# Machine Learning
## (機器學習)

### Lecture 3: Feasibility of Learning

Hsuan-Tien Lin (林軒田)

htlin@csie.ntu.edu.tw

Department of Computer Science
& Information Engineering

National Taiwan University
(國立台灣大學資訊工程系)

# Roadmap

**1** **When** Can Machines Learn?

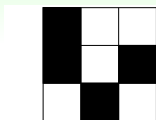*When/where learning is possible*

## Lecture 3: Feasibility of Learning

- Learning is Impossible?
- Probability to the Rescue
- Connection to Learning
- Connection to Real Learning
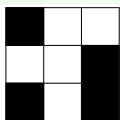- Feasibility of Learning Decomposed
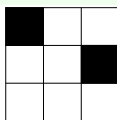
# A Learning Puzzle

Binary classification

training data

5    3    3

11



$y_n = -1$

$y_n = +1$

4    7    5

Batch supervised.

16

$g(\mathbf{x}) = ?$

my brain    $-1$

**let's test your 'human learning'
with 6 examples :-)**

## Two Controversial Answers

### whatever you say about $g(\mathbf{x})$,



$y_n = -1$

$y_n = +1$

$g(\mathbf{x}) = ?$

### truth $f(\mathbf{x}) = +1$ because . . .

### truth $f(\mathbf{x}) = -1$ because . . .

### which reason is **correct**?
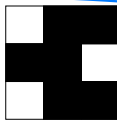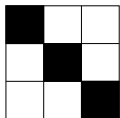
# Two Controversial Answers

**whatever you say about $g(\mathbf{x})$,**



$y_n = -1$

$y_n = +1$

$g(\mathbf{x}) = ?$

---

**truth $f(\mathbf{x}) = +1$ because . . .**

- symmetry $\Leftrightarrow$ +1 (線對稱)
- (black or white count = 3) or (black count = 4 and middle-top black) $\Leftrightarrow$ +1

**truth $f(\mathbf{x}) = -1$ because . . .**

- left-top black $\Leftrightarrow$ -1 (只看左上角)
- middle column contains at most 1 black and right-top white $\Leftrightarrow$ -1

規則隨人説

→所以需要假設

*lots of hypothesis*

all valid reasons, your **adversarial teacher** can always call you '**didn't learn**'. **:-(**

*teacher call always call you wrong.*

# A Brain-Storming Problem

$$(5, 3, 2) \rightarrow 151022, \quad (7, 2, 5) \rightarrow \textbf{?}$$

It is like a 'learning problem' with $N = 1$, $\mathbf{x}_1 = (5, 3, 2)$, $y_1 = 151022$. Learn a hypothesis from the one example to predict on $\mathbf{x} = (7, 2, 5)$. What is your answer?

*offset and sum*

**151026**

$g(\mathbf{x}) = 151012 + x_1 + x_2 + x_3$

$7 + 2 + 5 + 151012 = \underline{151026}$ ✦

*correct answer*

**143547**

$$g(\mathbf{x}) = x_1 \cdot x_2 \cdot 10000$$
$$+ \ x_1 \cdot x_3 \cdot 100$$
$$+ \ (x_1 \cdot x_2 + x_1 \cdot x_3 - x_2)$$

which one is the **smarter** answer that only top 2% people can think of?

# What is the Next Number?

1,4,1,5

# What is the Next Number?

## 1,4,1,5

1,4,1,5,**0**,-1,1,6
by $y_t = y_{t-4} - y_{t-2}$

1,4,1,5,**1**,6,1,7
by $y_t = y_{t-2} + [\![t \text{ is even}]\!]$

1,4,1,5,**2**,9,3,14
by $y_t = y_{t-4} + y_{t-2}$

**any number** can be the next!

# A 'Simple' Binary Classification Problem

| $\mathbf{x}_n$ | $y_n = f(\mathbf{x}_n)$ |
|:---:|:---:|
| 0 0 0 | ○ |
| 0 0 1 | ✗ |
| 0 1 0 | ✗ |
| 0 1 1 | ○ |
| 1 0 0 | ✗ |

*3個Bit的 training set*

*其實－某也只有8種可能
這裡已經限制了5種*

- $\mathcal{X} = \{0,1\}^3$, $\mathcal{Y} = \{\circ, \times\}$, can enumerate all candidate $f$ as $\mathcal{H}$

*g在所有 training data 都答對，但 g 接近 f 嗎?*

> pick $g \in \mathcal{H}$ with all $g(\mathbf{x}_n) = y_n$ (like PLA),
> **does $g \approx f$?**

