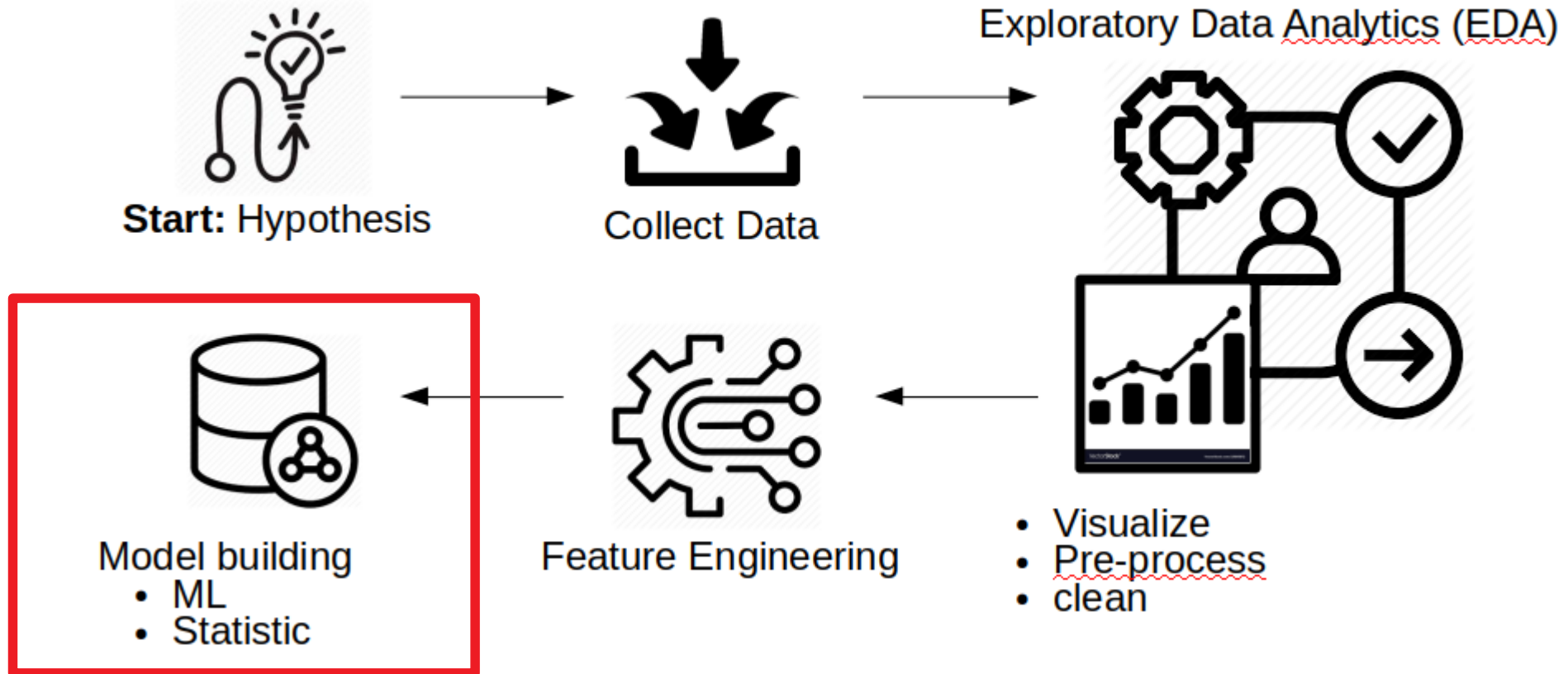


## Machine Learning I

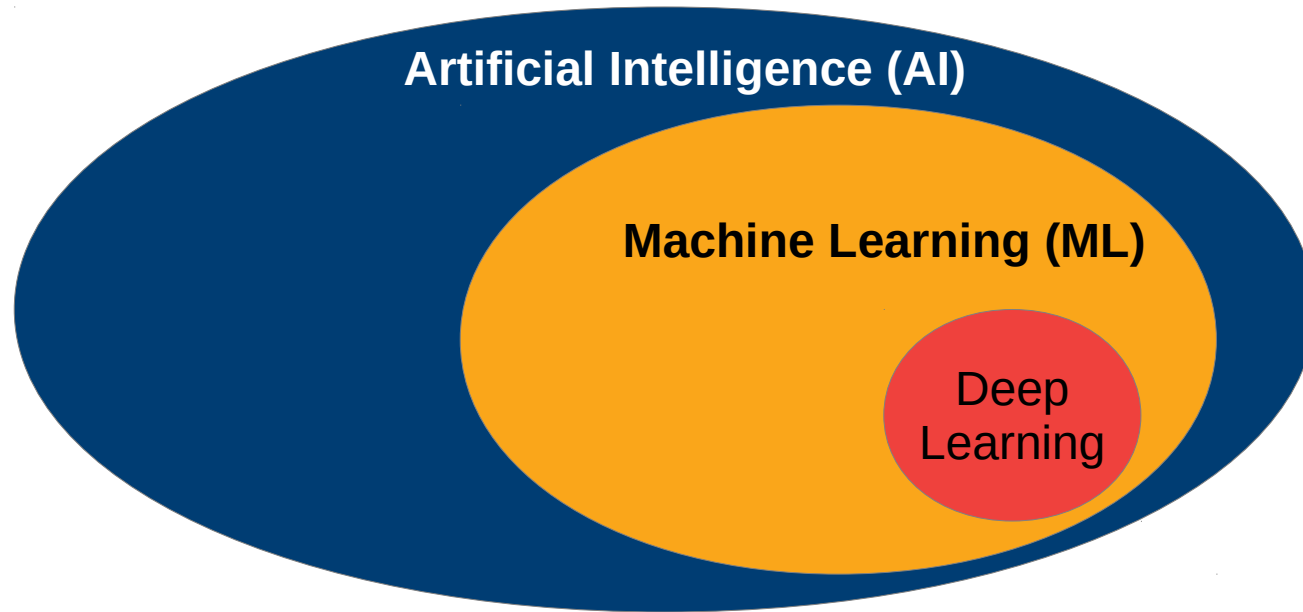
### Introduction and Overview



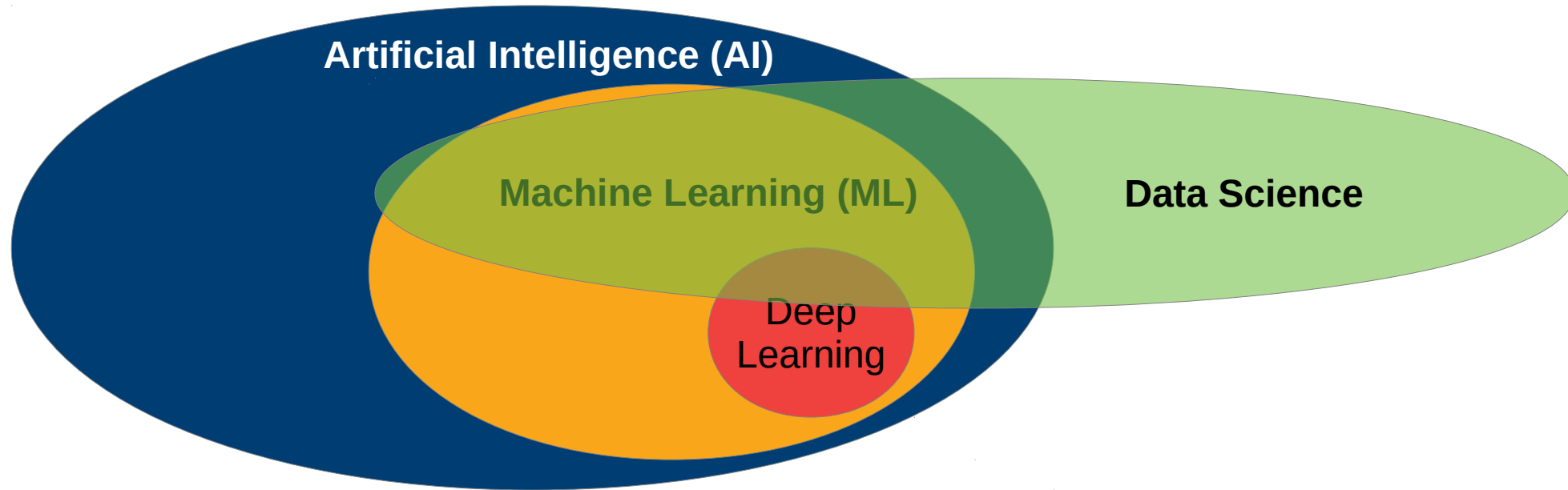
## Outline

- Introduction to ML
  - Basic Definitions and Terminology
    - Supervised Learning
    - Generalization and Overfitting
    - Unsupervised Learning

## Research and Application Fields

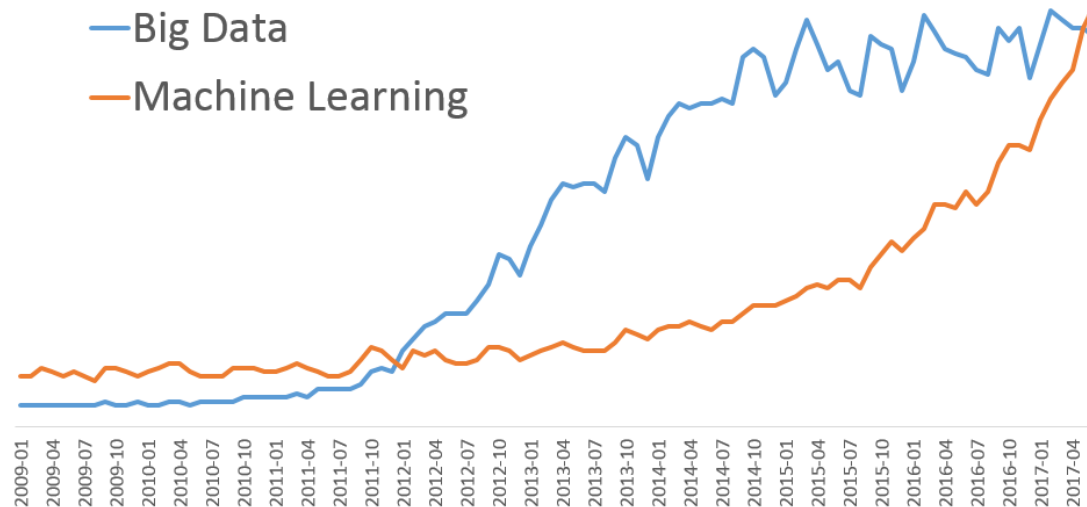


## Research and Application Fields



## The ML Hype

### Google Trends Worldwide



## Basic Types of Machine Learning Algorithms

**Supervised Learning**

**Unsupervised Learning**

**Reinforcement Learning**

## Basic Types of Machine Learning Algorithms

**Supervised Learning**

**Unsupervised Learning**

**Reinforcement Learning**

- Labeled data
- Direct and quantitative evaluation
- Learn model from „ground truth“ examples
- Predict unseen examples



