

CIE 5133 機器學習與深度學習導論

線上課程

開始之前 (09.29.2021)

- 請將你的麥克風靜音
- 請找個安全、舒適的空間
- 聽講時有任何問題請到 [slido #073374](#) 留言
- 我們會透過 Zoom, slido, 討論區，臉書社群來強化無法面對面所造成的互動不足，同學們有任何建議也請讓我們知道。
- Stay Home! Happy Learning!!

加簽？ 旁聽？

1. 想加簽的同學：如果人數不要太離譜，我會儘可能加簽。
2. 歡迎旁聽。
3. 請寄 email 給我
(dchen@ntu.edu.tw) 或大助教 (harry@caeca.net)。

Question? Zoom 回應舉手、Zoom 聊天留言、
slido #073374 留言

CIE 5133 機器學習與深度學習導論



Course FaceBook

Today ...

- Fundamentals and Landscape of Classical Machine Learning (I) and (II)
- HW1



Fundamentals and Landscape of Classical Machine Learning (I)

- Learning? What do we mean?
- Is learning feasible?

[https://www.sli.do/
#073374](https://www.sli.do/#073374)

Learning? What do we mean?

Traditional Programming

$+$ $-$
 \times $=$

Rules



Computer



The computer can perform the task it has been *instructed* to do.

Machine Learning



Data



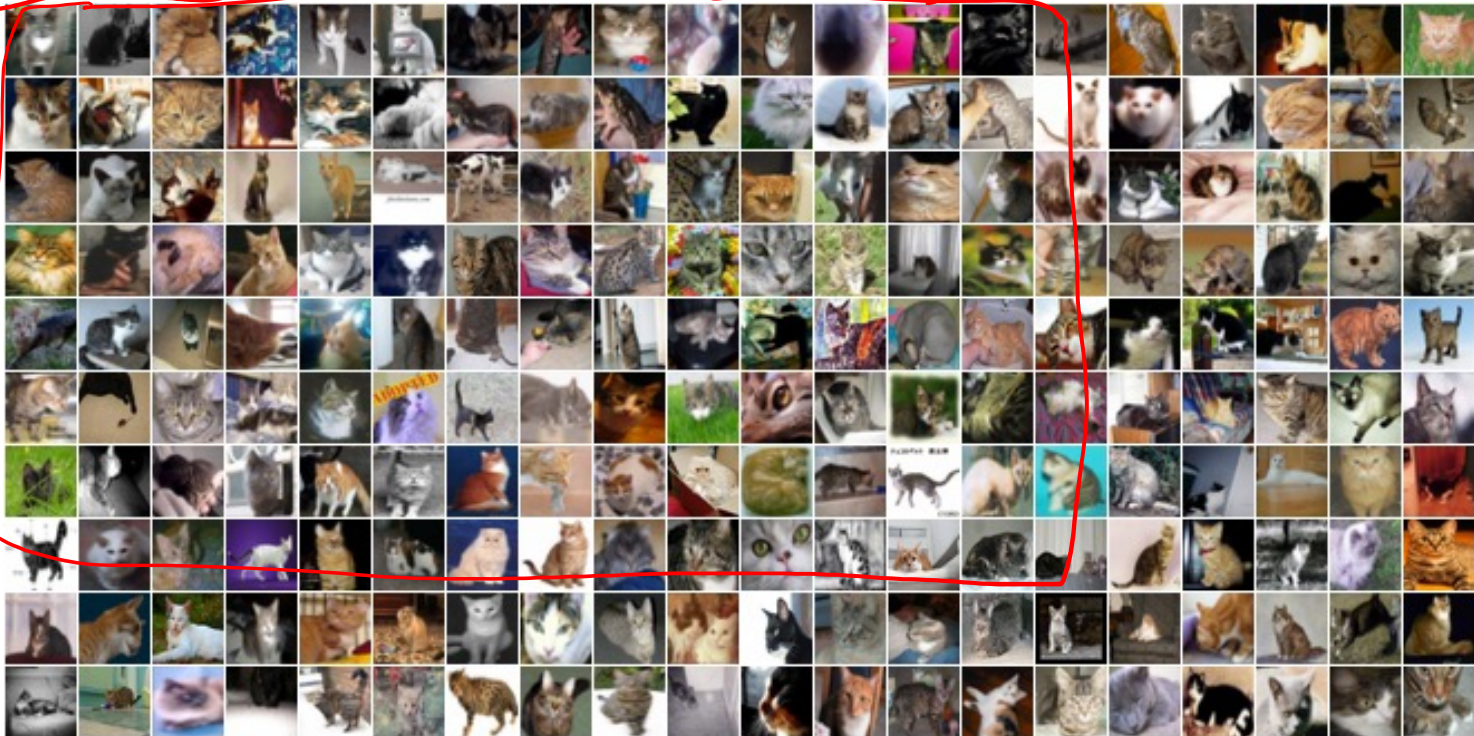
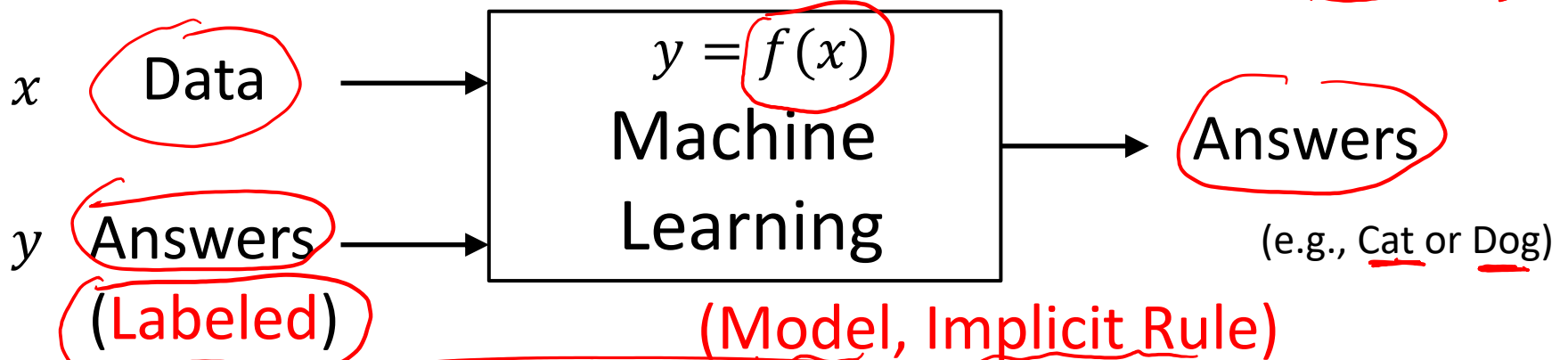
Computer



