

# ECE595 / STAT598: Machine Learning I

## Lecture 22 Is Learning Feasible?

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Learning Theory

- Welcome to Part 3 of ECE595 / STAT598!
- Here is what we have learned:
  - Part 1: The machine learning pipeline
  - Part 2: Classification methods
- What are we going to do in Part 3?
  - Now that we have a method, then what?
  - Will it do well? How well?
  - Will it fail? When?
  - Complex model = better?
  - More sample = better?
  - Can every problem be solved by learning?
  - When do you overfit?
  - How to avoid overfit?

# Outline

## Today's Lecture:

- What constitutes a learning problem?
  - Training and testing samples
  - Target and Hypothesis function
  - Learning Model
- Is learning feasible?
  - An example
  - The power of probability
- Training versus Testing
  - In-sample error
  - Out-sample error
  - Probability bound

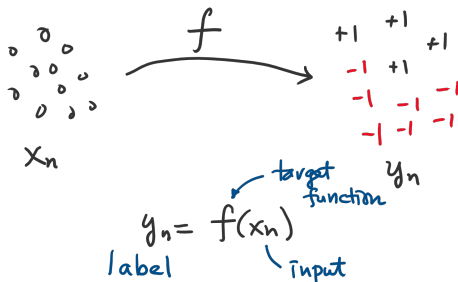
## Reference:

- Learning from Data, chapter 1.3

# Dataset

Let us first talk about a dataset:

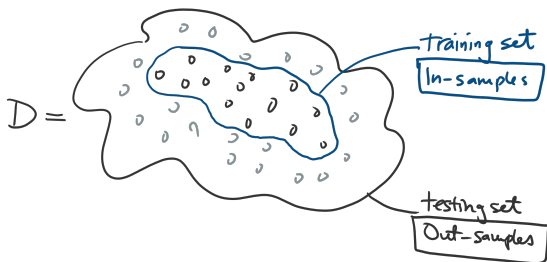
- Input vectors:  $\mathbf{x}_1, \dots, \mathbf{x}_N$
- Labels:  $y_1, \dots, y_N$
- Training set:  $\mathcal{D}$
- Target function  $f$ : Maps  $\mathbf{x}_n$  to  $y_n$
- Target function is always **unknown** to you



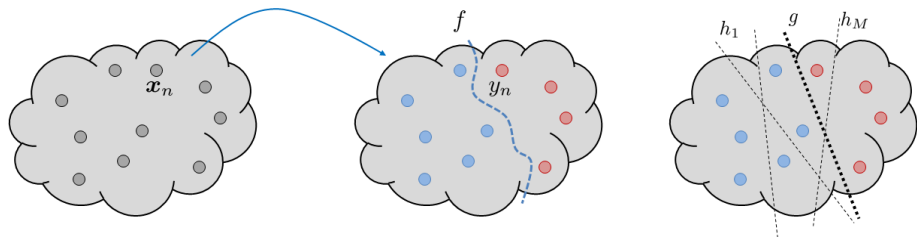
# Training and Testing Set

Let us first talk about a dataset:

- **In-sample:** Samples that are inside the training set
- **Out-sample:** Samples that are outside the training set

