

Lecture 1

- Big Picture of Statistical Learning
- Types of Learning Problems
- Statistical Decision Theory
- Nearest - Neighbor Methods
- Linear Regression Model

learning from data

next time

1

Statistical learning is a subfield of modern statistics that develops methods for understanding data, making predictions based on the data, and making data-informed decisions. Statistical learning plays a key role in many areas of science, engineering, finance, and medicine. As a result, people with SL skills are in high demand.

Nowadays the crucial role of data in scientific and engineering discoveries is obvious. But this has not been always the case!

- Plato (around 380 BC): astronomers should study stars ("decorations on a visible surface") by "reason and thought", without looking at the stars. *The Republic*, Book VII.
- Aristotle (384 - 322 BC): wrote a lot about biology and medicine, but never looked at the real data about human body. He believed that the heart, not the brain, was the seat of thought. ←

Of course, with time, people began to recognize the importance of data. One remarkable example of "learning from data":

- Johannes Kepler (1571 - 1630) discovered the laws of planetary motion (around 1610) from the extensive astronomical observations collected by Tycho Brahe (1546 - 1601).

This could be easily disproved by observing the effects of trauma to the brain.

(This was known by Hippocrates, who lived before Aristotle, and his school)

The field of Statistics, solely dedicated to data analysis, has emerged from the work of

- John Graunt (1620 - 1674), who estimated the population of London by analyzing the statistics of deaths. He published his results in 1662.

Since then, the field of classical statistics has been developed by many prominent scientists and mathematicians such as Francis Galton, Karl Pearson, Ronald Fisher, William Gosset (aka Student), Egon Pearson (son of Karl Pearson), and Jerzy Neyman.

The availability of cheap computing power and the explosion of data truly revolutionized classical statistics:

- "Classical Statistics": problems come from surveys of human populations, agricultural and industrial experiments; data sets are small; solutions are analytical, done by hand.
- "Modern Statistics": problems come from everywhere; data sets are large; solutions are numerical, done on a computer.

Since computers are equally available to all research communities and data appear everywhere, modern methods for learning from data have been and continue to be developed by researchers from various fields :

Classical Field Recent Subfield (focused on learning from data)

Statistics	→	Statistical Learning
Computer Science	→	Machine Learning
Engineering	→	Pattern Recognition
Biology/Medicine	→	Bioinformatics
Finance	→	Financial Machine Learning

All these subfield can be viewed as different facets of the same development.
Remark: For example, the wikipedia article on "Statistical Learning" redirects to "Machine Learning".
Remark: Each of those subfields has its own terminology, notation, and jargon. This complicates communication between researchers from different subfields.

Main Goal of the Course : to discuss the most fundamental ideas and methods for learning from data and explain them in a statistical framework.

The course can be viewed as :

- Natural continuation of my IDS/ACM/CS 157 (Statistical Inference).
- Statistical and theoretical analog of CMS/CS/CNS/EE/IDS 155 (ML & Data Mining)

Types of Learning Problems

There are three main types of learning:

① Supervised Learning

This is the most studied and most often used in applications type of learning.

input $X \rightarrow$ Complex System \rightarrow Y output

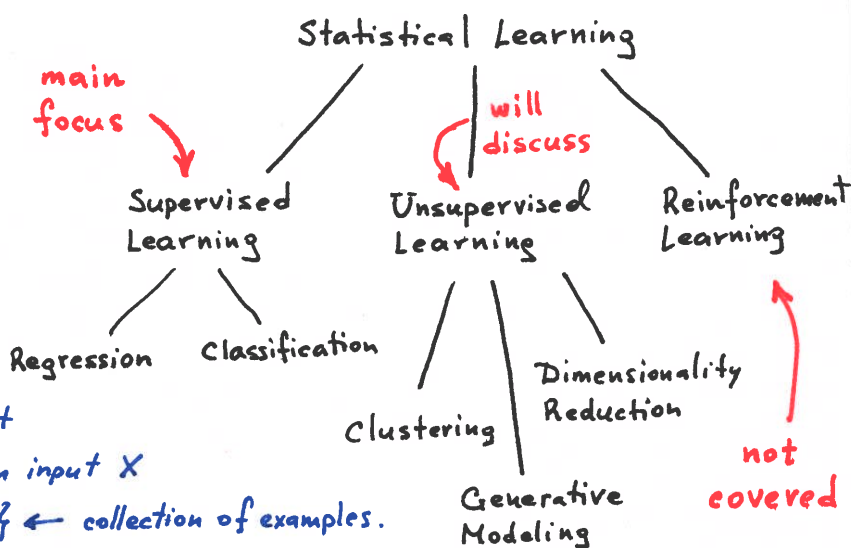
The goal of supervised learning is to construct a statistical model for predicting output Y from input X based on "training data" : $\{(X_1, Y_1) \dots (X_N, Y_N)\} \leftarrow$ collection of examples.