The computer can perform the task it has *learned* to do.

- Can computer automatically learn these rules from data?
- If so, can these rules be applied for new data to guess answers?

Machine Learning (Supervised Learning): rules can be implicitly



Fun time: which of the following is best suited for machine learning?

- (1) Throw a dice and predict its face value
- (2) Sort a few points in space
- (3) Decide credit card approval for a customer
- (4) Predict the next big earthquake

- (1) Random event, no pattern
- (2) Programmable task
- (3) Pattern: customer behavior

Not easily programmable

Data: history of bank operation

- (4) Arguably not enough data (yet)

[https://www.sli.do/
#073374](https://www.sli.do/#073374)

Components of Learning: Metaphor Using Credit Approval

Applicant Information

→	age	23 years
→	gender	female
→	annual salary	NTD 1,000,000
→	year in residence	1 year
→	year in job	0.5 year
→	current debt	200,000

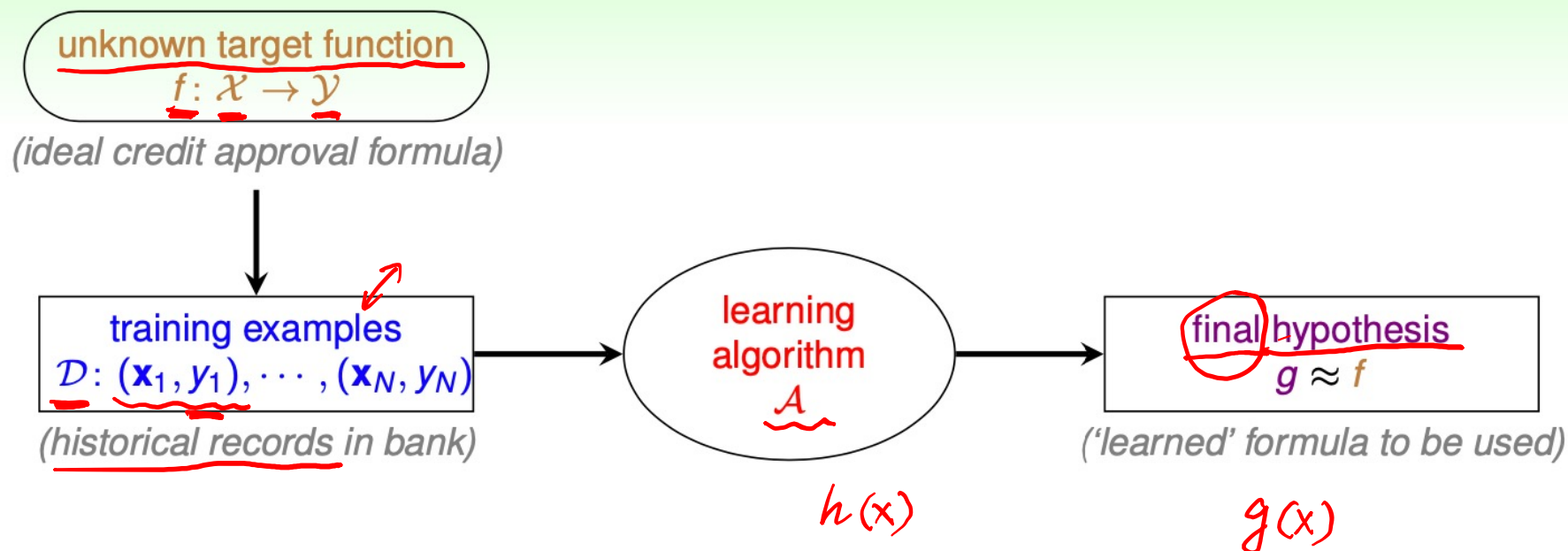
 x_1 $f(x)$

unknown pattern to be learned:

