

Concepts of Learning From Data

Ruwen Qin
Stony Brook University

August 21, 2025

1 An Example

The following is an example of image data with two classes: with and without crack.

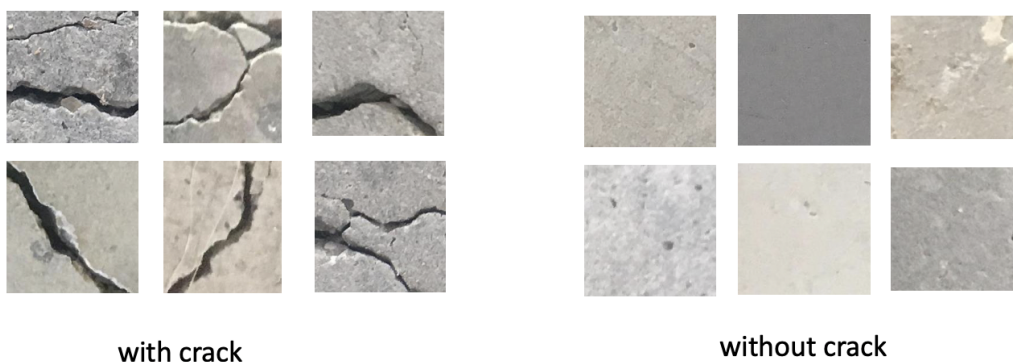


Figure 1: Image data of concrete surface with and without crack

Can you tell the class of each image below: with or without crack? Everyone including those having no domain knowledge can finish this task easily after reviewing the examples in [Figure 1](#).



Figure 2: Images to be classified

Can we teach machine (e.g., a computer) to perform the same task as you did just now? Learning from examples or comparison is a common way of gaining new knowledge. Machine learning is possible when a pattern exists but we cannot pin it down mathematically. If we have data, we can discover the pattern from the data.

This module, mainly developed based on the reference [\[1\]](#), introduces the basic concepts of machine learning to prepare you for studying related modules. Interested students can refer to that reference for details.

2 Learning Diagram

Figure 3 illustrates major components of (supervised) learning and solution:

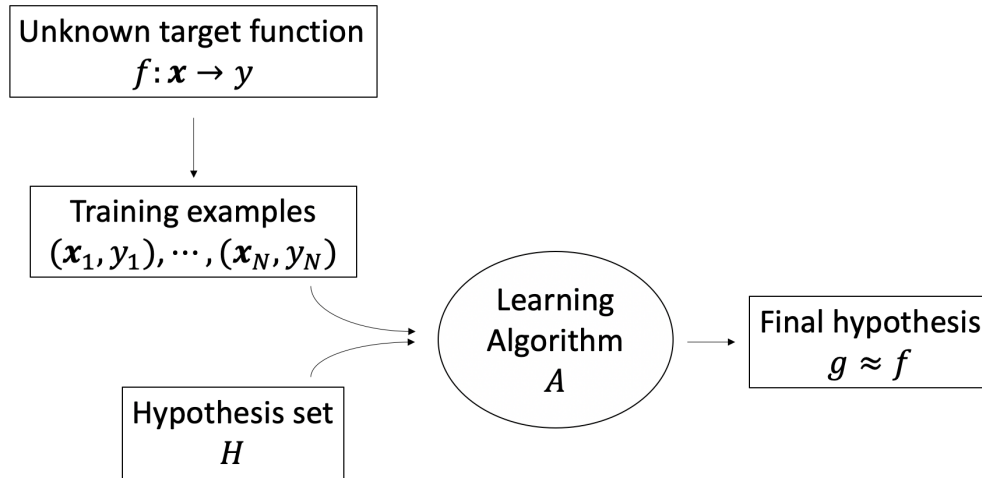


Figure 3: Learning and solution components

2.1 Learning Components

- Input, \mathbf{x}
- output, y
- Target function, f
- Training data: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
- Hypothesis: $g : \mathbf{x} \rightarrow y$

Regarding the example in Section 1, the input \mathbf{x} is image data, the output y is image class, the target function f is unknown and we would like to infer, and the training data $\{\mathbf{x}_n, y_n\}$ are illustrated in Figure 1. Given the training data, we will learn and find the final hypothesis g that predicts the output for an given input like those images in Figure 2.

