |
| The regression Algorithm can be further divided into Linear and Non-linear Regression. | The Classification algorithms can be divided into Binary Classifier and Multi-class Classifier. |

# *Prediction: Regression*

- Example: Price of a used car

- $x$ : car attributes

  $y$ : price

  $y = g\,(x\mid\theta\,)$

  $g\,(\ )$ model,

  $\theta$ parameters

$y = wx + w_0$

# *Regression Applications*

- Navigating a car: Angle of the steering wheel (CMU NavLab)
- Kinematics of a robot arm

$(x,y)$
$\alpha_1 = g_1(x,y)$
$\alpha_2 = g_2(x,y)$

$\alpha_2$

$\alpha_1$

# *Issues with Unsupervised Learning*

- Unsupervised Learning is harder as compared to Supervised

  Learning tasks..

- How do we know if results are meaningful since no answer

  labels are available?

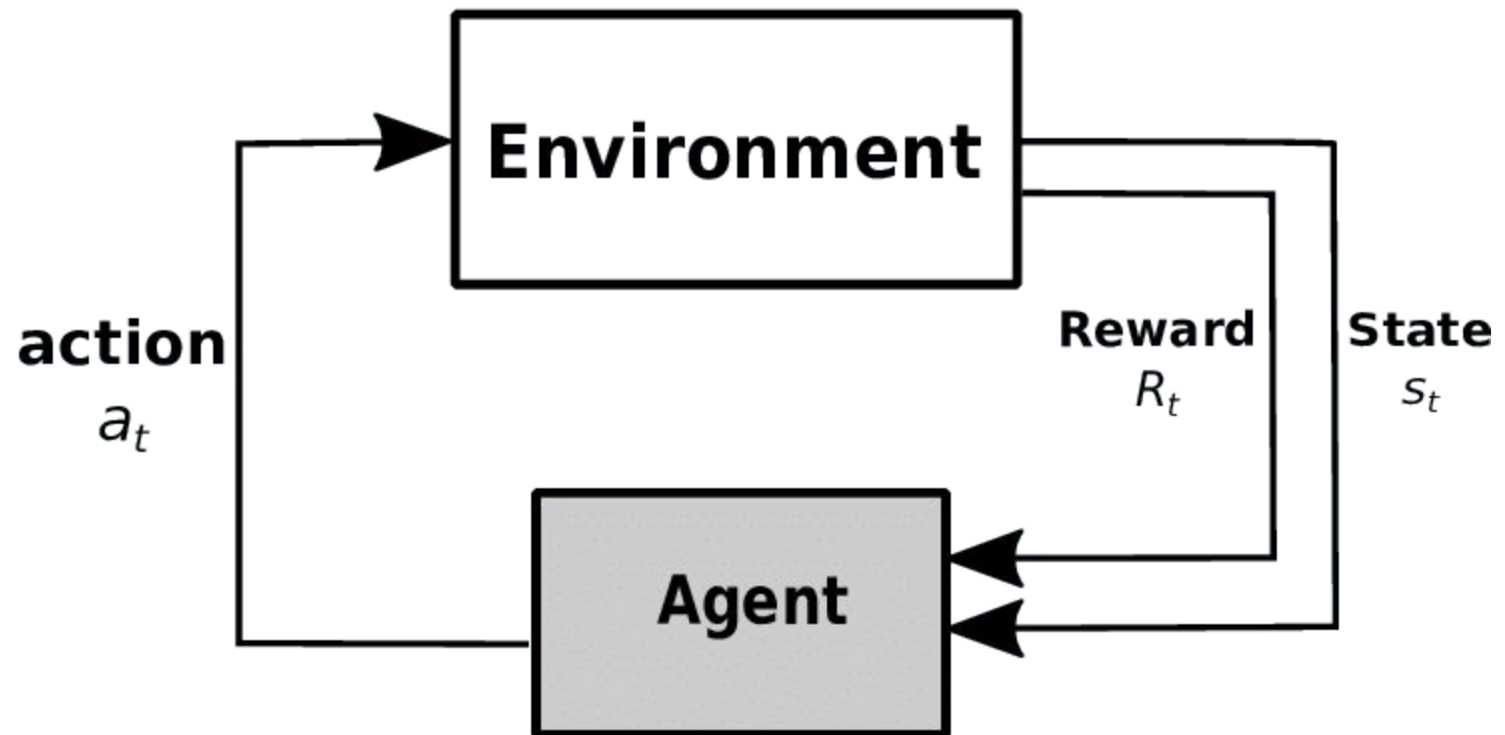- Let the expert look at the results (external evaluation)

# *Why is it still used?*

- Annotating large datasets is very costly and hence we can label only a few examples manually. Example: Speech Recognition

- There may be cases where we don't know how many/what classes is the data divided into. Example: Data Mining

- We may want to use clustering to gain some insight into the structure of the data before designing a classifier.

# *Reinforcement Learning*

- Topics:
  - Policies: what actions should an agent take in a particular situation
  - Utility estimation: how good is a state ($\square$used by policy)
- Credit assignment problem (what was responsible for the outcome)
- Applications:
  - Game playing
  - Robot in a maze

# *Types of Fits usually seen*



**Under-fitting**
(too simple to explain the variance)

**Appropirate-fitting**

**Over-fitting**
(forcefitting--too good to be true)