'approve credit card good for bank?'

 y f

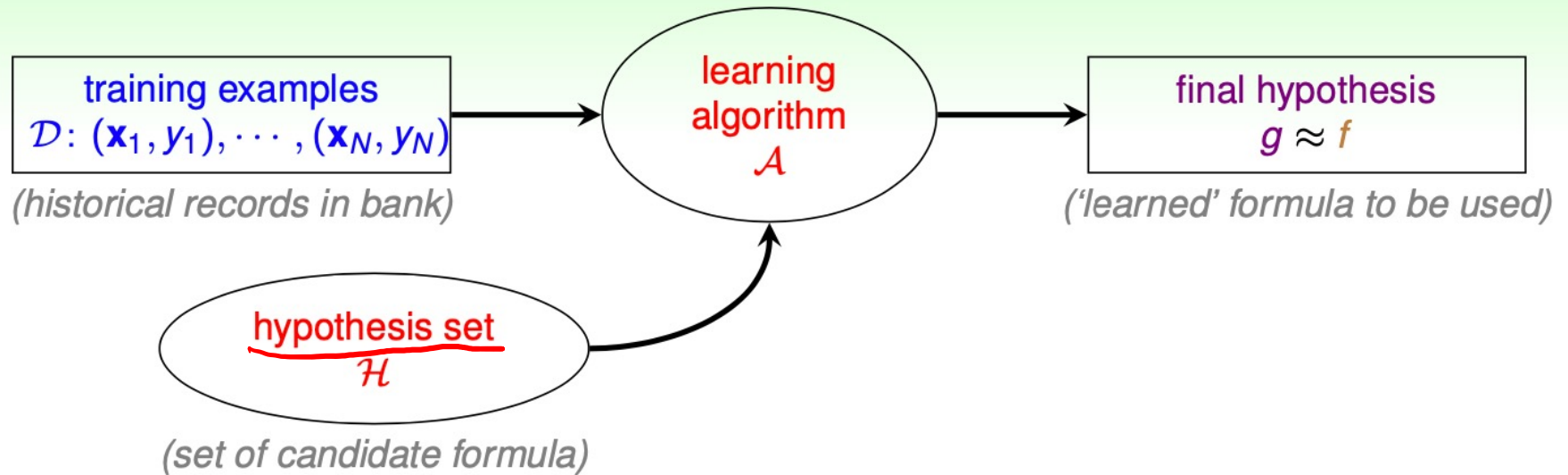
Learning Flow for Credit Approval



- target f **unknown**
 (i.e. no programmable definition)
- hypothesis g hopefully $\approx f$ but possibly **different** from f
 (perfection 'impossible' when f unknown)

What does g look like?

