

机器学习

→ 机器自动找 函数

For Example

1. 言音辨識

输入：聲音 信号

输出：語言辨識結果 (文字)

2. 影像辨識

输入 = 图片

输出 = 图片有什麼樣的物件 / 東西

3. Playing Go

输入：棋盤上黑子 & 白子的位置

输出：下一步落子的位置

4. 对话式系统

输入：你要对机器說的語言

输出：机器的回应

想要找什麼函数？

• Regression

输出 = 數值

Want 机器預測未來某個時刻的 PM 2.5

⇒ 找 function 输入 = 過去的臭氧濃度 / 天氣等
输出 = 明天某一時刻的 PM 2.5 數值

• binary classification

输出 = Yes / No

e.g. RNN 输入 = 一个句子

输出 = 句子是正面 / 负面

- Multi-class classification

輸入：讓機器做選擇題 (n個選項)

 出：選出一个 Ans.

 eg: CNN

 輸入一個圖片，它是一個麵包

B: 麵

C: 湯

which?

- Generation (生成)

→ 產生有結構的複雜東

→ 訓繡機器學習如何創造

eg: 翻譯：產生文句

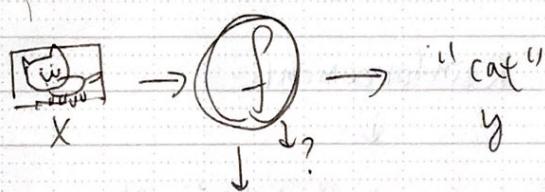
B: 產生二次元人物

那怎麼告訴機器去找什麼樣的函式？

- Supervised Learning

• suppose 已經想好 輸入是一張圖片：X

 輸入是這張圖片屬於什麼樣的
 類別：y



need 教練資料 (Labeled Data)

↳ need 收集 大量圖片

② 理想 / 正確的輸出是什麼

↳ 也就是 Labeled 的意思

→ 提供給機器有 labeled 的資料進行學習

 稱 "Supervised Learning"

 then 机器就可以評估一個函式的好坏 \Rightarrow Loss

Loss Function

Labeled Data

x_1		$y_1 = \text{cat.}$
x_2		$y_2 = \text{cat}$
x_3		$y_3 = \text{Dog}$
x_4		$y_4 = \text{Dog}$

 $"x_1" "x_2" "x_3" "x_4"$

$$\downarrow [f_1]$$

"Dog" "Dog"
 "Dog" "Dog"

$$\underline{\text{Loss}} = 50\%$$

↳ 越小越好

 $"x_1" "x_2" "x_3" "x_4"$

$$\downarrow [f_2]$$

"cat" "cat"
 "Dog" "Dog"

$$\underline{\text{Loss}} = 0\%$$

接下來機器會自動找出 Loss 最低的函式

↓ will use 什麼樣的演算法 找最低的 Loss

Reinforcement Learning

- alpha Go

Supervised V.S. Reinforcement



Now



Next move

"2-3"

機器 & 機器 play

⇒ win

機器 & somebody play

Now



Next

3-3

↑
 機器想辦法自己提高正確率

win / lose ⇒ Reward (31章)

• Unsupervised Learning

→ 給一大堆圖片，但沒有標注

机器怎麼找出你想要的函式

First: 級定函式尋找的範圍

ex A: linear function (前兩個 HW)

B: Network Architecture

Use RNN: seq to seq

CNN = a lot of

ex
 Meta Learning
 Unsupervised Learning
 Anomaly Detection
 Transfer Learning
 Explainable AI

函式尋找方法 - Gradient Descent

前兩個作業 = Implement the algorithm by yourself

前三步研究

- Explainable AI

- ↳ why your 程式 think the picture is Cat?