## Supervised Learning

Basic Notation:

Data is given as tuples

$$(X, Y) := \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

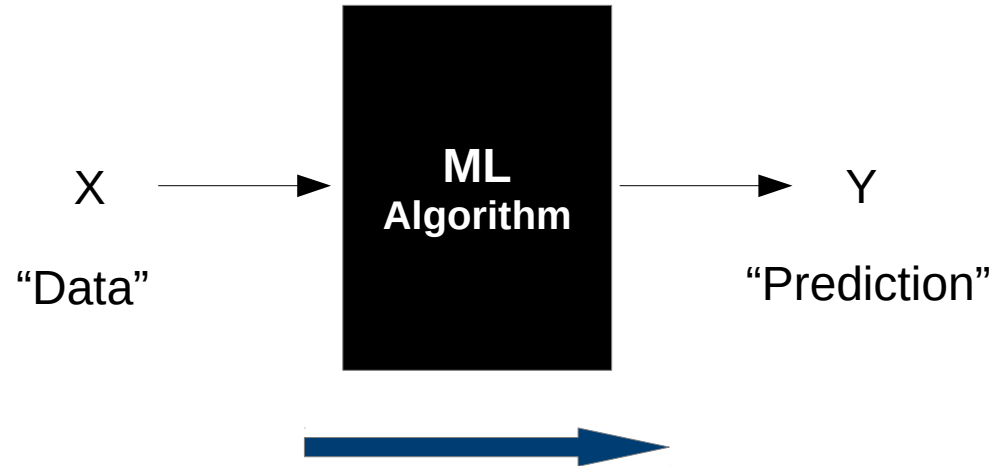
Where  $X$  is the actual **data** (sample) and  $y$  the associated **label**.

For most ML algorithms (**many Deep Learning algorithms are an exception**)

$$x_i \in \mathbb{R}^n, y_i \in \mathbb{R}$$

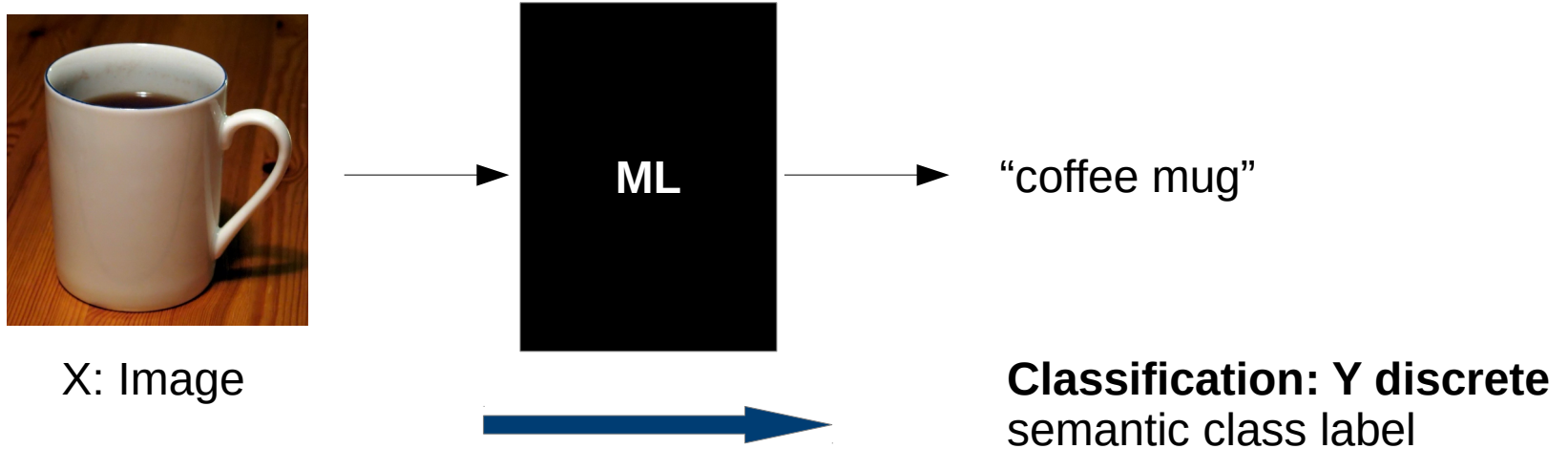
**The data has to be represented as vectors and the labels are scalars.**