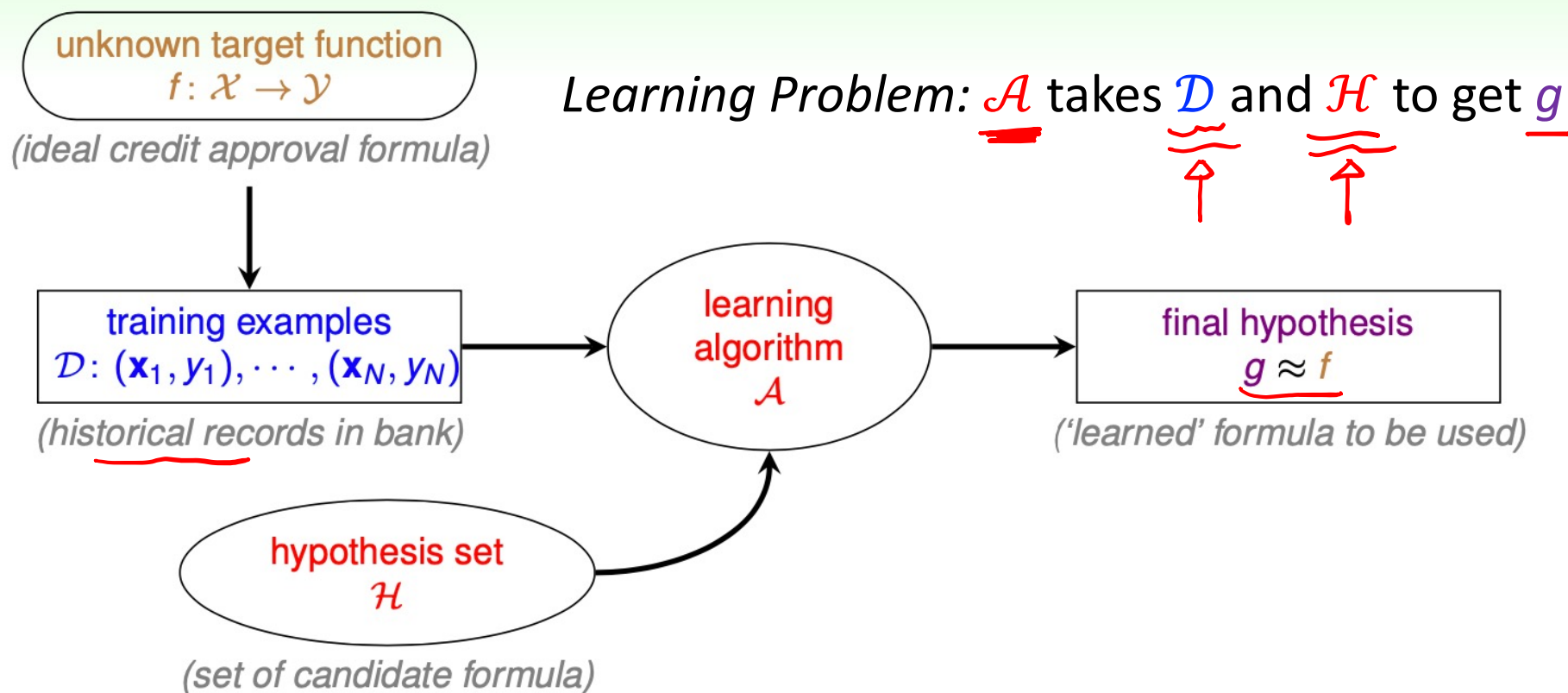
The Learning Model



- assume $g \in \mathcal{H} = \{h_k\}$, i.e. approving if
 - h_1 : annual salary > NTD 800,000
 - h_2 : debt > NTD 100,000 (really?)
 - h_3 : year in job ≤ 2 (really?)
- hypothesis set \mathcal{H} :
 - can contain good or bad hypotheses
 - up to \mathcal{A} to pick the ‘best’ one as g

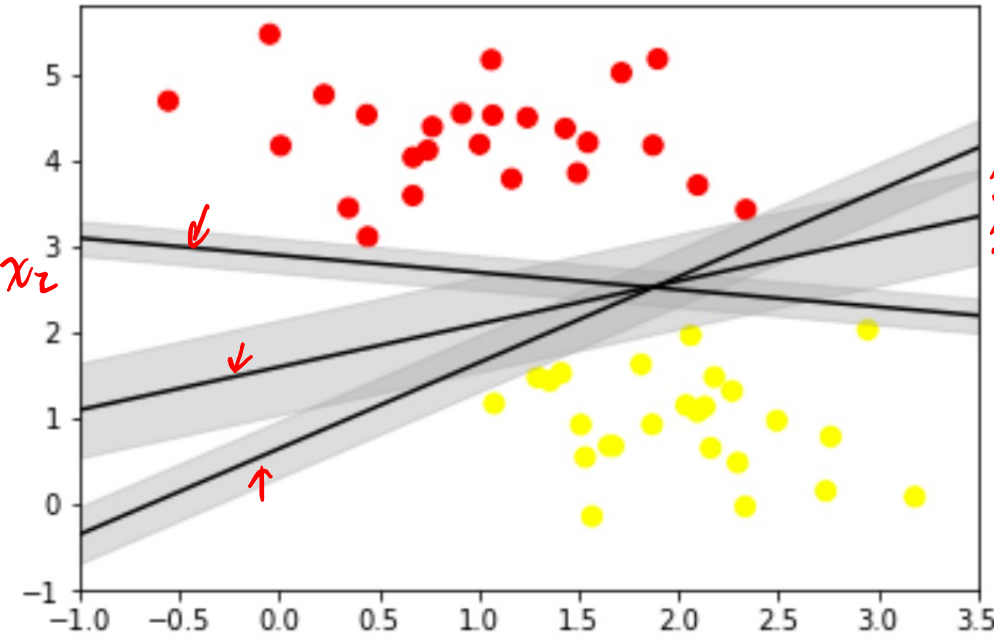
learning model = \mathcal{A} and \mathcal{H}

Practical Definition of Machine Learning



↙
machine learning:
 use **data** to compute **hypothesis g**
 that approximates **target f**

Learning Problem in Practice: **learning algorithm** \mathcal{A} takes training examples \mathcal{D} and **hypothesis set** \mathcal{H} to get final hypothesis g .