- Adversarial Attack

- ↳ If 被惡意攻擊 system 會發生什麼事

- ↳ If 加入的雜訊不是一般的，是被刻意製造的？

- Network compression

- ↳ 把很大的 Network 縮小

- 放入小的 page device 上

Use by CNN 結果

■ Anomaly Detection

if 看到訓練沒看到的，能否知道要 say "我不知道"

■ Domain Adversarial

[Training Data \Rightarrow have same distribution
 Testing Data]

ex 手寫辨認 \Rightarrow 會有很高的正確率

but if Training data is 黑白

Testing data is colorful

\Rightarrow 正確率 5%~5%. Very low

④ Meta Learning

→ learn to learn

(學習如何學習)

讓機器自己發明一個學習 Algo

⑤ Life-long Learning (終身學習)

i.e. Continuous Learning, Never Ending Learning

→ what is it. 難莫？

Date 110. 1. 23

Regression : Case study

For example.

1. Stock Market Forecast (股票預測系統)

$$f(\text{各种股票起伏的数据}) = \text{明天道瓊工業指數的數}$$

2. Self-driving car (无人車、自動車)

$$f(\text{无人車的紅外線感測的sensor}) = \text{方向盤的角度}$$

3. Recommendation (推薦 system)

$$f(\text{使用者A 商品B}) = \text{使用者A 購買商品B的可能性}$$

4. 預測 Pokemon 的 CP 值

$$f(\text{某一隻Pokemon}) = \text{進化后的CP值}$$

包含 information

$\downarrow y$

包含 X_{cp} , X_s , X_{hp} , X_w , X_n

物种 生命值 重量 高度

ML step

1. 找一个 model (function set)

2. 定义 function set 里面某个 function

找一个 function, then evaluation 做好

3. 找一个最好的 function

$$\text{Step1 Model.} = \underline{\underline{y}} = b + \underline{\underline{w}} \cdot \underline{\underline{X_{cp}}}$$

進化后 CP 常数 進化前 CP
 參數

$$\text{Linear model } y = b + \sum_j w_j x_j$$

bias weigh Input Xj's attribute
(ex. pokemon CP%)

Step 2 = Goodness of function.

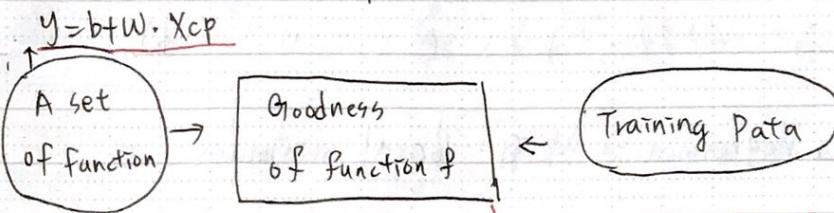
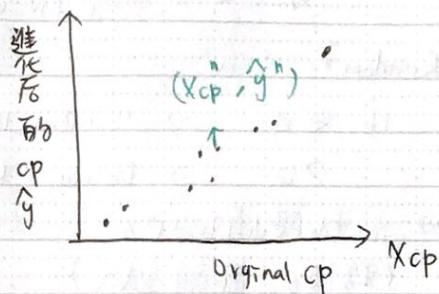
→ 收集 training data, 才能找到這個 function

function input: $x^1 = \text{水} \rightarrow x^2 = \text{伊布}$

output $\hat{y}^1 = 979 \quad \hat{y}^2 = 1420$

Training Data. = 10 pokemons.

$(x^1, \hat{y}^1), (x^2, \hat{y}^2), \dots, (x^{10}, \hat{y}^{10})$



Loss function L

[Input: a function.]

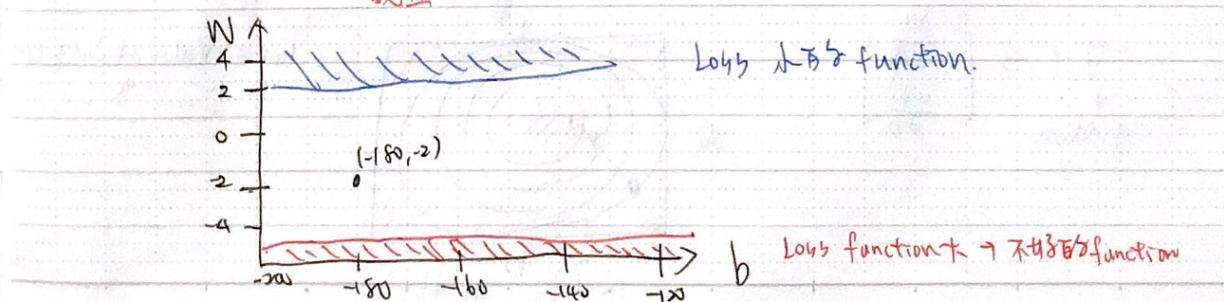
[Output: how bad it is]