# Infeasibility of Learning

| **x** | $y$ | $g$ | $f_1$ | $f_2$ | $f_3$ | $f_4$ | $f_5$ | $f_6$ | $f_7$ | $f_8$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 0 0 | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| 0 0 1 | × | × | × | × | × | × | × | × | × | × |
| 0 1 0 | × | × | × | × | × | × | × | × | × | × |
| 0 1 1 | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ | ○ |
| 1 0 0 | × | × | × | × | × | × | × | × | × | × |
| 1 0 1 | | ? | ○ | ○ | ○ | ○ | × | × | × | × |
| 1 1 0 | | ? | ○ | ○ | × | × | ○ | ○ | × | × |
| 1 1 1 | | ? | ○ | × | ○ | × | ○ | × | ○ | × |

*$\mathcal{D}$* — *training*, *testing*

*f 有 8 种可能，但任选一种都不代表学到了什么*

- $g \approx f$ inside $\mathcal{D}$: sure! ← zero loss → No free lunch
- $g \approx f$ outside $\mathcal{D}$: **No!** (but that's really what we want!)
  ↳ you can be totally wrong.

learning from $\mathcal{D}$ (to infer something outside $\mathcal{D}$)
is doomed if **any 'unknown' $f$ can happen**. :(
↳ 需要假設

# No Free Lunch Theorem for Machine Learning

*Without any assumptions on the learning problem on hand,*
           *all learning algorithms perform the same.*

if we allow everything.
to happen

所以用ML predict
random number 不可行!?



(CC-BY-SA 2.0 by Gaspar Torriero on Flickr)

It IS infeasible for ML to learn.

## no algorithm is best
for all learning problems

**Questions?**

# Inferring Something Unknown with Assumptions

difficult to infer **unknown target $f$ outside** $\mathcal{D}$ in learning;
can we infer **something unknown** in **other scenarios**?

統計学.

bin model

- consider a <u>bin</u> of many many <u>orange</u> and green marbles
- do we **know** the orange portion (probability)? **No!**

can you **infer** the orange probability?

# Statistics 101: Inferring **Orange** Probability



**sample** 抽樣

Assumption: Sample independently (uniform)

3 orange $\Rightarrow 03$

**bin**

| **bin** | **sample** |
|---|---|
| **assume** | **assume** $N$ marbles sampled independently: |
| orange probability = $\mu$, | orange fraction = $\nu$, |
| green probability = $1 - \mu$, | green fraction = $1 - \nu$, |
| with $\mu$ **unknown** ?? | now $\nu$ **known** |

does **in-sample** $\nu$ say anything about
out-of-sample $\mu$? 用 sample 能回推 bin 嗎？

# Possible versus Probable

does **in-sample** $\nu$ say anything about out-of-sample $\mu$?

## No!

possibly not: sample can be mostly
green while bin is mostly orange

## Yes!

probably yes: in-sample $\nu$ likely **close
to** unknown $\mu$

**sample**

**bin**

formally, **what does $\nu$ say about $\mu$?**

# Hoeffding's Inequality (1/2)

**sample of size $N$**

$\mu =$ orange
probability in bin

**bin**

$\nu =$ orange
fraction in sample

- in big sample ($N$ large), $\nu$ is probably close to $\mu$ (within $\epsilon$)

抽樣与母体的差距     ←抽樣的次數

$$\mathbb{P}\left[|\underline{\nu - \mu}| > \boxed{\epsilon}\right] \leq 2\exp\left(-2\epsilon^2\boxed{N}\right)$$

偏差的機率小於這個值

- called **Hoeffding's Inequality** for marbles, coin, polling, . . .

the statement '$\nu = \mu$' is    大概是對

**probably approximately correct** (PAC)

## Hoeffding's Inequality (2/2)

$$\mathbb{P}\left[|\nu - \mu| > \epsilon\right] \le 2 \exp\left(-2\epsilon^2 N\right)$$

- valid for all $N$ and $\epsilon$
- does not depend on $\mu$,
  **no need to 'know'** $\mu$
- larger sample size $N$ or
  looser gap $\epsilon$
  $\implies$ higher probability for '$\nu \approx \mu$'