Fun time: Quick Check (who is who)

- Learning algorithm \mathcal{A}
Support Vector Machines

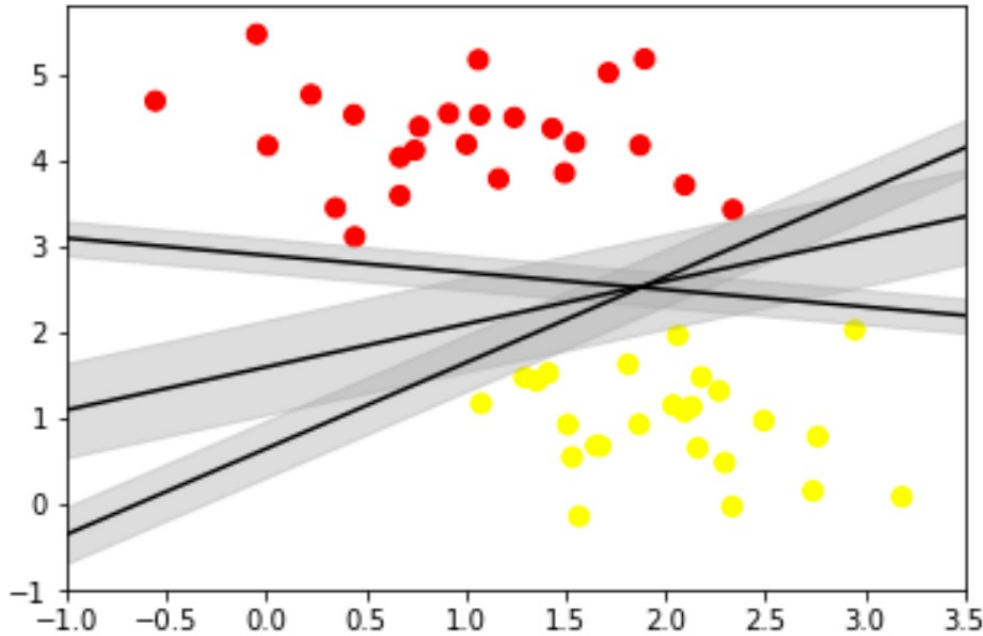
https://en.wikipedia.org/wiki/Support-vector_machine

$$w_1 x_1 + w_2 x_2 + b = 0$$

In support vector machines (SVM), the line that maximizes this margin is the one we will choose as the optimal model.

<https://www.sli.do/#073374>

Learning Problem in Practice: **learning algorithm \mathcal{A}** takes training examples \mathcal{D} and **hypothesis set \mathcal{H}** to get **final hypothesis g** .



$$w_1x_1 + w_2x_2 + b = 0$$

In support vector machines (SVM), the line that maximizes this margin is the one we will choose as the optimal model.

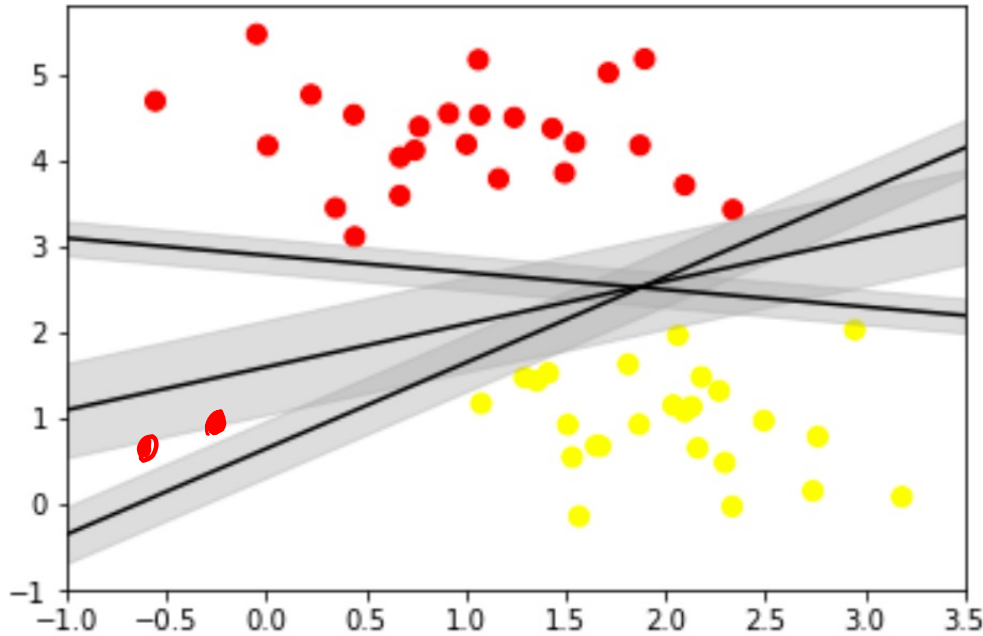
Fun time: Quick Check (who is who)

- Hypothesis set \mathcal{H}

All the lines from the real numbers w_1 , w_2 , b ∞

<https://www.sli.do/#073374>

Learning Problem in Practice: **learning algorithm \mathcal{A}** takes training examples \mathcal{D} and **hypothesis set \mathcal{H}** to get **final hypothesis g** .



