

Machine Learning Foundations (機器學習基石)



Lecture 4: 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?

Lecture 3: Types of Learning

focus: binary classification or regression from a batch of supervised data with concrete features

Lecture 4: Feasibility of Learning

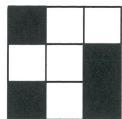
- Learning is Impossible?
- Probability to the Rescue
- Connection to Learning
- Connection to Real Learning

② Why Can Machines Learn?

③ How Can Machines Learn?

④ How Can Machines Learn Better?

A Learning Puzzle



$$y_n = -1$$



$$y_n = +1$$

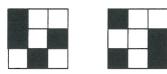


$$g(\mathbf{x}) = ?$$

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 ...

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

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

- left-top black $\Leftrightarrow -1$
- middle column contains at
most 1 black and right-top
white $\Leftrightarrow -1$

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

A 'Simple' Binary Classification Problem

\mathbf{x}_n	$y_n = f(\mathbf{x}_n)$
0 0 0	o
0 0 1	x
0 1 0	x
0 1 1	o
1 0 0	x

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

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

No Free Lunch

\mathbf{x}	y	g	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8
0 0 0	o	o	o	o	o	o	o	o	o	o
0 0 1	x	x	x	x	x	x	x	x	x	x
0 1 0	x	x	x	x	x	x	x	x	x	x
0 1 1	o	o	o	o	o	o	o	o	o	o
1 0 0	x	x	x	x	x	x	x	x	x	x
1 0 1		?	o	o	o	o	x	x	x	x
1 1 0		?	o	o	x	x	o	o	x	x
1 1 1		?	o	x	o	x	o	x	o	x

 \mathcal{D}

已知
未知

- $g \approx f$ inside \mathcal{D} : sure!
- $g \approx f$ outside \mathcal{D} : No! (but that's really what we want!)

learning from \mathcal{D} (to infer something outside \mathcal{D})
is doomed if any 'unknown' can happen. :-(

Fun Time

This is a popular ‘brain-storming’ problem, with a claim that 2% of the world’s cleverest population can crack its ‘hidden pattern’.

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

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?

- | | |
|----------|--|
| ① 151026 | ③ I need more examples to get the correct answer |
| ② 143547 | ④ there is no ‘correct’ answer |

4. 因在未知狀況
是無法確定的正確

Fun Time

This is a popular ‘brain-storming’ problem, with a claim that 2% of the world’s cleverest population can crack its ‘hidden pattern’.

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

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?

- | | |
|----------|--|
| ① 151026 | ③ I need more examples to get the correct answer |
| ② 143547 | ④ there is no ‘correct’ answer |

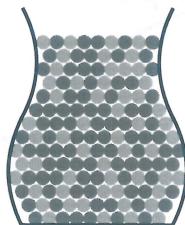
Reference Answer: ④

Following the same nature of the no-free-lunch problems discussed, we cannot hope to be correct under this ‘adversarial’ setting. BTW, ② is the designer’s answer: the first two digits = $x_1 \cdot x_2$; the next two digits = $x_1 \cdot x_3$; the last two digits = $(x_1 \cdot x_2 + x_1 \cdot x_3 - x_2)$.

未知想推論

Inferring Something Unknown

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

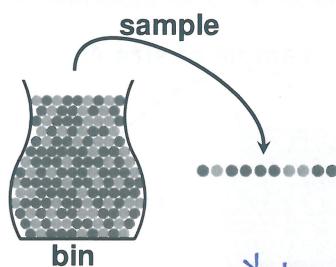


- consider a bin of many many orange and green marbles
- do we know the orange portion (probability)? No!

can you infer the orange probability?

抽樣

Statistics 101: Inferring Orange Probability



bin

assume

orange probability = μ ,
green probability = $1 - \mu$,
with μ unknown

sample 標本

N marbles sampled independently, with
orange fraction = ν ,
green fraction = $1 - \nu$,
now ν known

does in-sample ν say anything about
out-of-sample μ ?

試驗 試驗

Possible versus Probable

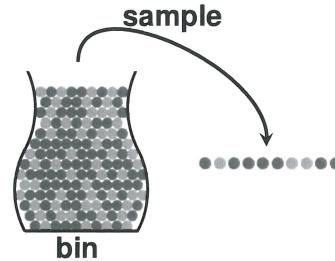
does in-sample ν say anything about out-of-sample μ ?