**sample of size $N$**



**bin**

if **large $N$**, can **probably** infer
unknown $\mu$ by known $\nu$
(under iid sampling assumption)

$$2 exp \left( -2(0.01) \times 10) \right)$$

## Questions?

# Connection to Learning

i.i.d = Independent and identically distributed

## bin

- unknown orange prob. $\mu$
- marble • ∈ bin
- orange •
- green •
- size-$N$ sample from bin

  of i.i.d. marbles

## learning

- fixed hypothesis $h(\mathbf{x}) \overset{?}{=}$ target $f(\mathbf{x})$

  marble
- $\boxed{\mathbf{x}} \in \boxed{\mathcal{X}}$   bin

  一樣
- $h$ is wrong $\Leftrightarrow \underline{h(\mathbf{x}) \neq f(\mathbf{x})}$ : 漆成 orange
- $h$ is right $\Leftrightarrow \underline{h(\mathbf{x}) = f(\mathbf{x})}$ : 漆成 green

  不一樣
- check $h$ on $\mathcal{D}_l = \{(\mathbf{x}_n, y_n)\}$

  檢查 h 在 D 上的表現   $f(\mathbf{x}_n)$ f(x) 的結果

  with i.i.d. $\mathbf{x}_n$

$X_n$ 是獨立平均的從 bin 裡抽出來

if **large $N$** & $\boxed{\text{i.i.d. } \mathbf{x}_l}$, can **probably** infer
unknown $[\![h(\mathbf{x}) \neq f(\mathbf{x})]\!]$ probability
by known $[\![h(\mathbf{x}_n) \neq y_n]\!]$ fraction

- $h(\mathbf{x}) \neq f(\mathbf{x})$
- $h(\mathbf{x}) = f(\mathbf{x})$

## Added Components

verify 確認是否相等



| unknown target function $f: \mathcal{X} \to \mathcal{Y}$ | | unknown $P$ on $\mathcal{X}$ ↓用P来抽樣 | | $h \overset{?}{\approx} f$ |
| --- | --- | --- | --- | --- |

*(ideal credit approval formula)*

training data
$\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N$

training examples
$\mathcal{D}: (\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_N, y_N)$

*(historical records in bank)*

learning algorithm
$\mathcal{A}$

final hypothesis
$g \approx f$

*('learned' formula to be used)*

hypothesis set
$\mathcal{H}$

*(set of candidate formula)*

fixed $h$

開上帝視角去 bin 裡抽彈珠,抽到 orange 的機率

for any fixed $h$, can probably infer

error rate    Exception    testing

**unknown** $E_{\text{out}}(h) = \underset{\mathbf{x} \sim P}{\mathcal{E}} [\![h(\mathbf{x}) \neq f(\mathbf{x})]\!]$  ← booling op.

by **known** $E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^{N} [\![h(\mathbf{x}_n) \neq y_n]\!]$  ← training.

(under iid sampling assumption)  *從已經抽樣的 training set 裡,抽到 orange 的機率

# The Formal Guarantee

for any fixed *h*, in 'big' data *(N large)*,

in-sample error $E_{in}(h)$ is probably close to

out-of-sample error $E_{out}(h)$ (within $\epsilon$)

# of sample

$$\mathbb{P}\left[\left|E_{in}(h) - E_{out}(h)\right| > \epsilon\right] \leq 2\exp\left(-2\epsilon^2 \underline{N}\right)$$

抽樣內的 error 與抽樣外的 error

## same as the 'bin' analogy . . . 的變題

- valid for all *N* and $\epsilon$
- does not depend on $E_{out}(h)$, **no need to 'know'** $E_{out}(h)$
  —*f* and *P* can stay unknown
- '$E_{in}(h) = E_{out}(h)$' is **probably approximately correct (PAC)**

N 如果夠大就大致上一樣。

if '$E_{in}(h) \approx E_{out}(h)$' and '$E_{in}(h)$ **small**'
$\implies E_{out}(h)$ small $\implies h \approx f$ with respect to *P*

achieve learning.

# Verification of One *h*

for any fixed *h*, when data large enough,

$$E_{in}(h) \approx E_{out}(\underline{h})$$

*g* 是 ML 選的 function.