## Supervised Learning as a Black Box



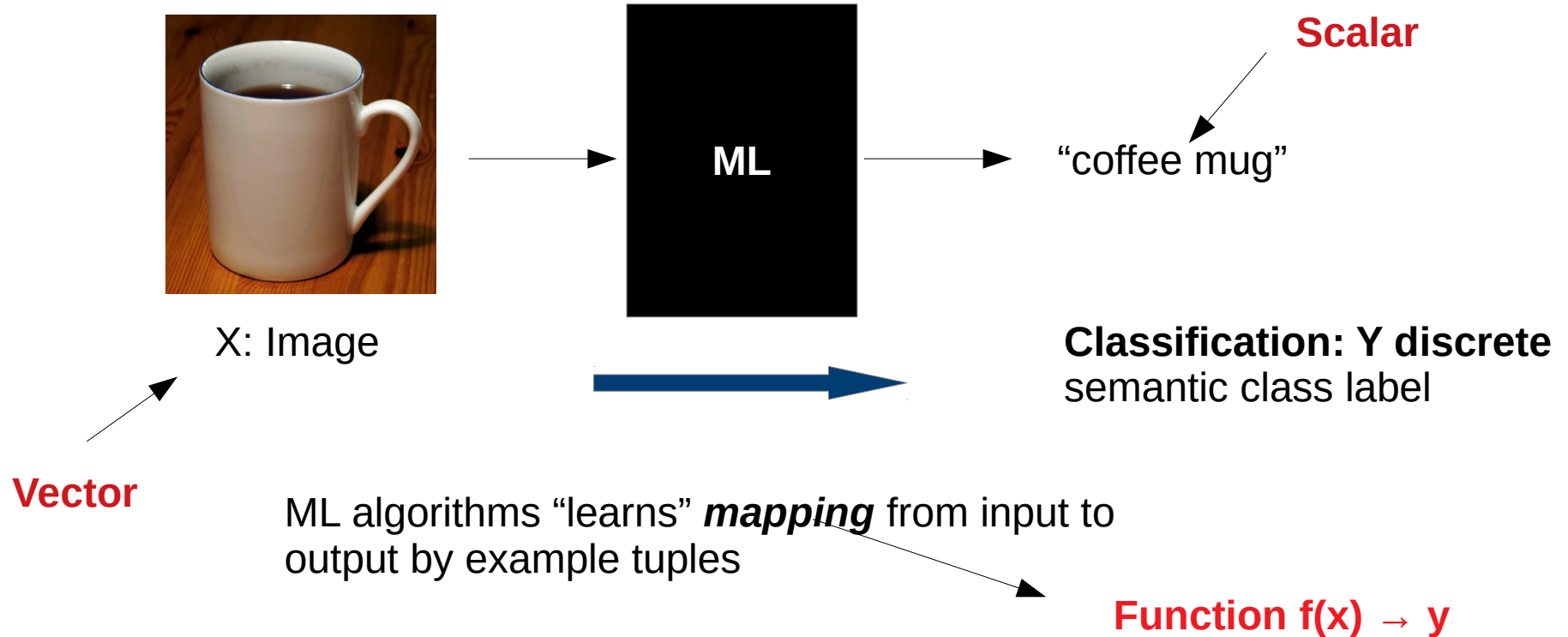
ML algorithms “learns” *mapping* from input to output by example tuples

## Supervised Learning: Example: Classification



ML algorithms “learns” *mapping* from input to output by example tuples

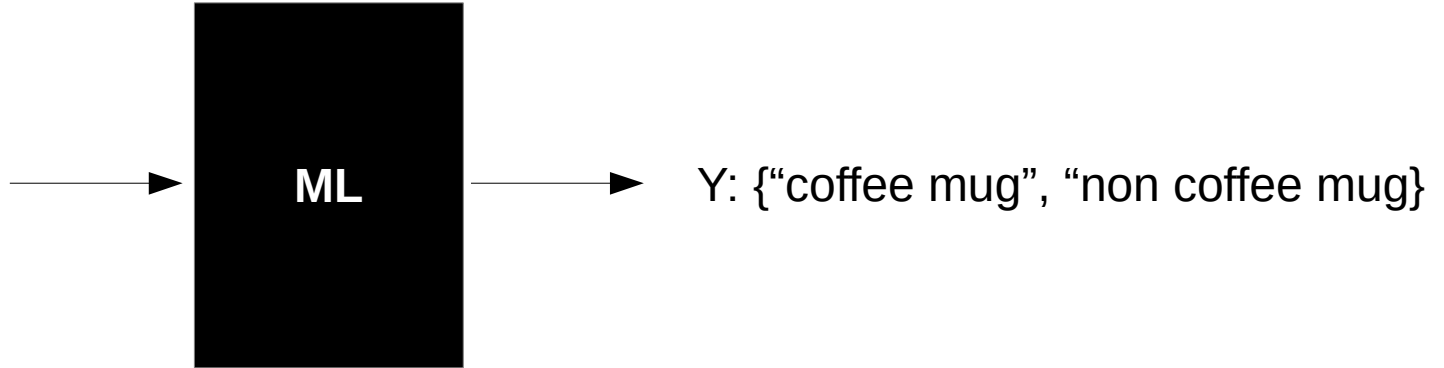
## Supervised Learning: Example: Classification