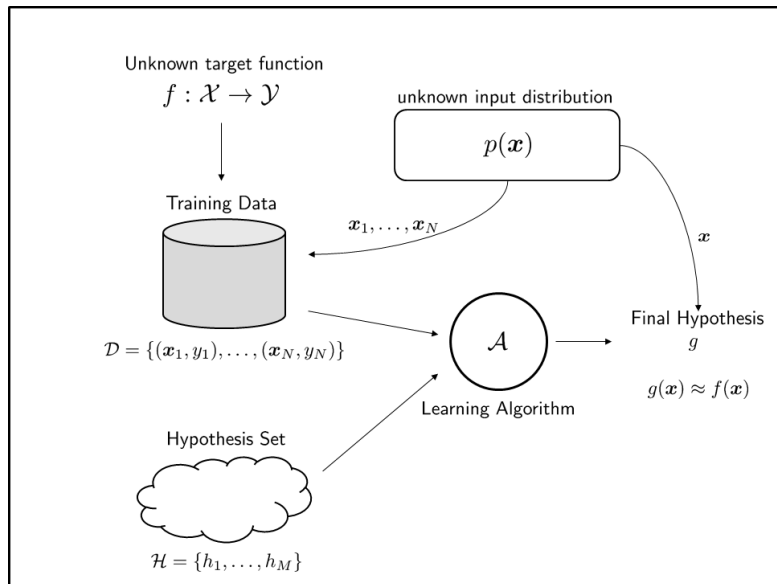


# Hypothesis Function



- Hypothesis set:  $\mathcal{H} = \{h_1, \dots, h_M\}$ : Possible decision boundaries
- Algorithm: Picks  $h_m$  from  $\mathcal{H}$
- Final hypothesis:  $g$ : The one you found

# Learning Model



## Today's Lecture:

- What constitutes a learning problem?
  - Training and testing samples
  - Target and Hypothesis function
  - Learning Model
- Is learning feasible?
  - An example
  - The power of probability
- Training versus Testing
  - In-sample error
  - Out-sample error
  - Probability bound



# Is Learning Feasible?

## In-sample and Out-sample:

- In-sample: Training Data
- Out-sample: Testing Data

## When can we claim "learning is feasible"?

Suppose we have a training set  $\mathcal{D}$ , can we learn the target function  $f$ ?

- "Learn" means: I use the data you give me to come up with an  $f$
- "Successful" means: All in-samples are correctly predicted
- And all out-samples are also correctly predicted
- If YES, then we are in business.
- Learning is feasible!
- If NO, then we can go home and sleep.
- There is just no way to learn  $f$  from  $\mathcal{D}$ .

## Example

- Let  $\mathcal{X} = \{0, 1\}^3$
- Each  $\mathbf{x} \in \mathcal{X}$  is a binary vector
- E.g.,  $\mathbf{x} = [0, 0, 1]^T$  or  $\mathbf{x} = [1, 0, 1]^T$
- How many possible vectors are there?  $2^3 = 8$
- Call them  $\mathbf{x}_1, \dots, \mathbf{x}_8$
- There is a target function  $f$
- $f$  maps every  $\mathbf{x}$  to a  $y$
- $y \in \{+1, -1\}$
- E.g.,  $f([0, 0, 1]) = +1$ ,  $f([0, 1, 1]) = -1$ , etc.
- How many possible  $f$ 's?
- You can think of  $f$  as a 8-bit vector
- E.g.,  $f = [+1, -1, -1, -1, +1, +1, +1, -1]$ .
- So there are  $2^8 = 256$  possible  $f$ 's.

## Example

- We have 8 input vectors:  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_8\}$
- We have 256 hypotheses:  $\mathcal{H} = \{h_1, \dots, h_{256}\}$
- Is learning feasible?
- Give me a subset  $\mathcal{D} \subset \mathcal{X}$ , can I find a hypothesis  $g \in \mathcal{H}$  such that  $g = f$ ?
- Suppose here is what you are given:  $\circ = -1$ ,  $\bullet = +1$ . You know 6 out of 8. These are the training data.

$\mathbf{x}_n$	$y_n$
0 0 0	$\circ$
0 0 1	$\bullet$
0 1 0	$\bullet$
0 1 1	$\circ$
1 0 0	$\bullet$
1 0 1	$\circ$
1 1 0	?
1 1 1	?