**Can we claim 'good learning' (*g* ≈ *f*)?** 想辦法選一個接近 *g*.

### Yes!

if $E_{in}(h)$ **small for the fixed** *h*
  and $\mathcal{A}$ **pick the** *h* **as** *g*
$\Longrightarrow$ '*g* = *f*' PAC

### No!

if $\mathcal{A}$ **forced to pick THE** *h* **as** *g*
$\Longrightarrow$ $E_{in}(h)$ **almost always not small**
$\Longrightarrow$ '*g* ≠ *f*' PAC

用來判斷 *g*, *f* 有沒有接近.

real learning:
  $\mathcal{A}$ shall **make choices** $\in \mathcal{H}$ (like PLA)
  rather than **being forced to pick one** *h*. :-(

這并不是 learning. 因為并不會做選擇

# The 'Verification' Flow (testing)



unknown target function
$f: \mathcal{X} \to \mathcal{Y}$

*(ideal credit approval formula)*

unknown
$P$ on $\mathcal{X}$

$\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N$

$\mathbf{x}$

**verifying** examples
$\mathcal{D}: (\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_N, y_N)$

*(historical records in bank)*

final hypothesis
$g \approx f$

**(given formula to be verified)**

one
**hypothesis**
$h$

**(one** *candidate formula)*

$g = h$

檢查看這個 hypothesis 是 好不好

can now use 'historical records' (data) to
**verify 'one candidate formula'** $h$

# Questions?

# Multiple $h$

green → correct
orange → Incorrect

one hypothesis for one bin



$h_1$      $h_2$      $h_M$

$E_{out}(h_1)$    $E_{out}(h_2)$    $E_{out}(h_M)$

. . . . . . . .

unlucky sample!

$E_{in}(h_1)$    $E_{in}(h_2)$ 选green多的 $h$ →    $E_{in}(h_M)$

But will $E_{out}$ be all gree?

real learning (say like PLA):
**BINGO** when getting ●●●●●●●●●?

# Coin Game

*BAD 就是 $E_{in}$, $E_{out}$ 不一樣*



Q: if everyone in size-400 NTU ML class flips a coin 5 times, and **one of the students gets 5 heads for her coin 'g'**. Is 'g' really magical?

*高機率不是作弊 coin*

A: No. Even if all coins are fair, the probability that **one of the coins** results in **5 heads** is $1 - \left(\frac{31}{32}\right)^{400} >$ 99%. *任一個學生丟出連續5個 head 機率其實很大*

**BAD sample:** *$E_{in}$* and *$E_{out}$* **far away**
**—can get worse when involving 'choice'**

## BAD Sample and BAD Data

### BAD Sample

e.g., $E_{out} = \frac{1}{2}$, but getting all heads ($E_{in} = 0$)! →BAD    Eout 和 Ein (差很遠)

丟 coin

### BAD Data for One $h$

$E_{out}(h)$ **and** $E_{in}(h)$ **far away**:
e.g., $E_{out}$ big (far from $f$), but $E_{in}$ small (correct on most examples)

不會有太多悟是不好的

|   | $\mathcal{D}_1$ | $\mathcal{D}_2$ | ... | $\mathcal{D}_{1126}$ | ... | $\mathcal{D}_{5678}$ | ... | Hoeffding |
|---|---|---|---|---|---|---|---|---|
| $h$ | **BAD** |  |  |  |  | **BAD** |  | $\mathbb{P}_{\mathcal{D}}$ [**BAD** $\mathcal{D}$ for $h$] $\leq$ ... |

*Hoeffding 保證 BAD data 發生機率小

Hoeffding: small

$$\left( \mathbb{P}_{\mathcal{D}} \left[ \textbf{BAD } \mathcal{D} \right] = \sum_{\text{all possible} \mathcal{D}} \mathbb{P}(\mathcal{D}) \cdot \llbracket \textbf{BAD } \mathcal{D} \rrbracket \right)$$

# BAD Data for Many $h$

**GOOD** data for many $h$ ⎤ *Good data can verify every hypothesis*
⟺ **GOOD** data for verifying any $h$ ⎦
⟺ there exists **no BAD** $h$ such that $E_{out}(h)$ and $E_{in}(h)$ far away
**there exists some** $h$ **such that** $E_{out}(h)$ and $E_{in}(h)$ far away
⟺ **BAD** data for many $h$