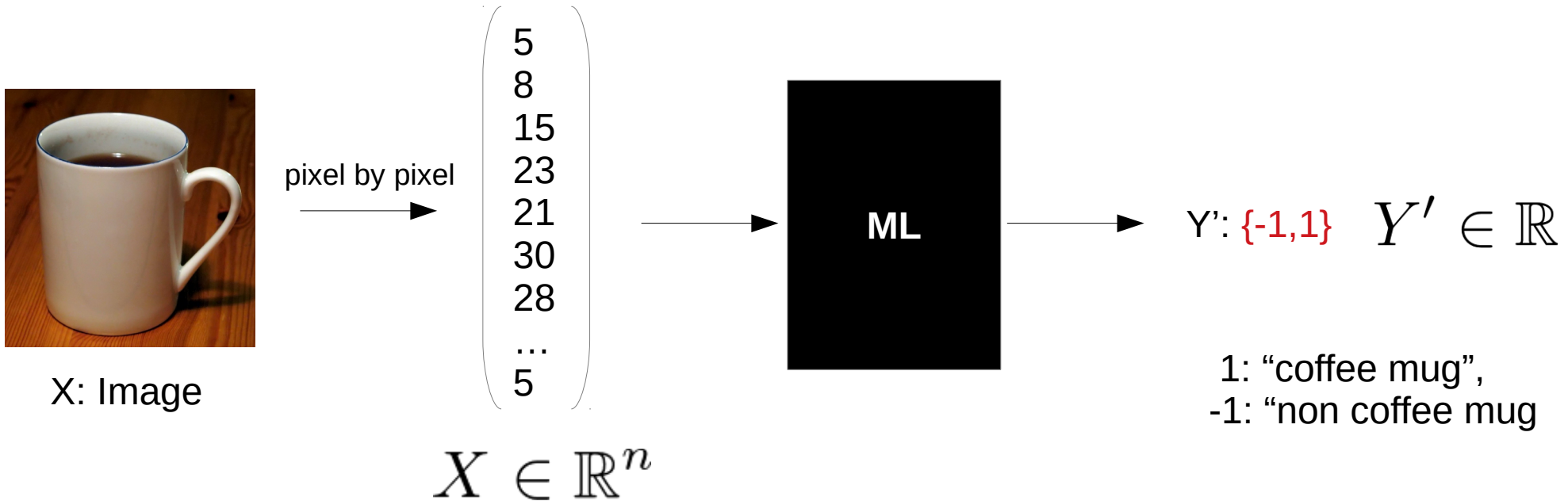
## Supervised Learning: Example: Classification



