

# 兩位跨域者的深度學習之路

游為翔 & 楊証琨

中研院資訊所資料洞察實驗室

@ 台灣人工智慧年會 – 2017/11/09



中央研究院  
Academia Sinica



資料洞察實驗室



# 講者



游為翔 Sean Yu



楊証琨 Jimmy Yang

背景

心理學 / 神經科學

土木工程 / 交通

專長

✓ 消費及使用者行為分析  
✓ 深度學習於影像辨識研究

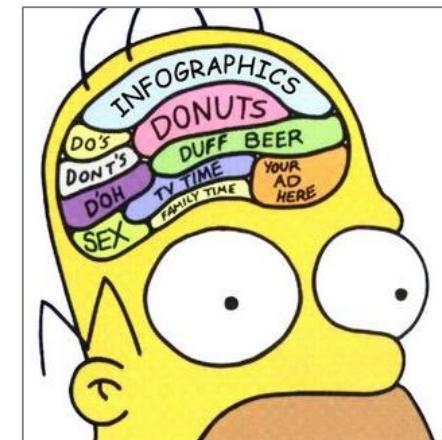
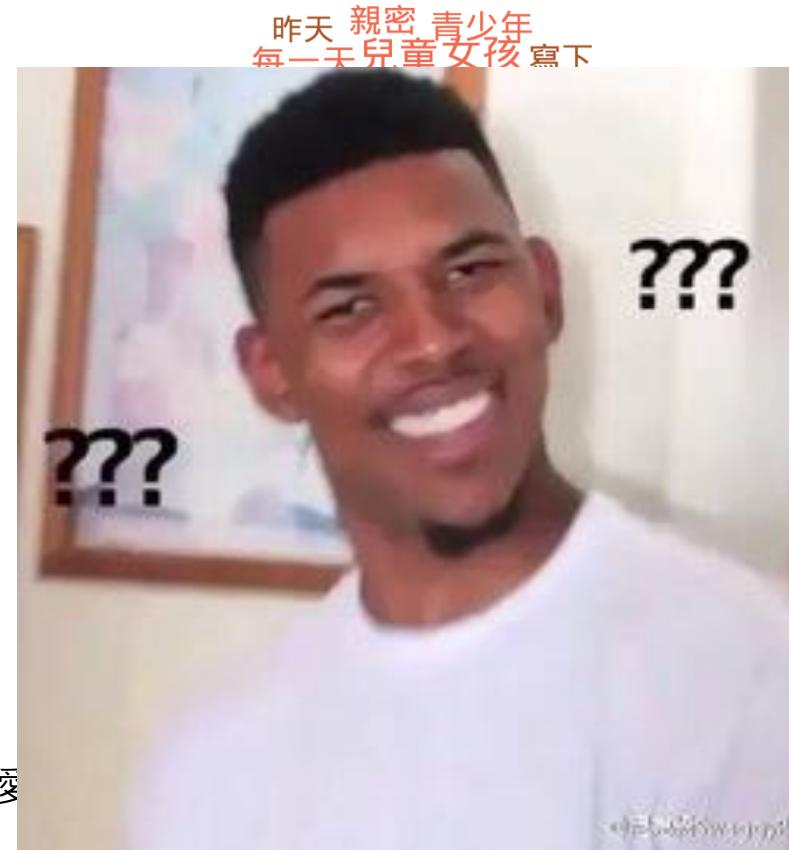
✓ 社群媒體資料處理分析  
✓ 製造工業影像資料分析

# 心理系到底在學什麼

激動 感悟  
細節 激勵 心理學家  
專家 溝通 大師 摆脫 潛意識  
魅力 圖解 任何人 天全 讓人  
教授 決斷 記憶 催眠 心術  
聊天 聖經 搭訕 把妹  
看穿 穿著 行為 史上 最強  
表露 逆境 語法 人生 情場  
最強 NLP 正妹 幽默  
圖法 人生 說服 動人心  
厚黑學 心理學 人心  
開發 誘導 思考 人心  
入門 謂點 厚黑  
必勝 問題 毒舌 做得到  
才是 超強