- If output Y is quantitative (continuous) \Rightarrow the problem is called regression.

Example : Predict the price of a stock Y in 1 monts from now, based on the input X that describes the company performance measures and economic data.

- If output Y is qualitative (discrete, categorical) \Rightarrow the problem is called classification.

Example : Predict whether a patient, hospitalized due to a heart attack, will have a second heart attack ($Y=1$) or not ($Y=0$), based on the input $X =$ (demographic, diet, clinical measurements).



② Unsupervised Learning

Here the situation is more challenging: the training data consists only of inputs $\{X_1, \dots, X_N\}$, there is no outputs Y_1, \dots, Y_N that could supervise our learning.

The goal of unsupervised learning is less well defined and more ambiguous: to understand the structure of the data, find possible patterns among inputs.

More specific problems include:

- Clustering: to divide $\{X_1, \dots, X_N\}$ into relatively distinct groups (clusters).

Example: market segmentation: divide customers into groups based on their characteristics (income, zip code, shopping habits) = X_i

- Dimensionality Reduction: to represent the data $\{X_1, \dots, X_N\}$ in a smaller space, $\{X_i \in \mathbb{R}^D\} \rightsquigarrow \{\tilde{X}_i \in \mathbb{R}^d\}$, $d \ll D$, for visualization or as a pre-processing step for supervised learning.



- Generative Modeling: to model a generating mechanism that produced $\{X_1, \dots, X_N\}$, in order to generate artificial inputs that are similar to real data. (synthetic data in the input space)

Example: synthetic earthquake generation (accelerograms similar to real ones X_1, \dots, X_N)

③ Reinforcement Learning

The goal of reinforcement learning is to infer optimal sequential decisions (actions) based on rewards or punishments received as a result of previous actions.

Example: training a robot to navigate a given environment in the presence of obstacles by penalizing decisions that result in collisions.

Example: training an algorithm to play chess (AlphaZero)

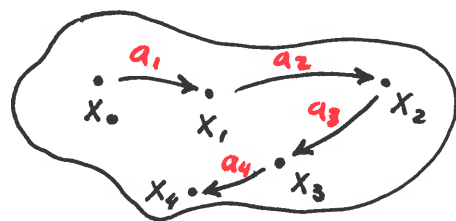
The following aspects make RL especially difficult:

- Choosing $a_i \in \mathcal{A}_{X_{i-1}}$ is done not based on optimization $R_{a_i}(X_{i-1}, X_i) \rightarrow \max$ but based on previous examples. $\mathcal{A}_{X_{i-1}}$ is often very large (or ∞).

After 9 hours of training via self-play, AlphaZero defeated Stockfish 8, the strongest chess engine at that time (2017)

Ref: D. Silver et al (2018), Science, vol 362, pp 1140-44.

Schematic Representation of RL



Input space (aka Environment)

X_0 is the initial input (aka state)

- Choose $a_1 \in \mathcal{A}_{X_0}$ ← set of available actions in state X_0
- Execute a_1 : $X_0 \xrightarrow{a_1} X_1$ (new state)
- Get reward $R_{a_1}(X_0, X_1) \in [-1, 1]$

bad good

- Executing a_i on X_{i-1} does not always bring us to deterministic state X_i ; Environment may change during execution. So, executing a_i on X_{i-1} results into a probability distribution $IP(X_i = x | X_{i-1}, a_i)$ on the input space.
- Action a_i may affect not only the immediate reward $R_{a_i}(X_{i-1}, X_i)$, but also rewards at all subsequent steps.

Let's start with a general discussion of supervised learning.

We are going to use the following notation and terminology:

- $X = (X_1, \dots, X_p)^T$ is a vector of inputs (aka predictors, features, independent variables)
 $p = \# \text{ inputs}$. We use X to refer to a generic input vector.

- $x_i = (x_{i1}, \dots, x_{ip})^T$ is the i^{th} observed value of X , $i = 1, \dots, N$.

$N = \# \text{ observations}$. $x_i \in \mathbb{R}^p$ is a p -vector.