*Super dataset*

*In real world - we only have one data set*

| | $\mathcal{D}_1$ | $\mathcal{D}_2$ | ... | $\mathcal{D}_{1126}$ | ... | $\mathcal{D}_{5678}$ | Hoeffding |
|---|---|---|---|---|---|---|---|
| $h_1$ | **BAD** | | | | | **BAD** | $\mathbb{P}_{\mathcal{D}}\left[\textbf{BAD }\mathcal{D}\text{ for }h_1\right] \leq \ldots$ |
| $h_2$ | | **BAD** | | | | | $\mathbb{P}_{\mathcal{D}}\left[\textbf{BAD }\mathcal{D}\text{ for }h_2\right] \leq \ldots$ |
| $h_3$ | **BAD** | **BAD** | | | | **BAD** | $\mathbb{P}_{\mathcal{D}}\left[\textbf{BAD }\mathcal{D}\text{ for }h_3\right] \leq \ldots$ |
| ... | | | | | | | |
| $h_M$ | **BAD** | | | | | **BAD** | $\mathbb{P}_{\mathcal{D}}\left[\textbf{BAD }\mathcal{D}\text{ for }h_M\right] \leq \ldots$ |
| all | **BAD** | **BAD** | | **GOOD** | | **BAD** | **?** |

*BAD for some hypothesis* | *对所有的 dataset 都 GOOD*

do *not* know if $\mathcal{D}$ is **BAD** or not;
wish $\mathbb{P}_{\mathcal{D}}[\textbf{BAD }\mathcal{D}]$ small & pray for **"GOOD luck"**

# Bound of BAD Data

到 dataset 只要對任一個 hypothesis BAD

$$\mathbb{P}_{\mathcal{D}}[\textbf{BAD } \mathcal{D}]$$

對任一個 hypothesis 不好的機率

$$= \mathbb{P}_{\mathcal{D}}[\textbf{BAD } \mathcal{D} \text{ for } h_1 \text{ or } \textbf{BAD } \mathcal{D} \text{ for } h_2 \text{ or } \ldots \text{ or } \textbf{BAD } \mathcal{D} \text{ for } h_M]$$

$$\leq \mathbb{P}_{\mathcal{D}}[\textbf{BAD } \mathcal{D} \text{ for } h_1] + \mathbb{P}_{\mathcal{D}}[\textbf{BAD } \mathcal{D} \text{ for } h_2] + \ldots + \mathbb{P}_{\mathcal{D}}[\textbf{BAD } \mathcal{D} \text{ for } h_M]$$

假設額不交集 (union bound) : not overlapping

$$\leq 2 \exp\left(-2\epsilon^2 N\right) + 2 \exp\left(-2\epsilon^2 N\right) + \ldots + 2 \exp\left(-2\epsilon^2 N\right)$$

$$= 2M \exp\left(-2\epsilon^2 N\right) \leftarrow \text{very loose upperbound.}$$

↳ # of hypothesis

- finite-bin version of Hoeffding, valid for all $M$, $N$ and $\epsilon$
- does not depend on any $E_{\text{out}}(h_m)$, **no need to 'know'** $E_{\text{out}}(h_m)$ —$f$ and $P$ can stay unknown

$E_{\text{in}}(h) \simeq E_{\text{out}}(h)$

- '$E_{\text{in}}(g) = E_{\text{out}}(g)$' is **PAC**, **regardless of** $\mathcal{A}$

'most reasonable' $\mathcal{A}$ (like PLA): 選最小 $E_{in}$ 的 h 作為 g.
pick the $h_m$ with **lowest $E_{\text{in}}(h_m)$** as $g$
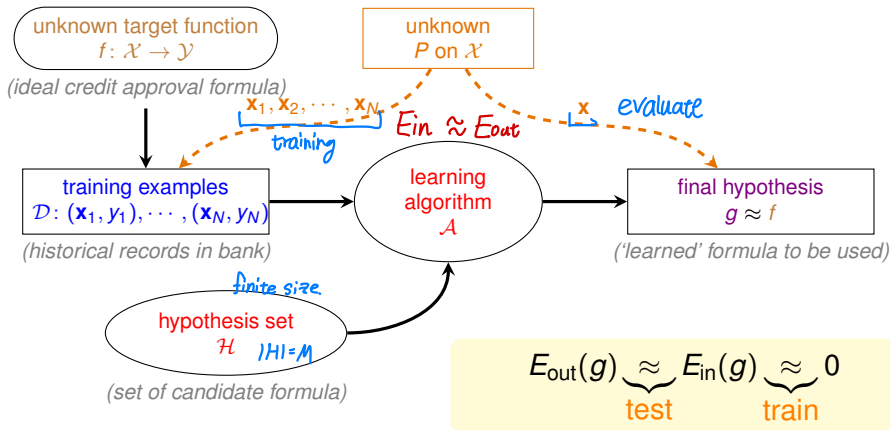
**Questions?**