## Possibility 1

$x_n$			$y_n$
0	0	0	○
0	0	1	●
0	1	0	●
0	1	1	○
1	0	0	●
1	0	1	○
1	1	0	○
1	1	1	○

- One 1's will give me ●; Others give me ○
- So the last two entries should be ○

## Possibility 2

$x_n$			$y_n$
0	0	0	○
0	0	1	●
0	1	0	●
0	1	1	○
1	0	0	●
1	0	1	○
1	1	0	○
1	1	1	●

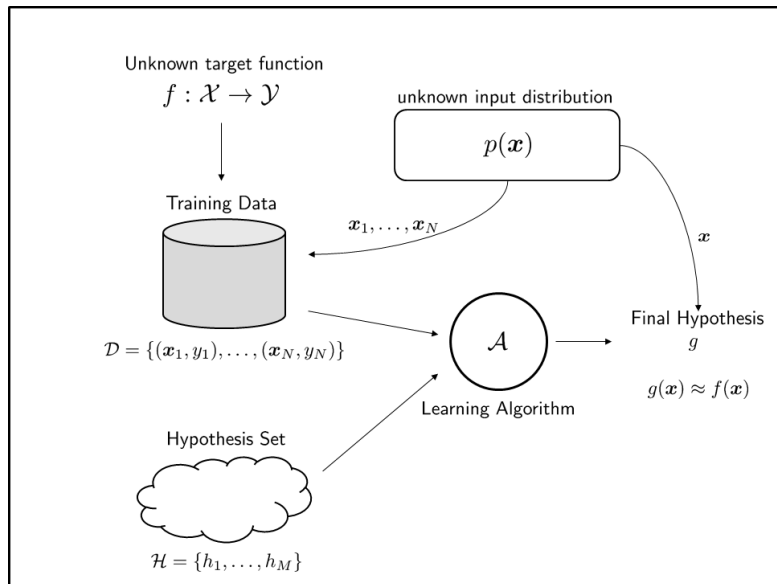
- Odd numbers of 1's give me ●
- Even numbers of 1's give me ○
- So [1 1 0] should be ○
- So [1 1 1] should be ●

## All the Possibilities

$x_n$	$y_n$	$g$	$f_1$	$f_2$	$f_3$	$f_4$
0 0 0	○	○	○	○	○	○
0 0 1	●	●	●	●	●	●
0 1 0	●	●	●	●	●	●
0 1 1	○	○	○	○	○	○
1 0 0	●	●	●	●	●	●
1 0 1	○	○	○	○	○	○
1 1 0		○/●	○	●	○	●
1 1 1		○/●	○	○	●	●

- $f_1, f_2, f_3, f_4$  are the only hypotheses you need to consider
- You just don't know which one out of the four to choose!
- You won't do better than random guess.
- So you haven't learned anything from the training data.
- Learning is infeasible.

# The Power of Probability



## Today's Lecture:

- What constitutes a learning problem?
  - Training and testing samples
  - Target and Hypothesis function
  - Learning Model
- Is learning feasible?
  - An example
  - The power of probability
- **Training versus Testing**
  - **In-sample error**
  - **Out-sample error**
  - **Probability bound**



## In-Sample Error

- Let  $\mathbf{x}_n$  be a *training* sample
- $h$ : Your hypothesis
- $f$ : The unknown target function
- If  $h(\mathbf{x}_n) = f(\mathbf{x}_n)$ , then say training sample  $\mathbf{x}_n$  is correctly classified.
- This will give you the **in-sample error**

### Definition (In-sample Error / Training Error)

Consider a training set  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and a target function  $f$ . The **in-sample error** (or the training error) of a hypothesis function  $h \in \mathcal{H}$  is the empirical average of  $\{h(\mathbf{x}_n) \neq f(\mathbf{x}_n)\}$ :

$$E_{\text{in}}(h) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{n=1}^N \llbracket h(\mathbf{x}_n) \neq f(\mathbf{x}_n) \rrbracket, \quad (1)$$

where  $\llbracket \cdot \rrbracket = 1$  if the statement inside the bracket is true, and  $= 0$  if the statement is false.

## Out-Sample Error

- Let  $\mathbf{x}$  be a *testing* sample drawn from  $p(\mathbf{x})$
- $h$ : Your hypothesis
- $f$ : The unknown target function
- If  $h(\mathbf{x}) = f(\mathbf{x})$ , then say testing sample  $\mathbf{x}$  is correctly classified.
- Since  $\mathbf{x} \sim p(\mathbf{x})$ , you need to compute the probability of error, called the **out-sample error**

### Definition (Out-sample Error / Testing Error)

Consider an input space  $\mathcal{X}$  containing elements  $\mathbf{x}$  drawn from a distribution  $p_{\mathcal{X}}(\mathbf{x})$ , and a target function  $f$ . The **out-sample error** (or the testing error) of a hypothesis function  $h \in \mathcal{H}$  is

$$E_{\text{out}}(h) \stackrel{\text{def}}{=} \mathbb{P}[h(\mathbf{x}) \neq f(\mathbf{x})], \quad (2)$$

where  $\mathbb{P}[\cdot]$  measures the probability of the statement based on the distribution  $p_{\mathcal{X}}(\mathbf{x})$ .