Fun time: Quick Check (who is who)

- Learning algorithm \mathcal{A}
- Hypothesis set \mathcal{H}
- Training Examples \mathcal{D}
- Final hypothesis g
- Target function f *unknown*

$$w_1x_1 + w_2x_2 + b = 0$$

In support vector machines (SVM), the line that maximizes this margin is the one we will choose as the optimal model.

Summary

Learning? What do we mean?

- Machine learning involves building mathematical models to help understand data.
- In practice, we use data to compute hypothesis g that approximate unknown target f .
- In practice, learning algorithm \mathcal{A} takes training examples \mathcal{D} and hypothesis set \mathcal{H} to get final hypothesis g .
- "Learning" enters the picture when we give these models tunable parameters that can be adapted to observed data; in this way the program can be considered to be "learning" from the data.

$w_1 \ w_2 \ b \leftarrow \text{SVM}$

Is learning feasible?

Traditional Programming

$+$ $-$
 \times $=$

Rules



Computer



The computer can perform the task it has been *instructed* to do.

Machine Learning



Data



Computer



The computer can perform the task it has *learned* to do.

- Can computer automatically learn these rules from data?
 - We use data to compute hypothesis g that approximates target f
- If so, can these rules be applied for new data to guess answers? (generalization)

Is learning feasible?

Fun Time: Can final hypothesis g predict new data?

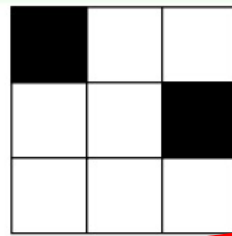
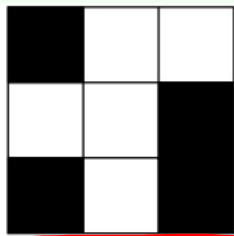
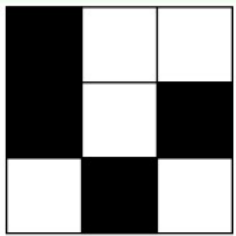
(1) Yes

(2) No

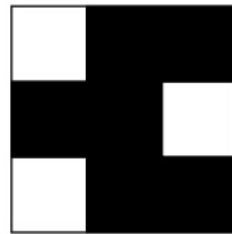
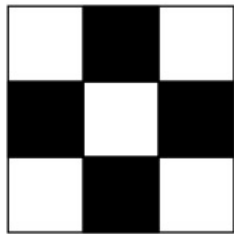
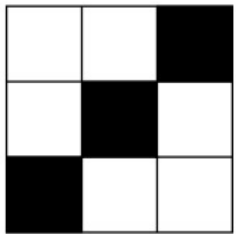
(3) Maybe

[https://www.sli.do/
#073374](https://www.sli.do/#073374)

A Learning Puzzle

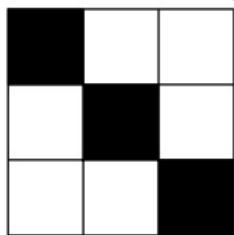


$$y_n = -1$$



$$y_n = +1$$

new



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

Fun Time: $g(\mathbf{x})$ prediction of new data

(1) -1

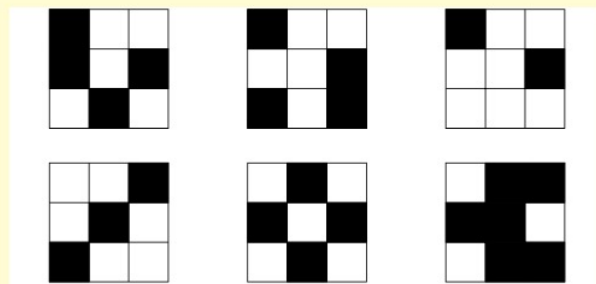
(2) +1

<https://www.sli.do/#073374>

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

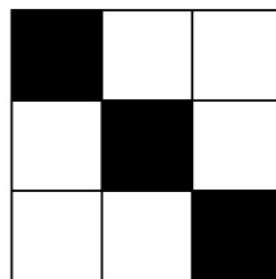
Two Controversial Answers

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



$$y_n = -1$$

$$y_n = +1$$



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

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

- symmetry $\Leftrightarrow +1$
- (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

Boolean

\mathbf{x}_n	$y_n = f(\mathbf{x}_n)$
<u>0</u> <u>0</u> <u>0</u>	○ 1 ●
<u>0</u> <u>0</u> <u>1</u>	× -1 ○
<u>0</u> <u>1</u> <u>0</u>	×
0 1 1	○
1 0 0	×
↑ ↑ ↑	

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

Let us do a two-bit case to explore (1) the dimension of the input space (2) all the samples in the input space (3) all the possible hypotheses (and we can pick up one as our target function)