$L(f) = L(w, b)$ // 衡量一組參數的好壞

$$= \sum_{n=1}^{10} (\hat{y}^n - (b + w \cdot x_{cp}^n))^2$$

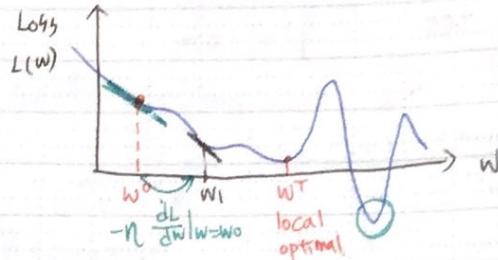
真正的
數值

現在的 function 的輸出
(預測的數值)
Estimation error:



Step 3: Gradient Descent:

Consider loss function $L(w)$ 可微



$$w^1 = w^0 - \eta \frac{dL}{dw} \Big|_{w=w^0}$$

和 learning rate
 增加 = 事先定好的數值
 是反向的
 越大 ↓ 調整速度

$$w^2 = w^1 - \eta \frac{dL}{dw} \Big|_{w=w^1}$$

- 找一个初值的点 (Random) : w^0
- Compute $\frac{dL}{dw} \Big|_{w=w^0}$ If $< 0 \rightarrow$ Increase w^0
 $> 0 \rightarrow$ Decrease w^0

If $\frac{dL}{dw} \Big|_{w=w^0}$ 越大 \rightarrow 在走陡峭的地方
 (移动的距离会加大)

↓ 小 (↓ 大)

* 在 linear regression = 没有 local minimum

two parameters $w^*, b^* = \arg \min_{w,b} L(w, b)$

- 找一个初值 w^0, b^0
- Compute $\frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0}, \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$

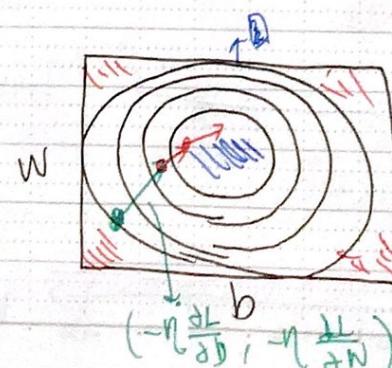
$$w^1 = w^0 - \eta \frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0}, \quad b^1 = b^0 - \eta \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$$

$$\text{Compute } \frac{\partial L}{\partial w} \Big|_{w=w^1, b=b^1}, \quad \frac{\partial L}{\partial b} \Big|_{w=w^1, b=b^1}$$

$$w^2 = w^1 - \eta \frac{\partial L}{\partial w} \Big|_{w=w^1, b=b^1}, \quad b^2 = b^1 - \eta \frac{\partial L}{\partial b} \Big|_{w=w^1, b=b^1}$$

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix}$$

gradient



Color = Value of Loss (w, b)
 ↓ 小 ↑ 大



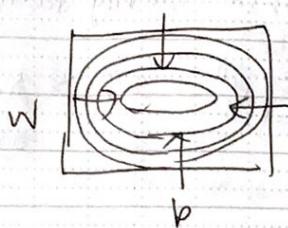
In linear regression,

the loss function L is convex.