X: Image

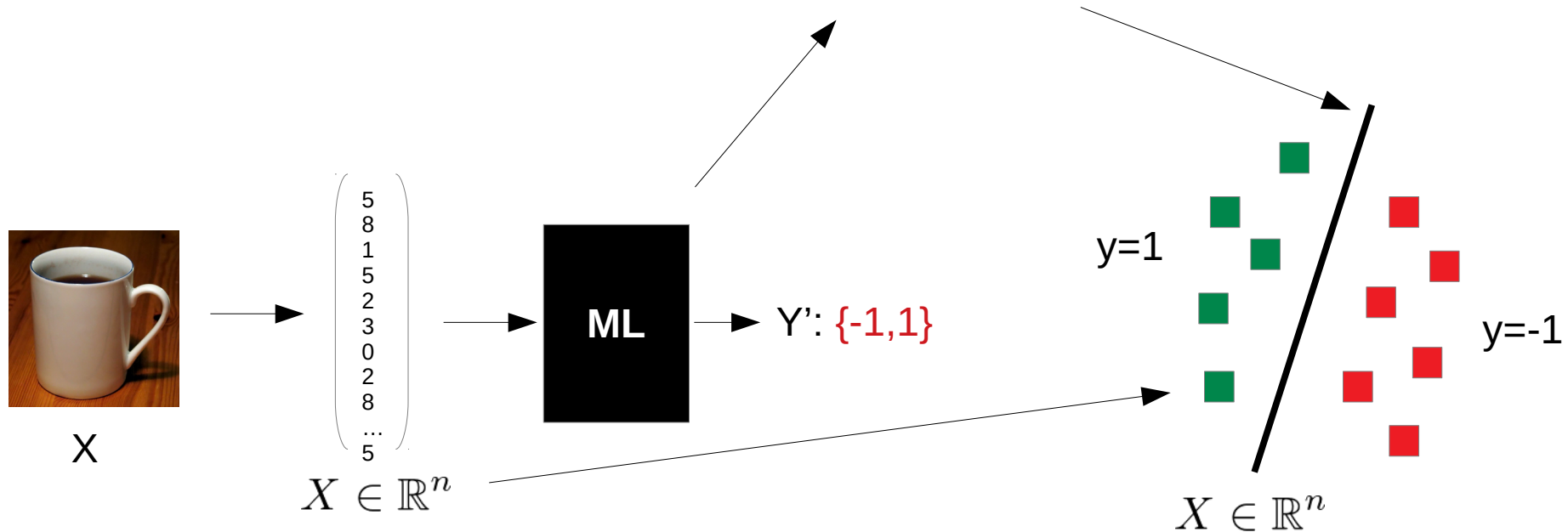


## Supervised Learning: Example: Classification



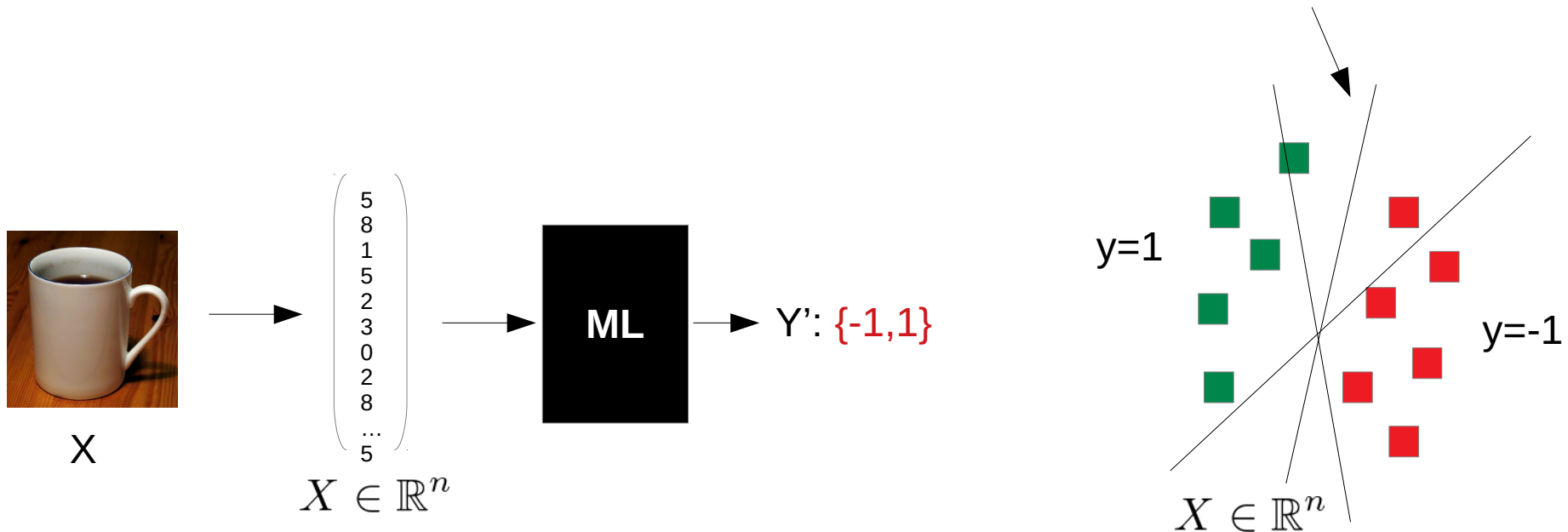
## Supervised Learning: Example: Classification

ML Model: function  $f$  separating mugs from rest



## Supervised Learning: Example: Classification

**LEARNING:** approximate „best“  $f$  for the given data





## Supervised Learning: Example: Classification

**LEARNING:** optimization problem:

$$\min(\|f(X, Y), Y'\|)$$



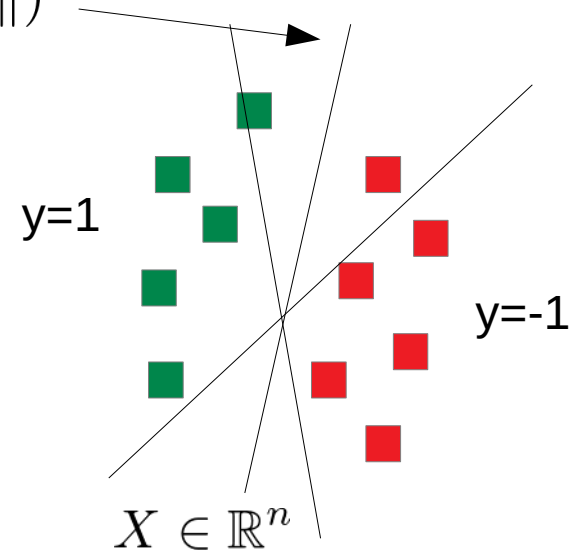
$X$

$$\begin{pmatrix} 5 \\ 8 \\ 1 \\ 5 \\ 5 \\ 2 \\ 3 \\ 0 \\ 2 \\ 8 \\ \dots \\ 5 \end{pmatrix}$$