A 'Simple' Binary Classification Problem

5

$2^3 = 8$

$2^{2^3} = 2^8 = 256$

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	○?	○?	○	×	○	×	○	×	○	×

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

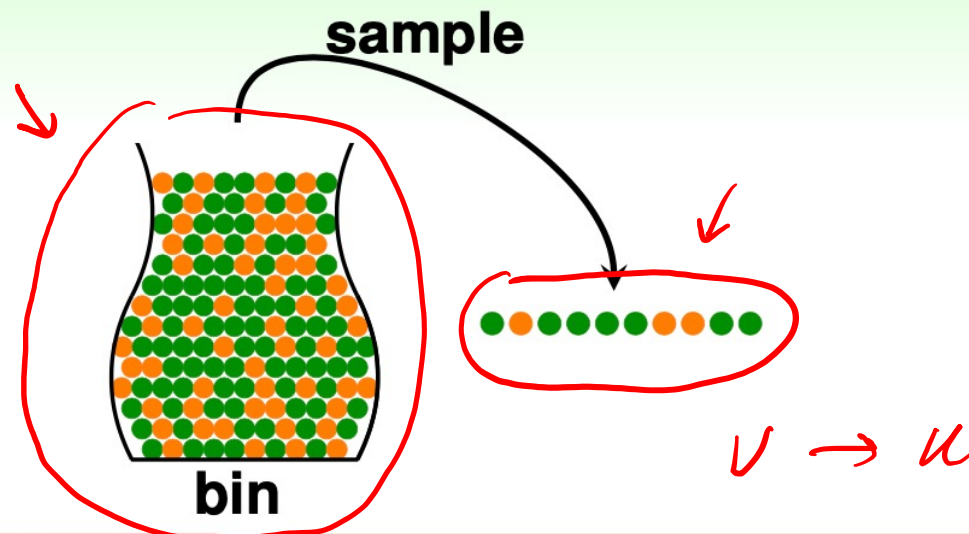
learning from \mathcal{D} (to infer something outside \mathcal{D})
is doomed if **any 'unknown' f can happen.** :-)

Is learning doomed (完蛋了)?
If so, this will be a very short course!!!



Probability to recuse!

Statistics 101: Inferring **Orange** Probability



bin

assume

orange probability = μ ,

green probability = $1 - \mu$,

with μ **unknown**

sample

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 ?

Possible versus Probable

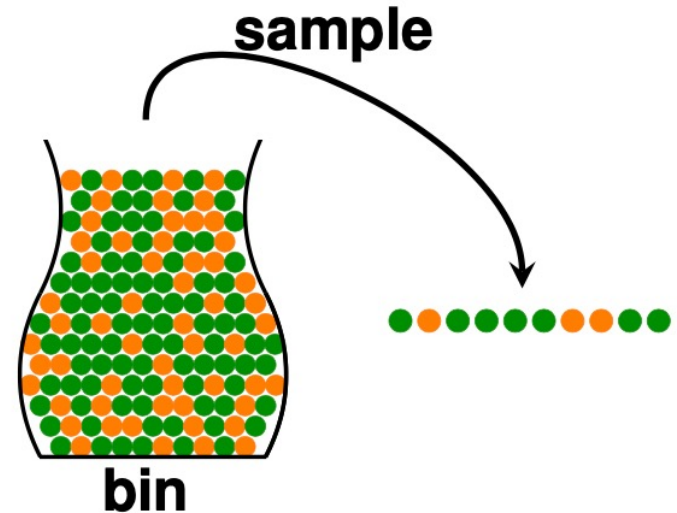
does **in-sample** ν say anything about out-of-sample μ ?

No!

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

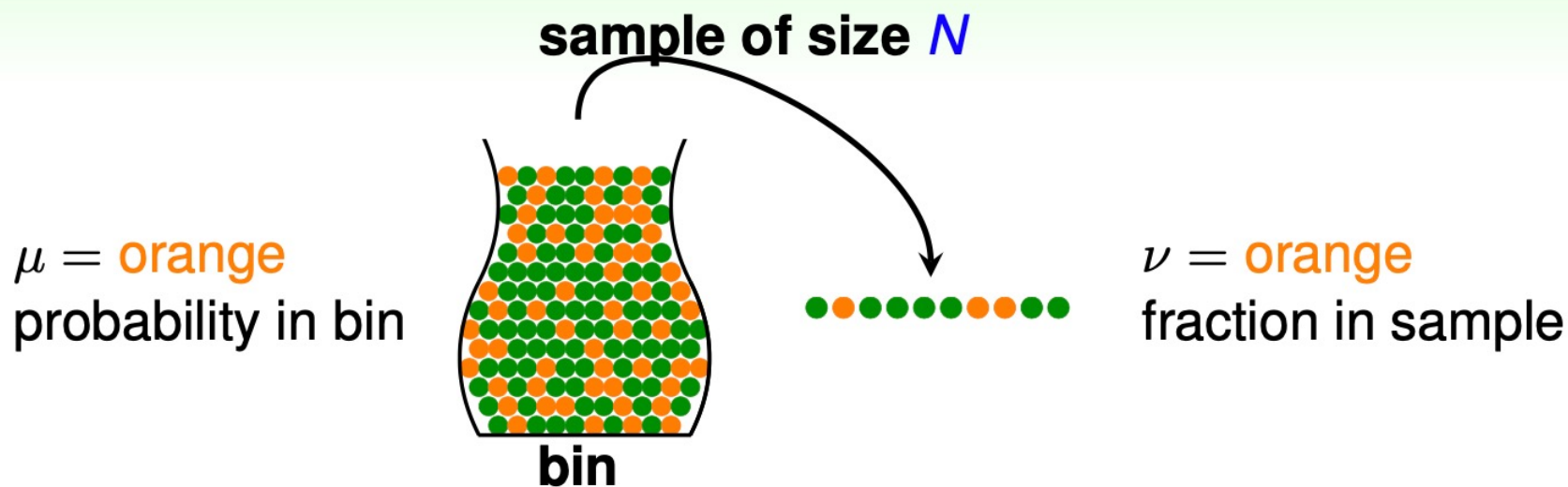
Yes!