2.2 Solution Components

The two solution components are:

- Hypothesis set: H
- Learning algorithm: A

In Figure 3, the hypothesis set and learning algorithm together are referred to as the learning model.

2.3 A Simple Hypothesis Set - Perceptron

Using the image classification problem as an exmaple. For an input, $\mathbf{x} = [x_1, \dots, x_d]$ (pixel values of an image),

$$y = \begin{cases} \text{crack}, & \sum_{i=1}^d w_i x_i + w_0 > 0 \\ \text{no-crack}, & \sum_{i=1}^d w_i x_i + w_0 \leq 0 \end{cases} \quad (1)$$

The function $h \in H$ can be written as

$$h = \text{sign}(\mathbf{w}^T \mathbf{x}) \quad (2)$$

where $\mathbf{w} = [w_0, w_1, \dots, w_d]^T$ are learnable parameters and $\mathbf{x} = [1, x_1, \dots, x_d]^T$ are the input.

2.4 Learning Algorithm

Each hypothesis is a hyper-plane in the space of input \mathbf{x} , which separates observations into two groups. The learning algorithm picks $g \approx h$ from the hypothesis set H . The final hypothesis g is the one that minimizes the number of misclassified points, as Figure 4 shows.

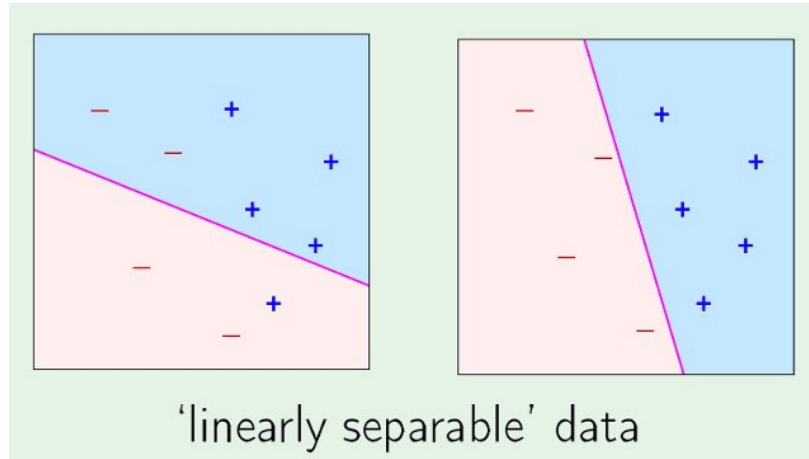


Figure 4: Linearly separable data

3 Learning Diagram with Error Measure

3.1 Error Measure

How well does a hypothesis $h(\mathbf{x})$ approximate y ? Point wise, the prediction error is the difference between the predicted output, $h(\mathbf{x}_n)$, and the correct output (i.e., ground truth), y_n :

$$e_n = h(\mathbf{x}_n) - y_n = h(\mathbf{x}_n) - f(\mathbf{x}_n). \quad (3)$$

A loss function measures the overall error. The in-sample loss function is

$$L_{in}(h) = \sum_{n=1}^N l(h(\mathbf{x}_n), f(\mathbf{x}_n)), \quad (4)$$

and the out-sample loss is

$$L_{out}(h) = \mathbb{E}_{\mathbf{x}}[l(h(\mathbf{x}), f(\mathbf{x}))]. \quad (5)$$

There are various choices of the function l that penalizes the prediction error. For example, if l is the first norm of errors, $\mathbf{e} = [e_1, \dots, e_N] = h(\mathbf{x}) - f(\mathbf{x})$, and then the loss function L is the sum of absolute errors.