No!

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

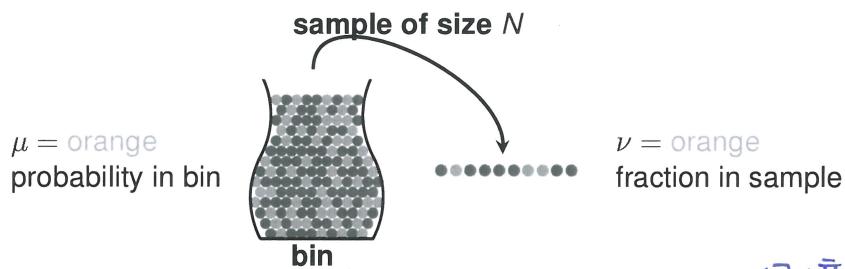
Yes!

probably yes: in-sample ν likely close to unknown μ



formally, what does ν say about μ ?

Hoeffding's Inequality (1/2)



- in big sample (N large), ν is probably close to μ (within ϵ)

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

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

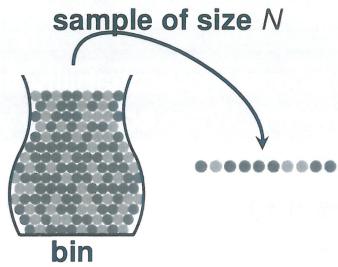
the statement ' $\nu = \mu$ ' is
probably approximately correct (PAC)

大概 差不多

Hoeffding's Inequality (2/2)

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

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



if large N , can probably infer
unknown μ by known ν

Fun Time

Let $\mu = 0.4$. Use Hoeffding's Inequality

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

to bound the probability that a sample of 10 marbles will have $\nu \leq 0.1$. What bound do you get?

- ① 0.67
- ② 0.40
- ③ 0.33
- ④ 0.05

Fun Time

Let $\mu = 0.4$. Use Hoeffding's Inequality

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

to bound the probability that a sample of 10 marbles will have $\nu \leq 0.1$. What bound do you get?

- ① 0.67
- ② 0.40
- ③ 0.33
- ④ 0.05

Reference Answer: ③

Set $N = 10$ and $\epsilon = 0.3$ and you get the answer. BTW, ④ is the actual probability and Hoeffding gives only an upper bound to that.

Connection to Learning

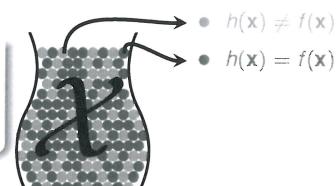
bin

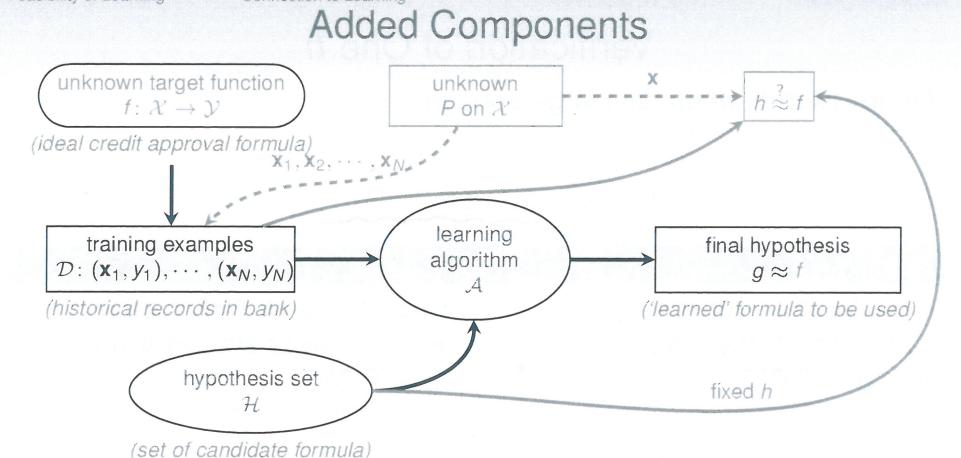
- unknown orange prob. μ
- marble $\bullet \in \text{bin}$
- orange \bullet
- green \bullet
- size- N sample from bin
of i.i.d. marbles

learning

- fixed hypothesis $h(\mathbf{x}) \stackrel{?}{=} \text{target } f(\mathbf{x})$
- $\mathbf{x} \in \mathcal{X}$
- h is wrong $\Leftrightarrow h(\mathbf{x}) \neq f(\mathbf{x})$
- h is right $\Leftrightarrow h(\mathbf{x}) = f(\mathbf{x})$
- check h on $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}$
獨立抽樣
with i.i.d. \mathbf{x}_n