* Convex $\Leftrightarrow \nabla^2 L \geq 0$

\rightarrow 沒有 local optima. 的位置

因為 $\nabla^2 L = -R$



$$L(w, b) = \sum_{n=1}^{10} (\hat{y}^n - (b + w \cdot x_{cp}^n))^2$$

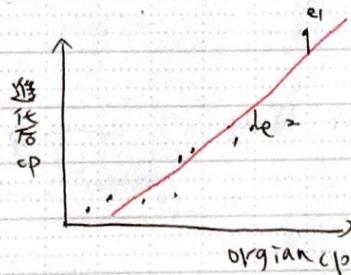
$$\frac{\partial L}{\partial w} = \sum_{n=1}^{10} 2 \cdot (-x_{cp}^n) (\hat{y}^n - (b + w \cdot x_{cp}^n))$$

$$\frac{\partial L}{\partial b} = \sum_{n=1}^{10} 2 \cdot (-1) (\hat{y}^n - (b + w \cdot x_{cp}^n))$$

How the results?

$$b = -188.4$$

$$w = 2.9$$



Average error on Training Data

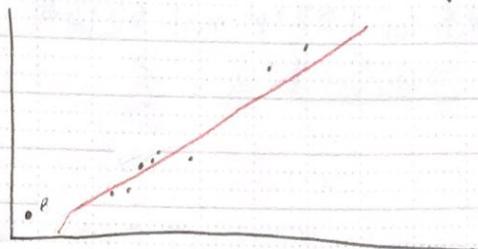
$$= \frac{1}{10} \sum_{n=1}^{10} e^n = 31.9$$

Generalization

(真正关心的是，沒有看過的新 data 其 error 是多少)
testing data

- Another 10 pokémon as testing data.

$$b = -188.4 \quad w = 2.17$$



Average Error on
testing data

$$= \sum_{n=1}^{10} e^n = 35 > \text{Average Error on Training Data}$$

How can we be better

→ Selecting another Model

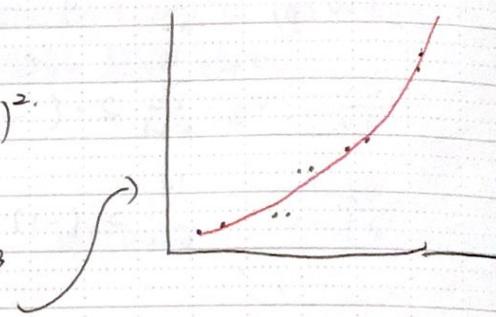


$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2$$

Best function.

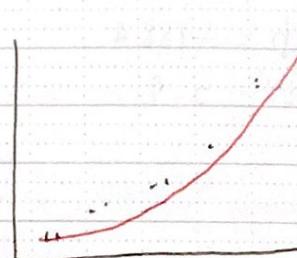
$$b = -10.3 \quad w_1 = 1 \quad w_2 = 2.7 \times 10^{-3}$$

$$\text{Average Error} = 15.4$$



Testing

$$\text{Average Error } 18.4$$



Underfitting mode

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3$$

Best Function

$$b = 6.4 \quad w_1 = 0.66 \quad w_2 = 4.3 \times 10^{-3}$$

$$w_3 = -1.8 \times 10^{-6}$$

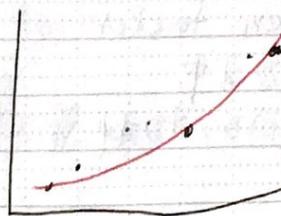
$$\text{Average Error} = 15.3$$



Testing

Average Error 18.1

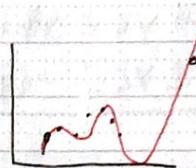
⇒ Slightly better.



$$\# \quad y = b + w_1 \cdot X_{cp} + \dots + w_5 (X_{cp})^5$$

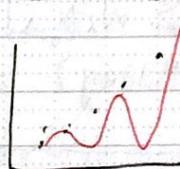
Best Function

Average Error = 12.8

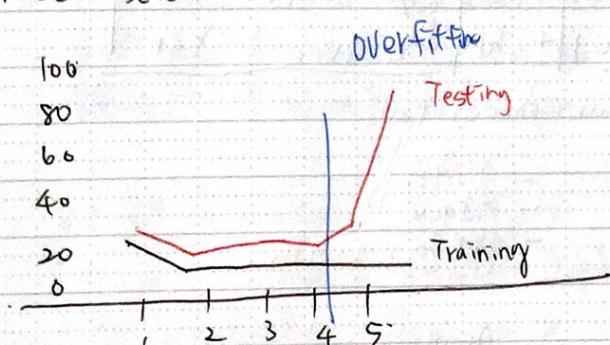


Testing

Average Error = 232.1



Model Selection



	Training	Testing
1	31.9	35
2	15.4	18.4
3	15.3	18.1
4	14.9	28.2
5	12.8	232.1

⇒ Overfitting: 更複雜的 model. 並沒有在 testing data.