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)

$$\mathbb{P} [\underbrace{|\nu - \mu|}_{\text{error}} > \epsilon] \leq \underbrace{2 \exp(-2\epsilon^2 N)}_{\text{probability of error}}$$

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

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

可能大概对

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

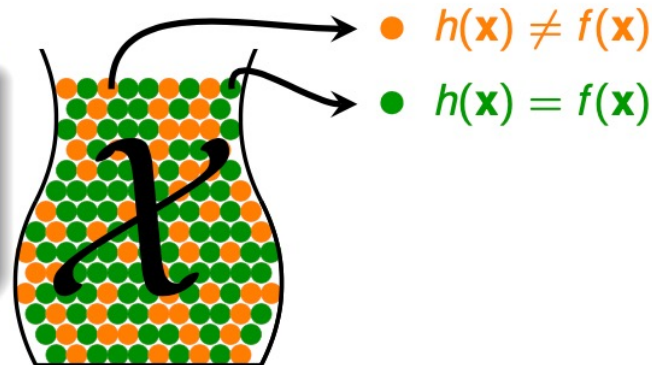
if **large** N & **i.i.d.** \mathbf{x}_n , can probably infer
unknown $\llbracket h(\mathbf{x}) \neq f(\mathbf{x}) \rrbracket$ probability
by known $\llbracket h(\mathbf{x}_n) \neq y_n \rrbracket$ fraction

i.i.d. independent and identically distributed

learning

- fixed hypothesis $h(\mathbf{x}) \stackrel{?}{=} \text{target } \underline{f(\mathbf{x})}$
- $\underline{\mathbf{x}} \in \mathcal{X}$
- h is **wrong** $\Leftrightarrow h(\mathbf{x}) \neq f(\mathbf{x})$
- h is **right** $\Leftrightarrow h(\mathbf{x}) = f(\mathbf{x})$
- check h on $\mathcal{D} = \{(\mathbf{x}_n, \underbrace{y_n}_{f(\mathbf{x}_n)})\}$

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



Boolean Learning Example: Take Two

Let us consider a Boolean target function (i.e., $\mathcal{Y} = \{0, 1\}$) over a four-bit vector representation of input space $\{0000, 0001, \dots, 0111, 1000, 1001, \dots, 1111\}$.

Q: For this example, what is the dimension of the input space \mathcal{X} ?

Q: For this example, how big is the entire input space \mathcal{X} ?

Q: For this example, how big is the entire Boolean hypothesis set \mathcal{H} ?

See [Boolean_Learning_Example.pdf](#)

[Boolean_Learning_Example.ipynb](#)

[https://www.sli.do/
#073374](https://www.sli.do/#073374)

Summary

Is learning feasible?

- Learning is only feasible in a *probabilistic* way and we can **predict** something useful outside the training set \mathcal{D} using only \mathcal{D} .
- We don't insist on using any particular probability distribution, or even on knowing what distribution is used.
However, whatever distribution we use for generating the samples, we must also use when we evaluate how well g approximates the *unknown* target function f .
- The hypothesis g is not fixed ahead of time before generating the data, because which hypothesis is selected to be g depends on the data.

unknown target function

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

(ideal credit approval formula)

training examples

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

(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)

Learning Problem: \mathcal{A} takes \mathcal{D} and \mathcal{H} to get g

Labels

Training Data

Machine Learning
Algorithm

New Data

Predictive Model

Prediction

Supervised Learning ✓

Summary

Learning? What do we mean?

Is learning feasible?

- Machine learning: use data to compute **hypothesis g** that approximate unknown **target f** .
- In practice, **learning algorithm \mathcal{A}** takes training examples **\mathcal{D}** and **hypothesis set \mathcal{H}** to get **final hypothesis g** .
- Learning is only feasible in a *probabilistic* way and we can predict something useful outside the training set \mathcal{D} using only \mathcal{D} .