3.2 The Role of Loss Function in Learning Algorithm

The loss function l provides a pointwise error measure of a hypothesis $h(\mathbf{x})$ with the ground truth $f(\mathbf{x})$. The learning algorithm uses the error function as a reference for searching the final hypothesis $g(\mathbf{x})$, as shown in Figure 5. The final hypothesis is the one that minimizes the loss function.

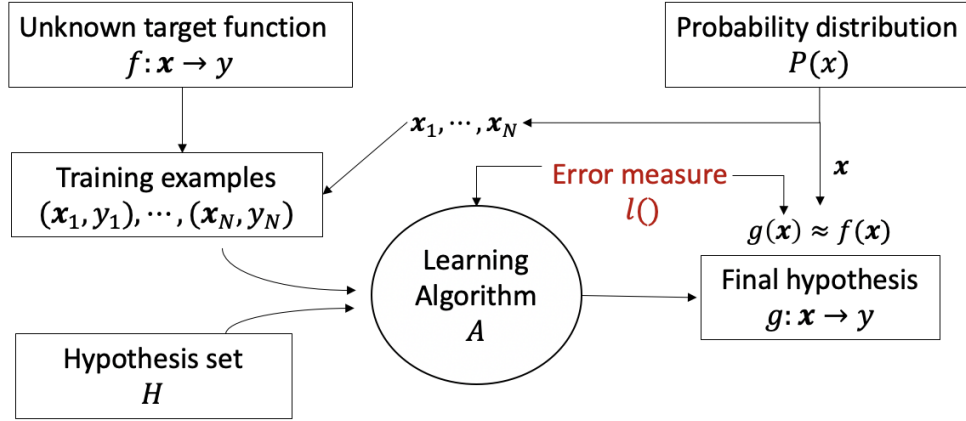


Figure 5: Learning diagram with error measure

4 Learning Diagram with a Noisy Target Function

The target function is not always a function¹. Mathematically speaking, two identical inputs may have different outputs. Therefore, we specify target distribution instead:

$$P(y|\mathbf{x}) \quad (6)$$

(\mathbf{x}, y) is generated from the joint distribution

$$P(y|\mathbf{x})P(\mathbf{x}). \quad (7)$$

Noisy target is the deterministic target function, $f(\mathbf{x}) = E[y|\mathbf{x}]$, plus noise, $y = f(\mathbf{x})$. The deterministic target function is a special case of the noisy target.

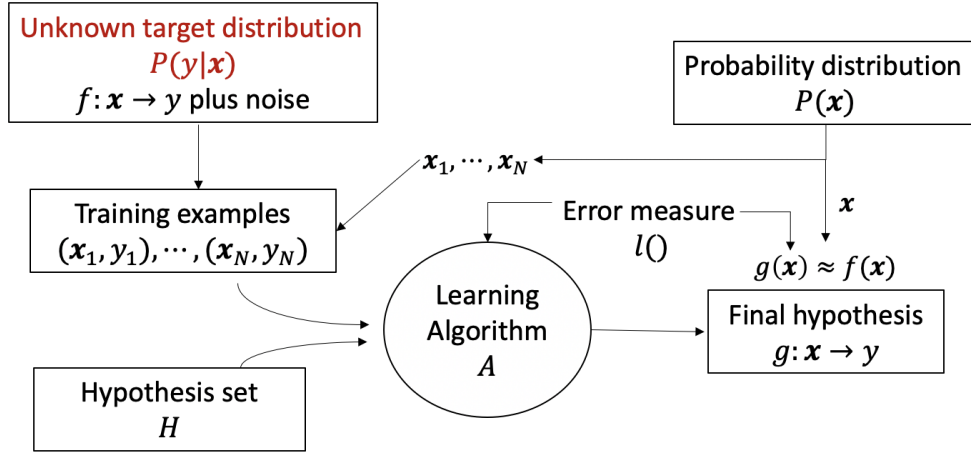


Figure 6: Learning diagram with noise

The target distribution $P(y|\mathbf{x})$ is what we are trying to learn. The input distribution $P(\mathbf{x})$ quantifies relative importance of \mathbf{x} . They both convey probabilistic aspects of input \mathbf{x} and output y .