有更好的表現

(在 best function Good but 在 testing data bad)

hidden factor

⇒ 物种

→ 不同的物种带不同的 linear function.

$$X_s = \text{species}, f X$$

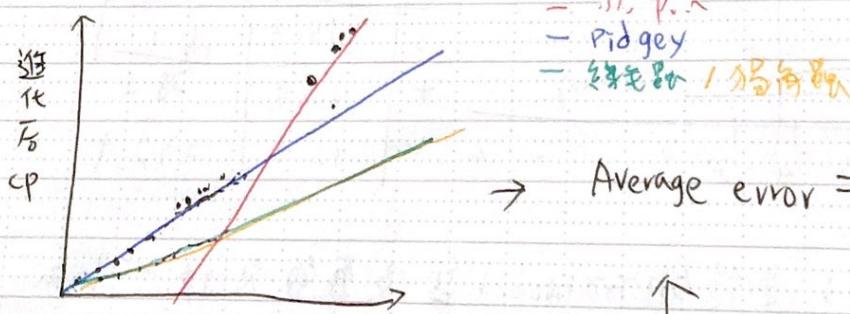
$$X \rightarrow \begin{cases} \text{if } X_s = \text{波波} & y = b_1 + w_1 \cdot CP \\ \text{if } X_s = \text{獨角獸} & y = b_2 + w_2 \cdot CP \\ \text{if } X_s = \text{綠毛蟲} & y = b_3 + w_3 \cdot CP \\ \text{if } X_s = \text{伊布} & y = b_4 + w_4 \cdot CP \end{cases} \rightarrow y$$

$$\Rightarrow y = b + \sum w_i \cdot x_i$$

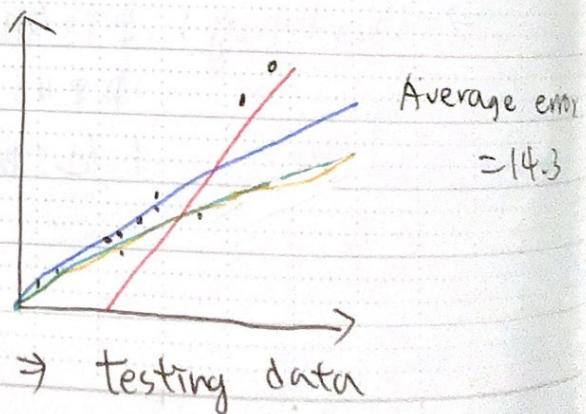
$$\delta(X_s = \text{Pidgey}) = \begin{cases} 1 & \text{if } X_s = \text{Pidgey} \\ 0 & \text{o.w.} \end{cases}$$

$$\Rightarrow y = b_1 \cdot \delta(X_s = \text{Pidgey}) + w_1 \cdot \delta(X_s = \text{Pidgey}) \cdot CP + b_2 \cdot \delta(X_s = \text{Weedle}) + w_2 \cdot \delta(X_s = \text{Weedle}) \cdot CP + b_3 \cdot \delta(X_s = \text{Caterpie}) + w_3 \cdot \delta(X_s = \text{Caterpie}) \cdot CP + b_4 \cdot \delta(X_s = \text{Eevee}) + w_4 \cdot \delta(X_s = \text{Eevee}) \cdot CP$$

→ is linear model too



→ Average error = 3.8.

 \Rightarrow training data fit by different properties


Are there any other hidden factors?

→ Weight? Height? HP?