- $x^{(j)} = (x_{1j}, \dots, x_{Nj})^T$ is the vector of all observations of input X_j , $j = 1, \dots, p$.

$$X = \underbrace{\begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & & \vdots \\ x_{N1} & \dots & x_{Np} \end{bmatrix}}_{N \times p \text{ matrix}} = \underbrace{\begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}}_{\text{inputs}} = \begin{bmatrix} x^{(1)} & \dots & x^{(p)} \end{bmatrix}$$

\uparrow observations of X_1 \uparrow observations of X_p

Remark: All vectors are assumed to be column vectors.

$x_{ij} = \text{value of } j^{\text{th}} \text{ input in the } i^{\text{th}} \text{ observation}$.

- Y is a quantitative output corresponding to X (aka response, dependent variable)
 $Y \in \mathbb{R}$

- y_i is the i^{th} observed value of Y that correspond to x_i .

- G is a qualitative output corresponding to X , $G \in \mathcal{G}$ set of groups/classes.

- g_i is the i^{th} observation of G that correspond to x_i .

Main Goal: Given the training data $\underbrace{(x_1, y_1), \dots, (x_N, y_N)}_{\text{regression}}$ (or $\underbrace{(x_1, g_1), \dots, (x_N, g_N)}_{\text{classification}}$) and the value of input X , make a good prediction \hat{Y} (or \hat{G}) of the output Y (or G).

Statistical Decision Theory provides a framework for developing powerful prediction methods.

Let's start with the case of regression (quantitative output).

Let $X \in \mathbb{R}^p$ be a random input and $Y \in \mathbb{R}$ be the corresponding random output.

We want to find a function (prediction rule) $f: \mathbb{R}^p \rightarrow \mathbb{R}$ such that

$\hat{Y} = f(X) \approx Y$ is an accurate prediction.

To measure the accuracy of prediction (the goodness of $f(x)$),

we need a loss function $L: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$. If the expected loss is

small, the prediction rule is good; otherwise it is bad:

$$E[L(Y, \hat{Y})] = E[L(Y, f(X))] = \begin{cases} \text{small} \Rightarrow f \text{ is good,} \\ \text{large} \Rightarrow f \text{ is bad.} \end{cases}$$

By far the most popular and often used loss function is the squared error loss : $L(Y, \hat{Y}) = (Y - \hat{Y})^2$ aka L_2 loss function.

The corresponding expected loss is then called the mean squared error (MSE) or expected prediction error (EPE) :

$$MSE(f) = IE[(Y - \hat{Y})^2] = IE[(Y - f(X))^2] \leftarrow \begin{array}{l} \text{this expected value is w.r.t.} \\ \text{the joint distribution of } X \text{ and } Y. \end{array}$$

So, our goal is to find the prediction rule f that minimizes $MSE(f)$.

We can find f by conditioning on X and using the law of total expectation :

$$MSE(f) = IE \left[\underbrace{IE[(Y - f(X))^2 | X]}_{\substack{\uparrow \text{wrt } X \\ \uparrow \text{wrt } Y}} \right] \quad \begin{array}{l} \text{this is a random variable, whose value is} \\ IE[(Y - f(x))^2 | X = x] \text{ when } X = x. \\ \text{conditional expectation} \end{array}$$

The inner conditional expectation :

$$IE[(Y - f(x))^2 | X = x] = \int (y - f(x))^2 \underbrace{p_{Y|X}(y|x)}_{\text{conditional PDF of } Y \text{ given } X} dy$$

$$= IE[(Y - IE[Y|X=x] + IE[Y|X=x] - f(x))^2 | X = x]$$

$$= \underbrace{IE[(Y - IE[Y|X=x])^2 | X = x]}_{V[Y|X=x]} + \underbrace{IE[(IE[Y|X=x] - f(x))^2 | X = x]}_{\text{constant}}$$

$$+ 2 IE[(Y - IE[Y|X=x]) \cdot \underbrace{(IE[Y|X=x] - f(x))}_{\text{constant}} | X = x]$$

$$= V[Y|X=x] + (IE[Y|X=x] - f(x))^2 + 2 (IE[Y|X=x] - f(x)) \underbrace{(IE[Y|X=x] - IE[Y|X=x])}_0$$

So : $IE[(Y - f(x))^2 | X = x] = V[Y|X=x] + (IE[Y|X=x] - f(x))^2$, and

$$MSE(f) = IE[V[Y|X] + (IE[Y|X] - f(X))^2]$$

Therefore, the prediction rule that minimizes $MSE(f)$ is

$$f(x) = IE[Y|X]$$

in the MSE sense
(L_2 sense)

In other words, if we observe that input $X=x$, then the best prediction \hat{Y} for

the output is $f(x) = IE[Y|X=x]$

this conditional expectation is called the regression function.