# In-sample VS Out-sample

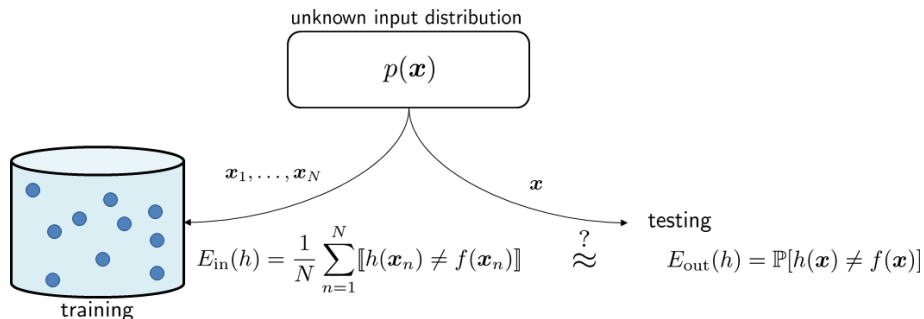
## In-Sample Error

$$E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}[h(\mathbf{x}_n) \neq f(\mathbf{x}_n)]$$

## Out-Sample Error

$$\begin{aligned} E_{\text{out}}(h) &= \mathbb{P}[h(\mathbf{x}) \neq f(\mathbf{x})] \\ &= \underbrace{\mathbb{I}[h(\mathbf{x}_n) \neq f(\mathbf{x}_n)]}_{=1} \mathbb{P}\{h(\mathbf{x}_n) \neq f(\mathbf{x}_n)\} \\ &\quad + \underbrace{\mathbb{I}[h(\mathbf{x}_n) = f(\mathbf{x}_n)]}_{=0} \left(1 - \mathbb{P}\{h(\mathbf{x}_n) \neq f(\mathbf{x}_n)\}\right) \\ &= \mathbb{E}\left\{\mathbb{I}[h(\mathbf{x}_n) \neq f(\mathbf{x}_n)]\right\} \end{aligned}$$

# The Role of $p(\mathbf{x})$



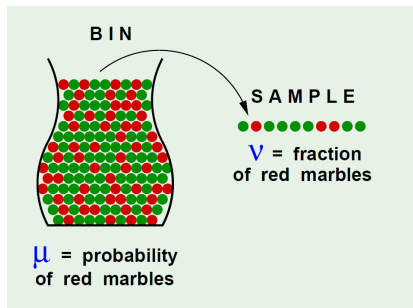
- Learning is feasible if  $\mathbf{x} \sim p(\mathbf{x})$
- $p(\mathbf{x})$  says: Training and testing are related
- If training and testing are unrelated, then hopeless – the deterministic example shown previously
- If you draw training and testing samples with different bias, then you will suffer

## When Will $E_{\text{in}} = E_{\text{out}}$ ?

### Theorem (Hoeffding Inequality)

Let  $X_1, \dots, X_N$  be a sequence of i.i.d. random variables such that  $0 \leq X_n \leq 1$  and  $\mathbb{E}[X_n] = \mu$ . Then, for any  $\epsilon > 0$ ,

$$\mathbb{P} \left[ \left| \frac{1}{N} \sum_{n=1}^N X_n - \mu \right| > \epsilon \right] \leq 2e^{-2\epsilon^2 N}. \quad (3)$$

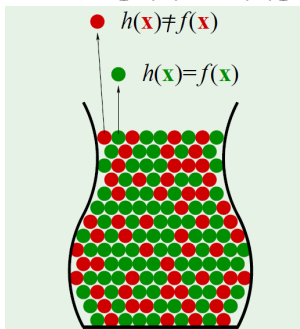


## When Will $E_{\text{in}} = E_{\text{out}}$ ?

- To us, the inequality can be stated as

$$\mathbb{P}[|E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}.$$

- $N$  = number of training samples
- $\epsilon$  = tolerance level
- Hoeffding is applicable because  $\mathbb{I}[h(\mathbf{x}) \neq f(\mathbf{x})]$  is either 1 or 0.



# Appendix