把妹達人  
正妹心理學  
心理學家的專業把妹術  
搭訕聖經  
正妹沒告訴你的事

貼心的女人，幸福無敵：改變男人的賀爾蒙，  
從貼心做起



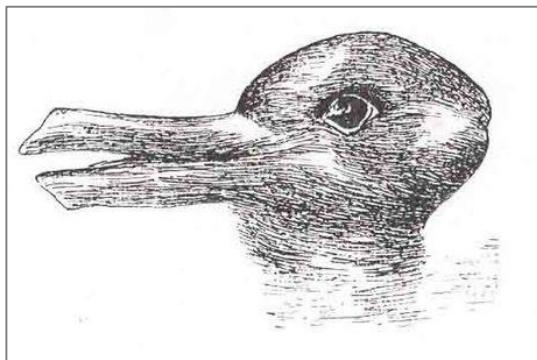
# 各種行為觀察與實驗

Do you find this smile to read?  
Because of the phenomenal power of the human mind, most people do.

你是是不是該嘲笑這個人？這個人的  
mind的哼哼唧唧唧唧唧唧唧唧唧唧  
的。

Nǐ shì bùshì gāi cháoxiào zhège rén? Zhège rén de mnid de hēng  
hēng ji jījījījījī jī jī jī de.

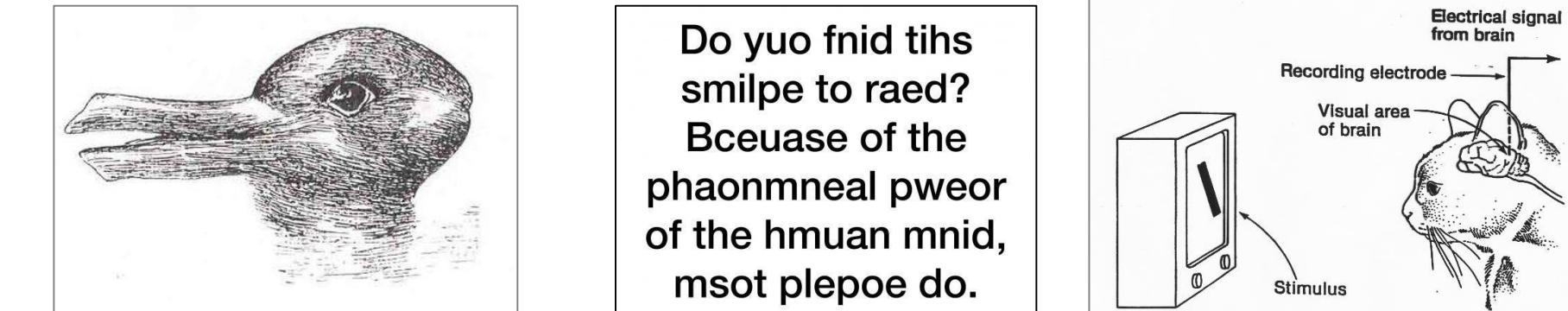
C 你是不是要查： Do you find this smile to *read*?  
*Because of the phenomenal power of the human mind,*  
*most people do.*



