

# Machine Learning Basics

A self-study materials for PRML [1]

Jisung Lim<sup>1</sup>

<sup>1</sup>B.S. Candidate of Industrial Engineering  
Yonsei University, South Korea.

28th January, 2017

# Summary

## 1 Introduction

- What is Machine Learning?
- Polynomial Curve Fitting
- Chapter Objectives

## 2 Probability Theory

- subsection name

## 3 Reference

# Optical Character Recognition

## OCR Problem and Some Approaches



FIGURE 1: The MNIST data-base.

- **Input** :  $28 \times 28$  pixel image, represented by a vector  $x$  comprising 784 real numbers.
- **Goal** : To build a machine that will take such a vector  $x$  as input and that will produce the identity of the digit  $0, \dots, 9$  as the output.

Consider the example of recognizing handwritten digits, illustrated in Figure 1. This is the nontrivial problem due to the wide variability of handwriting. And we can tackle this problem using following approaches :

- 1 Handcrafted rules or heuristics
- 2 Machine learning methods

In practice, the former leads to a proliferation of rules and invariably gives poor results. For better results, the latter can be considered an alternative.

# Machine Learning Approach

**Machine Learning Approach** A machine learning approach can be divided into three major stages : **training**, **testing**, and **predicting**.

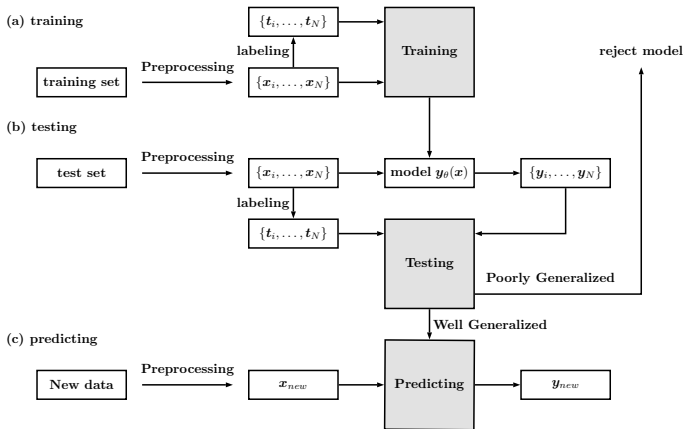


FIGURE 2: The description of full machine learning approach.



# ML Approach — Input

First of all, we need to define inputs more precisely. In OCR, inputs are handwriting images of a digit. In practical applications, the original input images may have different sizes, ratios, angles, or even colors. Hence, the images are typically scaled, rotated, translated, and grey-scaled so that each digit is contained within a box of a fixed-size, says  $28 \times 28$ , and is toned with greyscale color. And then, the preprocessed input images can be represented as a vector  $\mathbf{x} = (x_1, \dots, x_{784}) \in \mathbb{R}^n$  where the greyscale color of  $i$ th pixel is a  $x_i$ . Finally, you should label the true number  $t$  for each input image  $\mathbf{x}$ .

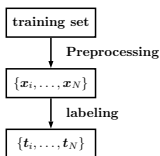


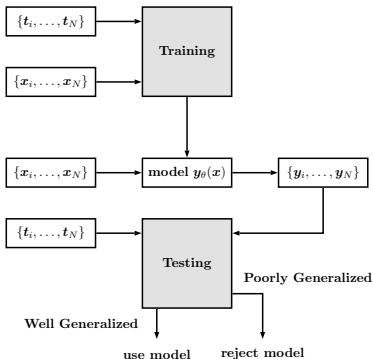
FIGURE 3

**Preprocess :** To transform the original input variables into some new space of variables so that the problem becomes easier to solve.

**Example  $\mathbf{x}_i$  :** An example is a collection of features. Typically, an **example** will be represented as a vector  $\mathbf{x}_i \in \mathbb{R}^n$  where each entry  $x_i$  is a **feature**.

**Label  $t_i$  :** A **label** represents the identity of the corresponding digit. Typically, the labels are hand-labelled by inspecting each image individually. Note that there is one such **label  $t_i$**  for each **example  $\mathbf{x}_i$** .

# ML Approach — Train and Test



**FIGURE 4:** The description of training and test.

## ■ Training

To tune the parameters of an adaptive model, it uses a large set of  $N$  **examples**  $\{x_1, \dots, x_N\}$  with its corresponding **labels**  $\{t_1, \dots, t_N\}$ . The model can be expressed as a function  $y(x)$  which takes a new digit image  $x_{new}$  as input and that generates an output vector  $y_{new}$ .

## ■ Testing

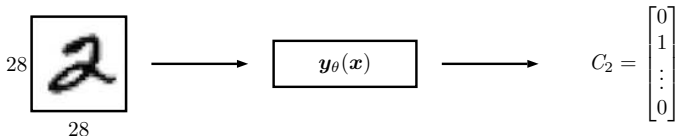
Once the model is trained it can then determine the identity of new digit images, which are said to comprise a test set. The ability to categorize correctly new examples that differ from those used for training is known as generalization. In practical applications, the variability of the input vectors will be such that the training data can comprise only a tiny fraction of all possible input vectors, and so generalization is a central goal in pattern recognition.

# ML Approach — Prediction

Now the model can work for new, unlabeled inputs. In OCR, prediction stage can be described as below.

**Input:** example  $x_{new}$

**output:** label  $\hat{t} = y_{new}$



**FIGURE 5:** The description of prediction in OCR.



# ML Approach — SL, UL, and RL

## SL Supervised Learning :

Corresponding target vectors are known.

Ex. Regression, Classification.

## UL Unsupervised Learning :

The training data consists of a set of input vector  $x$  without any corresponding target values.

Ex. Clustering, Density estimation.

## RL Reinforcement Learning :

Finding **suitable actions** to take in a given situation in order to maximize a reward. It discovers the optimal by a process of **trial and error**.

A general feature of reinforcement learning is the trade-off between **exploitation**, in which the system **tries out new kinds of actions** to see how effective they are, and **exploitation**, in which the system **makes use of actions** that are known to yield a high reward.





# Precise Definition of ML

## Definition of machine learning

A computer program is said to **learn** from **experience**  $E$  with respect to some class of **tasks**  $T$  and **performance measure**  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .



# Polynomial Curve Fitting

## Given Situation

- $N$  observations  $\mathbf{X} \equiv (x_1, \dots, x_N)^T$ , together with corresponding labels  $\mathbf{t} \equiv (t_1, \dots, t_N)^T$
- $(x_i, t_i)$  possess underlying regularity, which we wish to learn.
- Individual  $t_i$  observations are corrupted by random noise, which might arise from intrinsic stochasticity due to there being source of variability.

## Our Goal and Intrinsic Difficulty

- **Goal**  
exploit this given training set in order to make prediction of the value  $\hat{t}$  of the target variable for some new value  $\hat{x}$  of the input variable.
- **Intrinsic Difficulty**  
This is intrinsically a difficult problem not only as we have to generalize from a **finite data** set but also, for a given  $\hat{x}$  there is **uncertainty** as to the appropriate value for  $\hat{t}$ .

- 1 minimizing an error function.

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2$$

$$\text{where } y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots$$

- 2 use regularization for preventing overfitting.

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$



# Chapter Objectives

## ■ Probability Theory

Provides a consistent framework for expressing such uncertainty in a precise and quantitative manner.

## ■ Decision Theory

allows us to exploit this probabilistic representation to make optimal prediction according to appropriate criteria.

## ■ Information Theory





Christopher M. BISHOP. **Pattern Recognition and Machine Learning**. Springer, 2006.