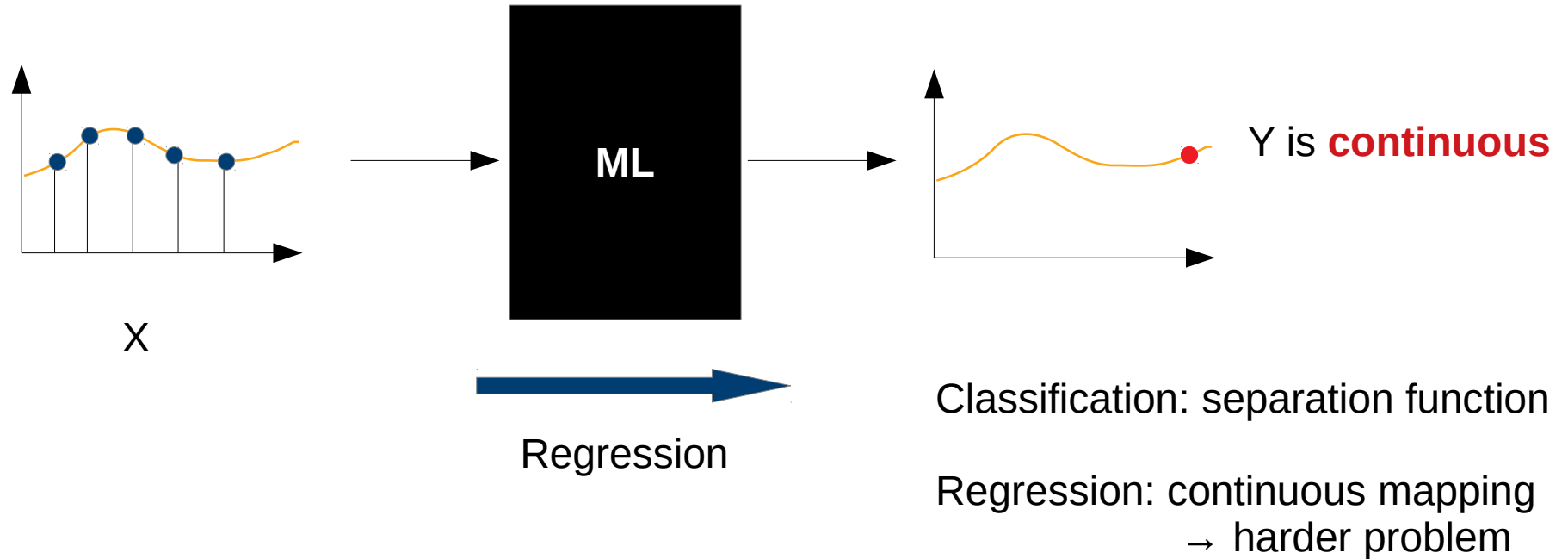
$X \in \mathbb{R}^n$



$Y': \{-1, 1\}$



## Supervised Learning: Example: Regression

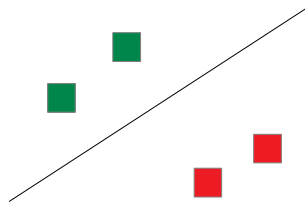


## Challenges of Supervised Learning

- Not only need data – also need to have  $Y \rightarrow$  human annotation
  - Getting “enough” labeled data is expensive
  - Sometimes impossible

UNDERFITTING

$$\min(\|f(X, Y), Y'\|)$$



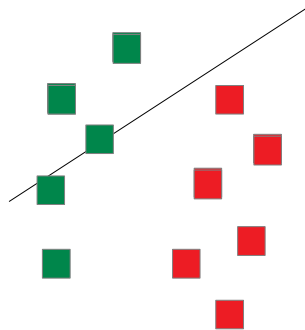
Training model  
On little data

## Challenges of Supervised Learning

- Not only need data – also need to have  $Y \rightarrow$  human annotation
  - Getting “enough” labeled data is expensive
  - Sometimes impossible

### UNDERFITTING

$$\min(\|f(X, Y), Y'\|)$$



→ bad sampling  
Of the data distribution

## Challenges of Supervised Learning

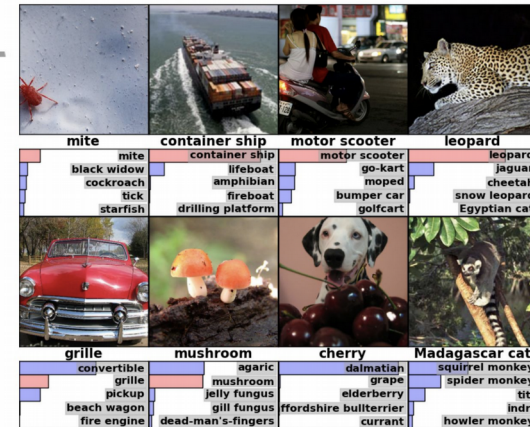
- Not only need data – also need to have  $Y \rightarrow$  human annotation
  - Getting “enough” labeled data is expensive
  - Sometimes impossible

## ImageNet Challenge

Example:

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.

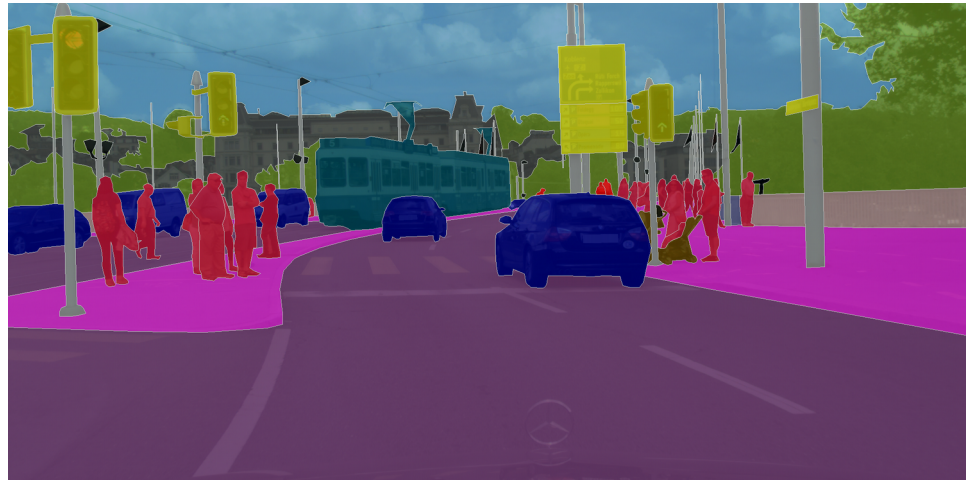
IMAGENET



## Challenges of Supervised Learning

- Not only need data – also need to have  $Y \rightarrow$  human annotation
  - Getting “enough” labeled data is expensive
  - Sometimes impossible

Example:



## Challenges of Supervised Learning

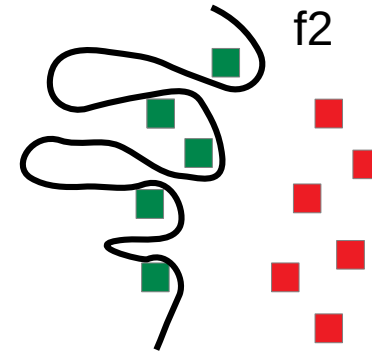
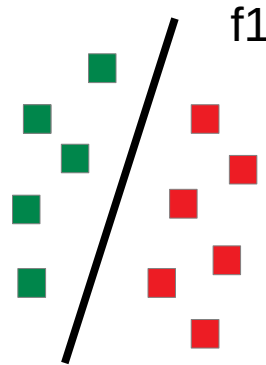
- Not only need data – also need to have  $Y \rightarrow$  human annotation
  - Getting “enough” labeled data is expensive
  - Sometimes impossible
- Training data is **only a sample**: prediction must work on **all data** → **generalization**