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





for any fixed h , can probably infer

$$\text{unknown } E_{\text{out}}(h) = \mathbb{E}_{x \sim P} [h(x) \neq f(x)]$$

by known $E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N [h(x_n) \neq y_n]$.

母体

標本

The Formal Guarantee

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

$E_{\text{in}}(h)$ is probably close to
 $E_{\text{out}}(h)$ (within ϵ)

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

same as the 'bin' analogy ...

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

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

Verification of One h

for any fixed h , when data large enough,

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

Can we claim 'good learning' ($g \approx f$)?

Yes!

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

No!

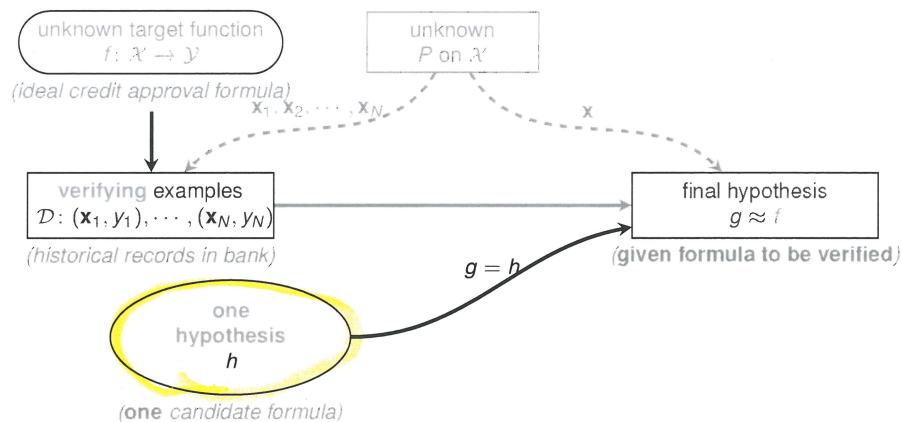
if \mathcal{A} forced to pick THE h as g
 $\Rightarrow E_{\text{in}}(h)$ almost always not small
 $\Rightarrow 'g \neq f'$ PAC!

real learning:

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

預測表現好不好

The 'Verification' Flow



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

Fun Time

Your friend tells you her secret rule in investing in a particular stock: 'Whenever the stock goes down in the morning, it will go up in the afternoon; vice versa.' **To verify the rule, you chose 100 days uniformly at random from the past 10 years of stock data, and found that 80 of them satisfy the rule.** What is the best guarantee that you can get from the verification?

- ① You'll definitely be rich by exploiting the rule in the next 100 days.
- ② You'll likely be rich by exploiting the rule in the next 100 days, if the market behaves similarly to the last 10 years.
- ③ You'll likely be rich by exploiting the 'best rule' from 20 more friends in the next 100 days.
- ④ You'd definitely have been rich if you had exploited the rule in the past 10 years.

Fun Time

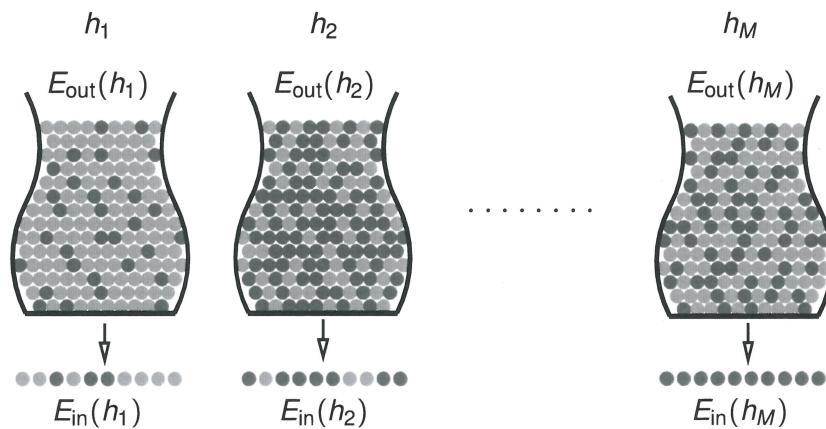
Your friend tells you her secret rule in investing in a particular stock: 'Whenever the stock goes down in the morning, it will go up in the afternoon; vice versa.' **To verify the rule, you chose 100 days uniformly at random from the past 10 years of stock data, and found that 80 of them satisfy the rule.** What is the best guarantee that you can get from the verification?