Back to step 1. (Redesign the model again)

$$\begin{aligned}
 & \text{If } X_S = \text{Pidgey} \quad y' = b_1 + w_1 X_{CP} + w_5 (X_{CP})^2 \\
 & \text{If } X_S = \text{Weedle} \quad y' = b_2 + w_2 X_{CP} + w_6 (X_{CP})^2 \\
 & \text{If } X_S = \text{Caterpie} \quad y' = b_3 + w_3 X_{CP} + w_7 (X_{CP})^2 \\
 & \text{If } X_S = \text{Eevee} \quad y' = b_4 + w_4 X_{CP} + w_8 (X_{CP})^2
 \end{aligned}
 \rightarrow y$$

$$y = y' + w_9 X_{HP} + w_{10} (X_{HP})^2 + w_{11} X_H + w_{12} (X_H)^2 + w_{13} X_W + w_{14} (X_W)^2$$

Training error = 1.9
 Testing error = 102.3. \Rightarrow Overfitting ^{now} \Rightarrow Regularization

Back to step 2 = Regularization

\Rightarrow 重新 redefine 我們的 loss function, 把一些 knowledge 放入
 讓我們找到比較好的 function.

$$y = b + \sum w_i x_i \quad \text{The function with smaller } w_i \text{ is better}$$

$$L = \sum (y - (b + \sum w_i x_i))^2 + \lambda \sum (w_i)^2$$

\Rightarrow 期待參數趨近 0 的 function.

\because function 較平滑

λ 的 output 對輸入的變化是比較不敏感的

$$\text{if } y = b + \sum w_i x_i$$

$$\text{now } y + \frac{\Delta y}{\lambda} x_i = b + \sum w_i (x_i + \Delta x_i)$$

$\frac{\Delta y}{\lambda}$ if it 趨近 0, 變化輸入

\Rightarrow if now have 雜訊, the effect 不太會受影響

Something to know.

λ	Training	Testing
0	1.9	102.3
1	2.3	68.9
10	3.5	25.7
100	4.1	11.1
1000	5.6	12.8
10000	6.3	18.9
100000	8.5	26.8

λ 越大 $\propto \sum(w_i)^2$
 \Rightarrow smooth 的那个 regularization 那部分
 的影响力越大
 \Rightarrow 找到的 function 越平滑
 \Rightarrow training data 上的 error 越大
 \Rightarrow 越倾向考虑 w 本来的值
 而减少考虑 error
 \Rightarrow testing data 上的 error 可能较小

Hence we prefer smooth function
 but don't be too smooth (ex: 平坦)

How smooth?

\Rightarrow Select λ obtaining the "best" model

$$\text{Ans: } \lambda = 100$$

* 調整 bias (b) 的大小和 function 的平滑程度无关

↓
∴ only 上下 move

Conclusion

- 原本的 CP & 物种几乎决定进化后 CP 值
(but 1% 有 hidden factors?)
- Gradient descent

Next What does the error come from?
lecture → validation

Where does the error come from?

error due to

$L \leftarrow$ bias
variance

⇒ then can improve your model.

Estimator

$$\hat{y} = \hat{f}(\text{pokemon})$$

→ only the company know \hat{f} .

→ from the training data, we can find f^*
(f^* is the estimator of \hat{f})

① for example

② estimator of mean of variable X

→ Assume the mean = M , variance = σ^2 .

③ estimator of mean M .

→ Sample N points = { X^1, X^2, \dots, X^N }

We know that $m = \frac{1}{N} \sum_{n=1}^N X^n \neq M$

$$\begin{aligned} \text{but } E(m) &= E\left[\frac{1}{N} \sum_{n=1}^N X^n\right] \\ &= \frac{1}{N} E\left(\sum_{n=1}^N X^n\right) \\ &= M \end{aligned}$$