¹The domain of a function, X , is a set of all values for which the function is defined. The codomain of the function Y , is the set into which all of the output of the function is constrained to fall. The function $f : X \rightarrow Y$ is a one-to-one mapping relationship between an element in the domain and an element in the codomain. That is, $\forall \mathbf{x} \in X$, f maps \mathbf{x} to one and only one $y \in Y$.

5 Nonlinear Hypothesis

Figure 7 contains two classes of data points, which are indicated by red cross and blue circle symbols. Can we separate those two classes of data points? Certainly we can, but the hypothesis that separates them is nonlinear.

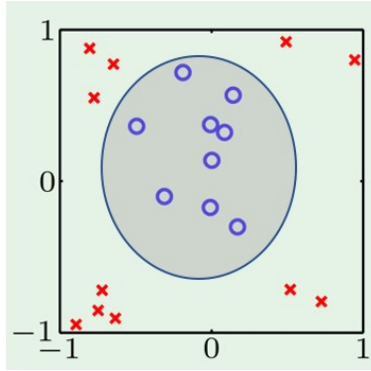


Figure 7: A Nonlinear Hypothesis

For some cases, we may transform the data so that we can separate them using a hyperplane.

- original input data $\mathbf{x}_1, \dots, \mathbf{x}_N$
- transform the data: $\mathbf{z}_n = \Phi(\mathbf{x}_n)$, for $n = 1, \dots, N$
- separate data in \mathcal{Z} -space with the hypothesis $\tilde{g}(\mathbf{z}) = \text{sign}(\tilde{\mathbf{w}}^T \mathbf{z})$
- The hypothesis in the \mathcal{X} -space: $g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x})) = \text{sign}(\tilde{\mathbf{w}}^T \Phi(\mathbf{x}))$