- ① You'll definitely be rich by exploiting the rule in the next 100 days.
- ② You'll likely be rich by exploiting the rule in the next 100 days, if the market behaves similarly to the last 10 years.
- ③ You'll likely be rich by exploiting the 'best rule' from 20 more friends in the next 100 days.
- ④ You'd definitely have been rich if you had exploited the rule in the past 10 years.

Reference Answer: (2)

(1): no free lunch; (3): no 'learning' guarantee in verification; (4): verifying with only 100 days, possible that the rule is mostly wrong for whole 10 years.

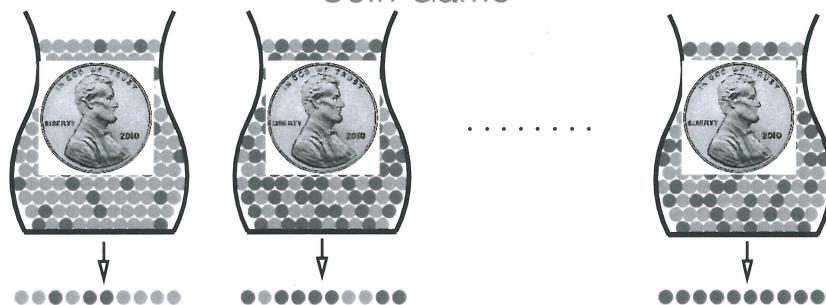
Multiple h



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

全对

Coin Game



Q: if everyone in size-150 NTU ML class flips a coin 5 times, and one of the students gets 5 heads for her coin 'g'. Is 'g' really magical?

A: No. Even if all coins are fair, the probability that one of the coins results in 5 heads is $1 - \left(\frac{31}{32}\right)^{150} > 99\%$.

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

有選擇就有多樣

BAD Sample and BAD Data

BAD Sample

e.g., $E_{\text{out}} = \frac{1}{2}$, but getting all heads ($E_{\text{in}} = 0$)!

BAD Data for One h

$E_{\text{out}}(h)$ and $E_{\text{in}}(h)$ far away:

e.g., E_{out} big (far from f), but E_{in} small (correct on most examples)

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

Hoeffding: small

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

BAD Data for Many h

BAD data for many h

\iff no 'freedom of choice' by \mathcal{A}

\iff there exists some h such that $E_{\text{out}}(h)$ and $E_{\text{in}}(h)$ far away

	\mathcal{D}_1	\mathcal{D}_2	...	\mathcal{D}_{1126}	...	\mathcal{D}_{5678}	Hoeffding
h_1	BAD					BAD	$\mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_1] \leq \dots$
h_2		BAD					$\mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_2] \leq \dots$
h_3	BAD	BAD				BAD	$\mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_3] \leq \dots$
...							
h_M	BAD					BAD	$\mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_M] \leq \dots$
all	BAD	BAD				BAD	?

for M hypotheses, bound of $\mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D}]$?

Bound of BAD Data

$$\begin{aligned}
 & \mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D}] \\
 &= \mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_1 \text{ or } \text{BAD } \mathcal{D} \text{ for } h_2 \text{ or } \dots \text{ or } \text{BAD } \mathcal{D} \text{ for } h_M] \\
 &\leq \mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_1] + \mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_2] + \dots + \mathbb{P}_{\mathcal{D}}[\text{BAD } \mathcal{D} \text{ for } h_M] \\
 &\quad (\text{union bound}) \\
 &\leq 2 \exp(-2\epsilon^2 N) + 2 \exp(-2\epsilon^2 N) + \dots + 2 \exp(-2\epsilon^2 N) \\
 &= 2M \exp(-2\epsilon^2 N)
 \end{aligned}$$