## Challenges of Supervised Learning

- Training data is **only a sample**: prediction must work on **all data** → **generalization**

Which model is better?

$$\min(\|f(X, Y), Y'\|)$$



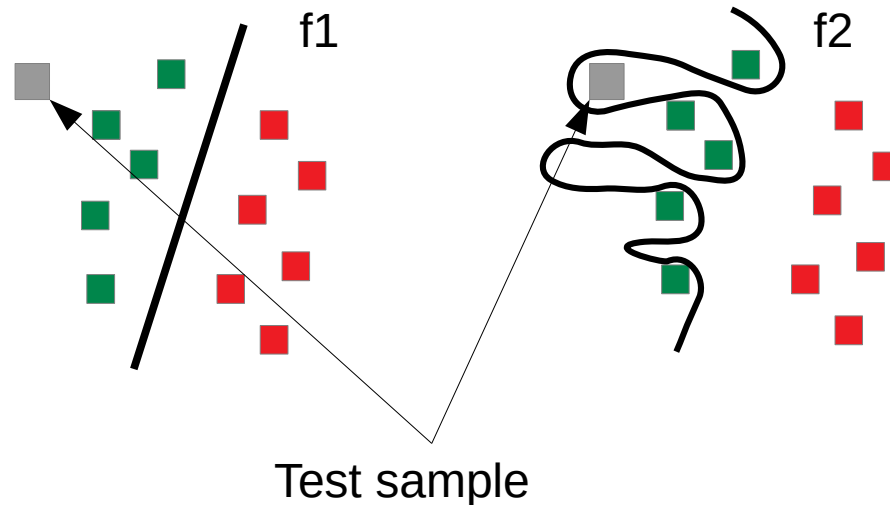


## Challenges of Supervised Learning

- Training data is **only a sample**: prediction must work on **all data** → **generalization**

Which model is better?

$$\min(\|f(X, Y), Y'\|)$$



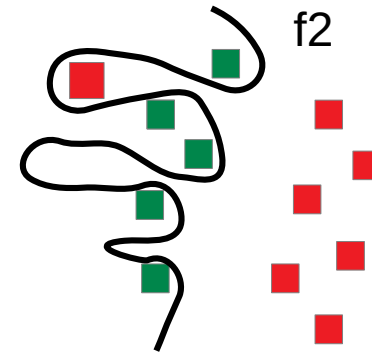
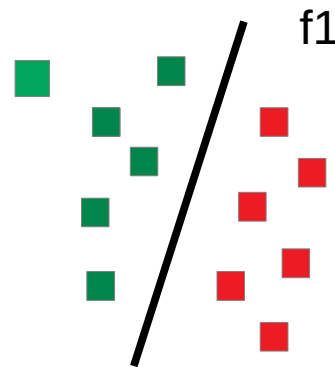
## Challenges of Supervised Learning

- Training data is **only a sample**: prediction must work on **all data** → **generalization**

### OVERFITTING

Model “to close” to train data

Very likely to happen in practice.  
→ we need to work against this...



## Data Preparation: Split into Train, Test, and Validate

A basic technique (we will learn more later) to at least detect overfitting is to split the available data into two or three subsets:

- Use unbiased **test set** for final evaluation of a model
- Use **train set** for model training
- **Validation set** (part of train set) can be used to optimize hyper parameters of the model

**Caution:** sets must be unbiased! (→ random sampling)  
In practice it can be hard to guarantee clean train/test sets:  
e.g. how to treat possible variance different data sources?  
→ statistical analysis needed!

**Basic evaluation** (more techniques to come)

**Train error:** measure of how well the model predicts the given labels

$$Err_{train} := \frac{1}{|X_{train}|} \sum_{x_i \in X_{train}} |f(x_i) - y_i|$$

low train error is the **necessary condition** for a “good” model

**Test error:** same as train error: low test error is the **sufficient condition**

$$Err_{test} := \frac{1}{|X_{test}|} \sum_{x_i \in X_{test}} |f(x_i) - y_i|$$

---

## Basic Types of Machine Learning Algorithms

**Supervised Learning**

**Unsupervised Learning**

**Reinforcement Learning**

- NO Labeled data
- NO Direct and quantitative evaluation
- Explore structure of data

## Unsupervised Learning

Data without “labels”  $(x_1, x_2, \dots, x_n)$

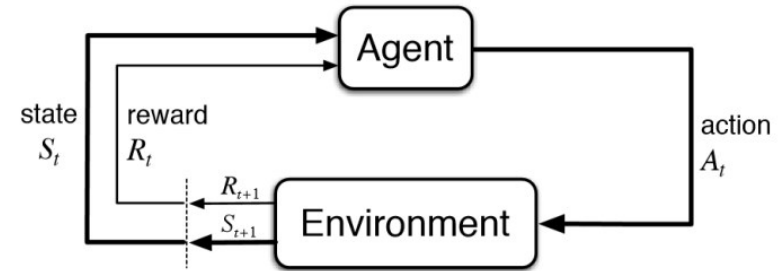
- Clustering
- Outlier Detection (e.g. Defect or Intrusion detection)

## Basic Types of Machine Learning Algorithms

Supervised Learning

Unsupervised Learning

Reinforcement Learning



- Learning decisions in an interactive environment
- State  $\leftrightarrow$  Action learning
- Game playing and robotics
- Hardly use in Data Science

Libraries used in this lecture:



Introduction in this week's lab



... introduction in Block 8