# The 'Statistical' Learning Flow

*# of hypothesis*    *# of data*

if $|\mathcal{H}| = M$ finite, $N$ large enough,

for whatever $g$ picked by $\mathcal{A}$, $E_{out}(g) \approx E_{in}(g)$

if $\mathcal{A}$ finds one $g$ with $E_{in}(g) \approx 0$,

(PAC) guarantee for $E_{out}(g) \approx 0 \implies$ **learning possible :-)**

unknown target function
$f: \mathcal{X} \to \mathcal{Y}$

*(ideal credit approval formula)*

unknown
$P$ on $\mathcal{X}$

$\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N$

*training*

$E_{in} \approx E_{out}$

$\mathbf{x}$    *evaluate*

training examples
$\mathcal{D}: (\mathbf{x}_1, y_1), \cdots, (\mathbf{x}_N, y_N)$

*(historical records in bank)*

learning
algorithm
$\mathcal{A}$

final hypothesis
$g \approx f$

*('learned' formula to be used)*

*finite size*

hypothesis set
$\mathcal{H}$    $|H| = M$

*(set of candidate formula)*

$$E_{out}(g) \underbrace{\approx}_{\text{test}} E_{in}(g) \underbrace{\approx}_{\text{train}} 0$$

# Two Central Questions

for batch & supervised binary classification, $g \approx f \iff E_{\text{out}}(g) \approx 0$.

lecture 2                                       lecture 1

achieved through $E_{\text{out}}(g) \approx E_{\text{in}}(g)$  and  $E_{\text{in}}(g) \approx 0$ PLA

lecture 3                          lecture 1

*if $M$ is finite*

learning split to two central questions:

1. can we make sure that $E_{\text{out}}(g)$ is close enough to $E_{\text{in}}(g)$? (test/generalize) *Can we extend from training → testing data*

2. can we make $E_{\text{in}}(g)$ small enough? (train/optimize) *↳ Learn well on training set!*

what role does $M$ play for the two questions?

$|\mathcal{H}|$

# Trade-off on *M*

1. can we make sure that $E_{out}(g)$ is close enough to $E_{in}(g)$?
2. can we make $E_{in}(g)$ small enough?

**small *M***

1. Yes!, 沒有太多 hypothesis 跟你 screwed over
   $\mathbb{P}[\textbf{BAD}] \leq 2 \cdot M \cdot \exp(\ldots)$
2. No!, too few choices

但很有可能學不起來

**large *M***    ← # of hypothesis

有可能迷到 BAD

1. No!,
   $\mathbb{P}[\textbf{BAD}] \leq 2 \cdot M \cdot \exp(\ldots)$
2. Yes!, many choices

using the right *M* (or $\mathcal{H}$) is important
$M = \infty$ **doomed?**

M很重要,適中比較好

# Preview

### Known

$$\mathbb{P}\left[\left|E_{\text{in}}(g) - E_{\text{out}}(g)\right| > \epsilon\right] \leq 2 \cdot \underline{M} \cdot \exp\left(-2\epsilon^2 N\right)$$

把 M 换成一个有限的数量: $m_{\mathcal{H}}$

### Todo

- establish **a finite quantity** that replaces $M$

  a finite quantity

$$\mathbb{P}\left[\left|E_{\text{in}}(g) - E_{\text{out}}(g)\right| > \epsilon\right] \stackrel{?}{\leq} 2 \cdot \boxed{m_{\mathcal{H}}} \cdot \exp\left(-2\epsilon^2 N\right)$$

- justify the feasibility of learning for infinite $M$
- study $m_{\mathcal{H}}$ to understand its trade-off for 'right' $\mathcal{H}$, just like $M$

> mysterious PLA to be fully resolved
> **"soon" :-)**

# Questions?

# Summary

**1** **When** Can Machines Learn?

> ## Lecture 2: The Learning Problems
>
> ## Lecture 3: Feasibility of Learning
>
> - Learning is Impossible?
>      **absolutely no free lunch outside** $\mathcal{D}$
> - Probability to the Rescue
>      **probably approximately correct outside** $\mathcal{D}$
> - Connection to Learning
>   **verification possible if** $E_{in}(h)$ **small for fixed** $h$
> - Connection to Real Learning
>   **learning possible if** $|\mathcal{H}|$ **finite and** $E_{in}(g)$ **small**
> - Feasibility of Learning Decomposed
>   **two questions:** $E_{out}(g) \approx E_{in}(g)$, **and** $E_{in}(g) \approx 0$

**2** Why Can Machines Learn?

- **next: what if** $|\mathcal{H}| = \infty$**?**