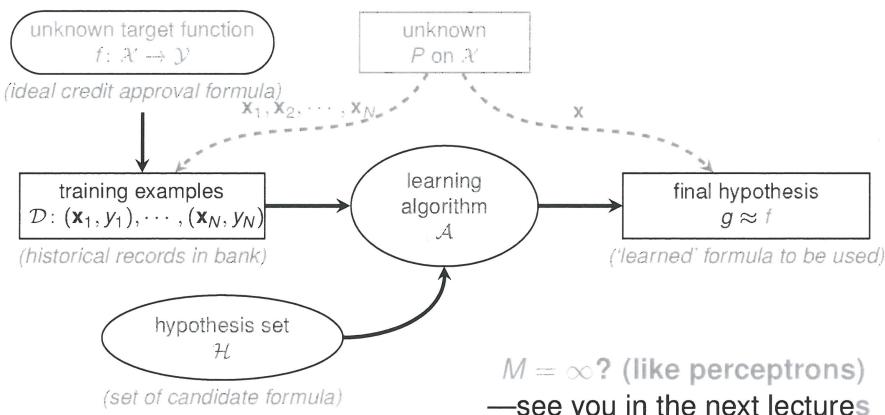
- finite-bin version of Hoeffding, valid for all M, N and ϵ
- does not depend on any $E_{\text{out}}(h_m)$, no need to ‘know’ $E_{\text{out}}(h_m)$
— f and P can stay unknown
- ‘ $E_{\text{in}}(g) = E_{\text{out}}(g)$ ’ is PAC, regardless of \mathcal{A}

‘most reasonable’ \mathcal{A} (like PLA/pocket):
pick the h_m with lowest $E_{\text{in}}(h_m)$ as g

The ‘Statistical’ Learning Flow

if $|\mathcal{H}| = M$ finite, N large enough,
for whatever g picked by \mathcal{A} , $E_{\text{out}}(g) \approx E_{\text{in}}(g)$

if \mathcal{A} finds one g with $E_{\text{in}}(g) \approx 0$,
PAC guarantee for $E_{\text{out}}(g) \approx 0 \implies$ learning possible :-)



Fun Time

Consider 4 hypotheses.

$$h_1(\mathbf{x}) = \text{sign}(x_1), \quad h_2(\mathbf{x}) = \text{sign}(x_2),$$

$$h_3(\mathbf{x}) = \text{sign}(-x_1), \quad h_4(\mathbf{x}) = \text{sign}(-x_2).$$

For any N and ϵ , which of the following statement is not true?

- ① the **BAD** data of h_1 and the **BAD** data of h_2 are exactly the same
- ② the **BAD** data of h_1 and the **BAD** data of h_3 are exactly the same
- ③ $\mathbb{P}_{\mathcal{D}}[\text{BAD for some } h_k] \leq 8 \exp(-2\epsilon^2 N)$
- ④ $\mathbb{P}_{\mathcal{D}}[\text{BAD for some } h_k] \leq 4 \exp(-2\epsilon^2 N)$

Fun Time

Consider 4 hypotheses.

$$h_1(\mathbf{x}) = \text{sign}(x_1), \quad h_2(\mathbf{x}) = \text{sign}(x_2),$$

$$h_3(\mathbf{x}) = \text{sign}(-x_1), \quad h_4(\mathbf{x}) = \text{sign}(-x_2).$$

For any N and ϵ , which of the following statement is not true?

- ① the **BAD** data of h_1 and the **BAD** data of h_2 are exactly the same
- ② the **BAD** data of h_1 and the **BAD** data of h_3 are exactly the same
- ③ $\mathbb{P}_{\mathcal{D}}[\text{BAD for some } h_k] \leq 8 \exp(-2\epsilon^2 N)$
- ④ $\mathbb{P}_{\mathcal{D}}[\text{BAD for some } h_k] \leq 4 \exp(-2\epsilon^2 N)$

Reference Answer: ①

The important thing is to note that ② is true, which implies that ④ is true if you revisit the union bound. Similar ideas will be used to conquer the $M = \infty$ case.

Summary

➊ When Can Machines Learn?

Lecture 3: Types of Learning

Lecture 4: 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

➋ Why Can Machines Learn?

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

➌ How Can Machines Learn?

➍ How Can Machines Learn Better?