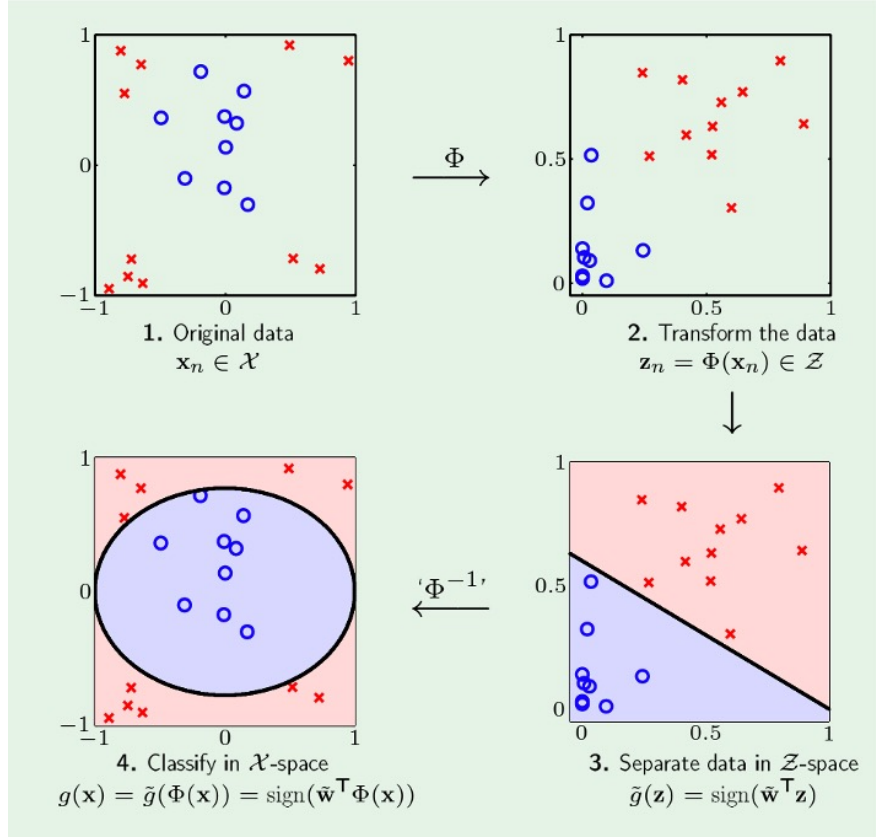


Figure 8: Transform function $\phi = \mathbf{x}^2$

6 Types of Learning

- Supervised learning: in the training dataset, every input \mathbf{x}_n comes with an output label y_n .
- Unsupervised learning: learning without y_n .
- Semi-supervised learning: Learning with some y_n .
- Reinforced learning: Explicit y_n is not available, but implicit \hat{y}_n with its goodness.

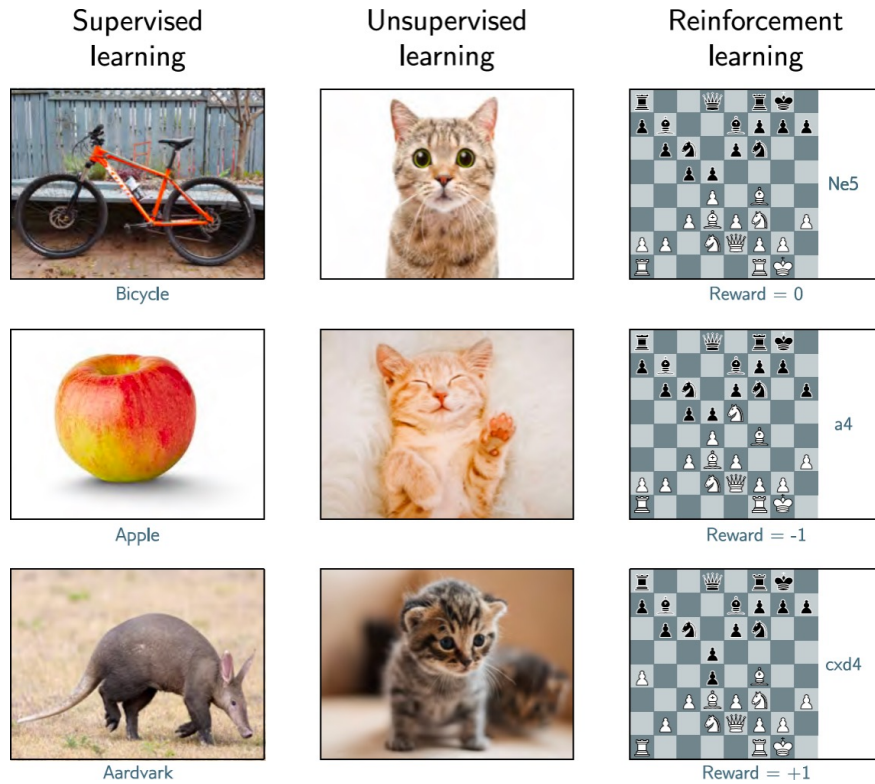


Figure 9: Different type of learning (image source: [2])

7 Some Machine Learning Problems

- *Classification*: this is the problem of assigning a category to each item.
- *Regression*: this is the problem of predicting a real value for each item.
- *Dimensionality reduction*: this problem consists of transforming an initial representation of items into a lower-dimensional representation while preserving some properties of the initial representation.
- *Clustering*: this is the problem of partitioning a set of items into homogeneous subsets.

References

- [1] Yaser S Abu-Mostafa, Malik Magdon-Ismail, and Hsuan-Tien Lin. *Learning From Data*, volume 4. AMLBook New York, 2012.
- [2] Simon JD Prince. *Understanding Deep Learning*. MIT press, 2023.