⇒ unbiased (不偏估計)

$$\text{Var}(m) = \frac{\sigma^2}{N}$$

↑
分散的降低度

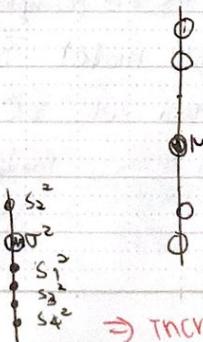
⇒ N 越大會越集中

$$S^2 = \frac{1}{N} \sum_n (X^n - m)^2$$

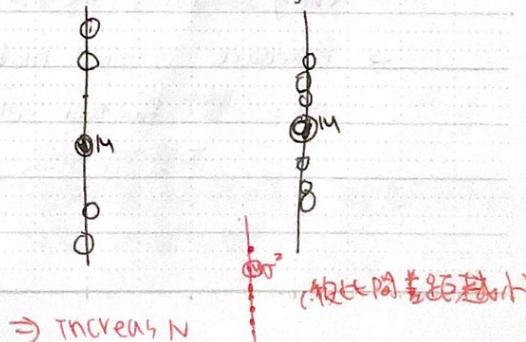
$$E(S^2) = \frac{N-1}{N} \sigma^2 \neq \sigma^2$$

(S^2 is smaller than σ^2)

Smaller N

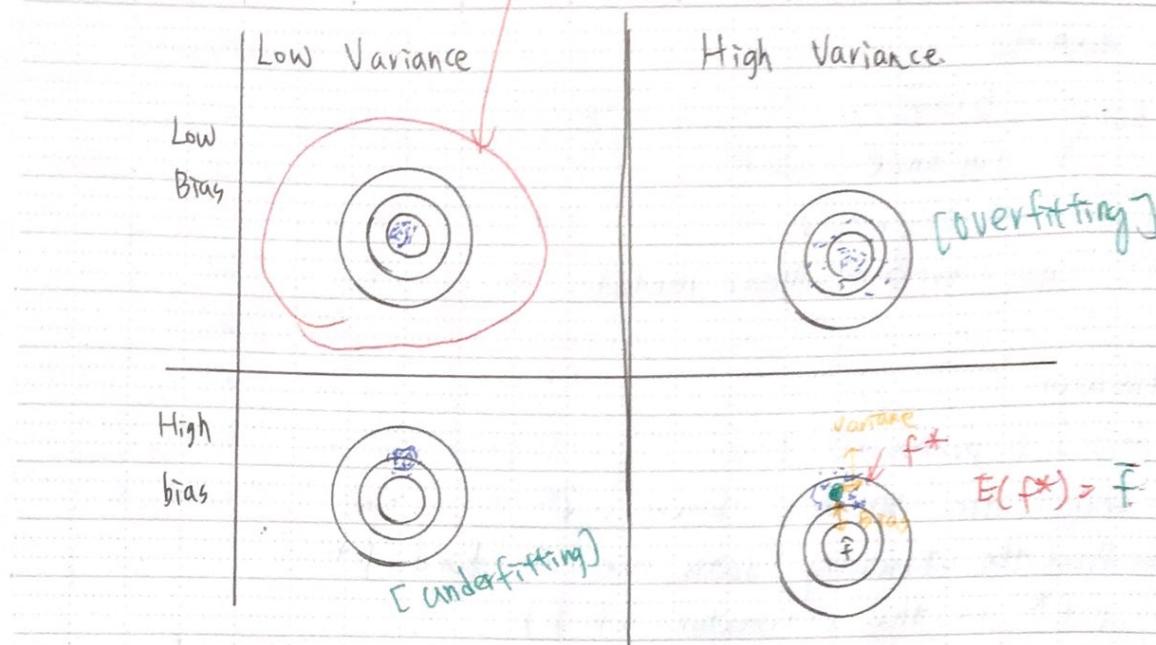


larger N



(相对偏差较小)

期待的情況



Error 取決兩件事

1. 瞄準的位置在哪裡 [bias有多大]

→ estimator 是不是 bias (有沒有偏準)

2. 瞄準的這個位置，子彈射出後仍有偏移 [Variance有多大]

→ f^* 和 f 中間的距離 = Variance

平行宇宙 (做很多次實驗)

In different universes, we use same model.

but obtain different f^*

* Use simple model, variance is small



Why? why 複雜 (complex) model, variance is bigger?

→ Because simple model 受到不同 data 的 impact 是小的

e.g.: 整个 function set. 就一个 function $f(x) = C$

, 木管抓几次, model 都 same, variance = 0

Bias

→ If we average all f^* , is it close to \hat{f} .

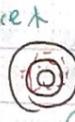
$$E(f^*) = \bar{f}$$

Large Bias Small Bias

[Underfitting]



Small variance



[Overfitting]



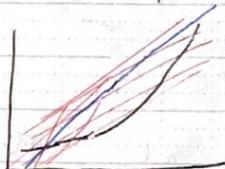
Large variance

function space \propto 有包含 target.

不同 function 間的 bias 有多大?

⇒ We can't evaluate it, because we don't know \hat{f}

⇒ 假設 \hat{f}



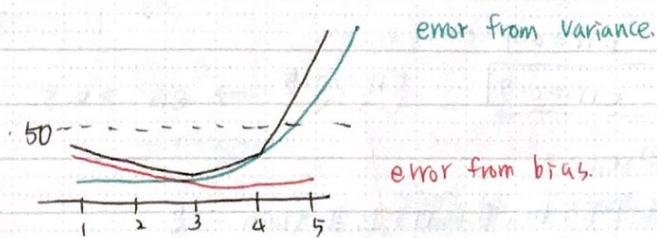
- \bar{f}

- \hat{f}

- $5000 \approx f^*$

* 比較 simple 模型 : bias 大 (A)

complex. " : " 小 (B)



future work. 要 improve model 要走什麼.

• how is bias \propto ? variance \propto ?

bias \propto : model 沒有辦法 fit 你的 training 的 examples

Variance \propto : you can fit the training data., but in testing data. have large error.

→ need redesign your model :: model 沒有包含 target

↳ add more feature. in input more complex model

→ ① More data. (generate 假的 data ex 手寫辨識角度轉15度, 反轉圖片左右顛倒語言, use 麥克風)

② Regularization ⇒ 曲線會更平滑

Model Selection

- 在 bias & variance 之間做一些 trade-off
- Select a model, its variance is small, bias is small
⇒ 合起來給我們最小的 testing data Err
- We should Not do

Training set.

$\left[\begin{array}{l} \text{model 1} \rightarrow \text{Err} = 0.9 \\ \text{model 2} \rightarrow \text{Err} = 0.7 \\ \text{model 3} \rightarrow \text{Err} = 0.4 \end{array} \right]$

Testing set

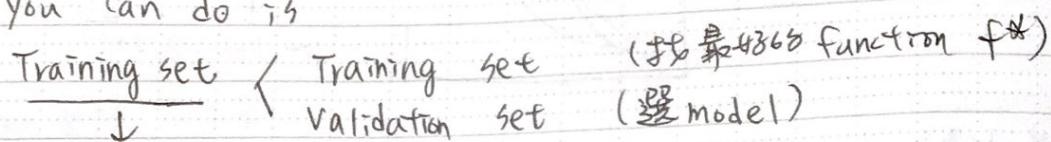
(your) public
see

(real!) Testing set

private set

$\rightarrow \text{Err} > 0.5$

\Rightarrow you can do is



Training set	Validation set	public Testing set	private Testing set
Model 1 $\rightarrow \text{err} = 0.9$			
Model 2 $\rightarrow \text{err} = 0.7$			
Model 3 $\rightarrow \text{err} = 0.5$		<u>$\text{err} > 0.5$</u>	$\rightarrow \text{err} > 0.5$

是要用 Model 3,

但用全部的 data 在 model 3 上面 train \rightarrow

don't recommend

\rightarrow because you will 把 public testing set 的 bias 考虑进去

\rightarrow 在 public set 上 找到的 performance 无法反应

你在 private set 的 performance

* public set 不是真正的结果

N-fold Cross Validation

Training set	Model 1	Model 2	Model 3
Train	Train	Val	
Train	Val	Train	
Val	Train	Train	
	Err = 0.2	Err = 0.4	Err = 0.4
	Err = 0.4	Err = 0.5	Err = 0.5
	Err = 0.3	Err = 0.6	Err = 0.3
	Avg Err = 0.3	Avg Err = 0.5	Avg Err = 0.4