Shifting Gestalt  
Jastow, 1899

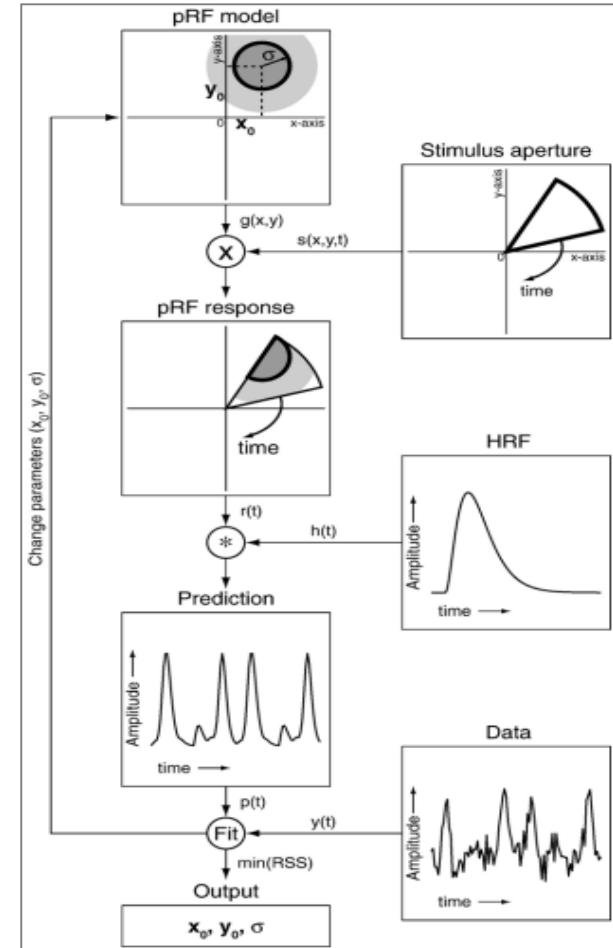
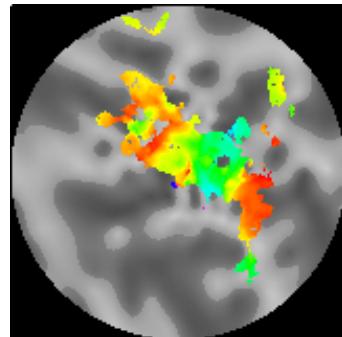
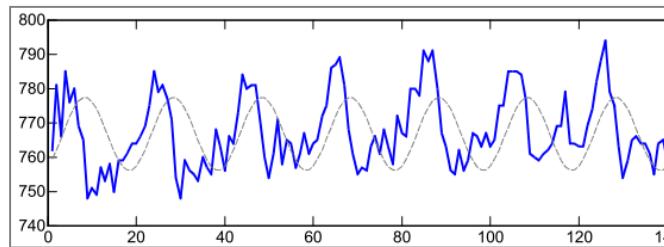
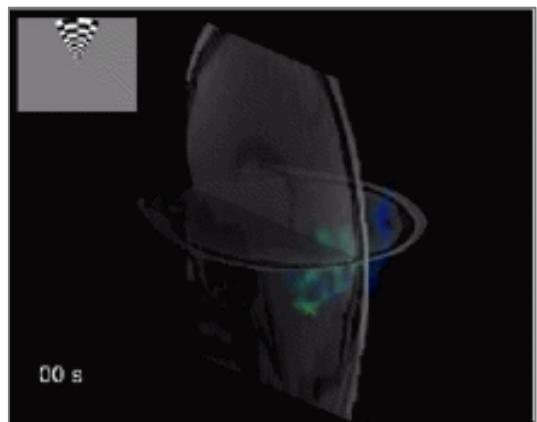
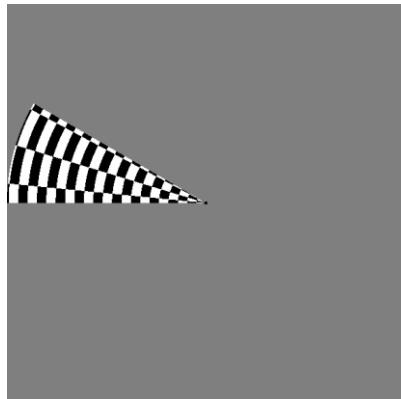
Do you find this  
smile to read?  
Because of the  
phenomenal power  
of the human mind,  
most people do.

Jumbled letter  
Warrington, 1980



Receptive fields  
Hubel & Wiesel, 1959

# 實際上，我過去的實驗室日常 ...



# 為何來到這裡？



**Bridge the Gap between Human Brain and Artificial Intelligence**

# 另一個故事

---

# 過去在做甚麼

- 背景

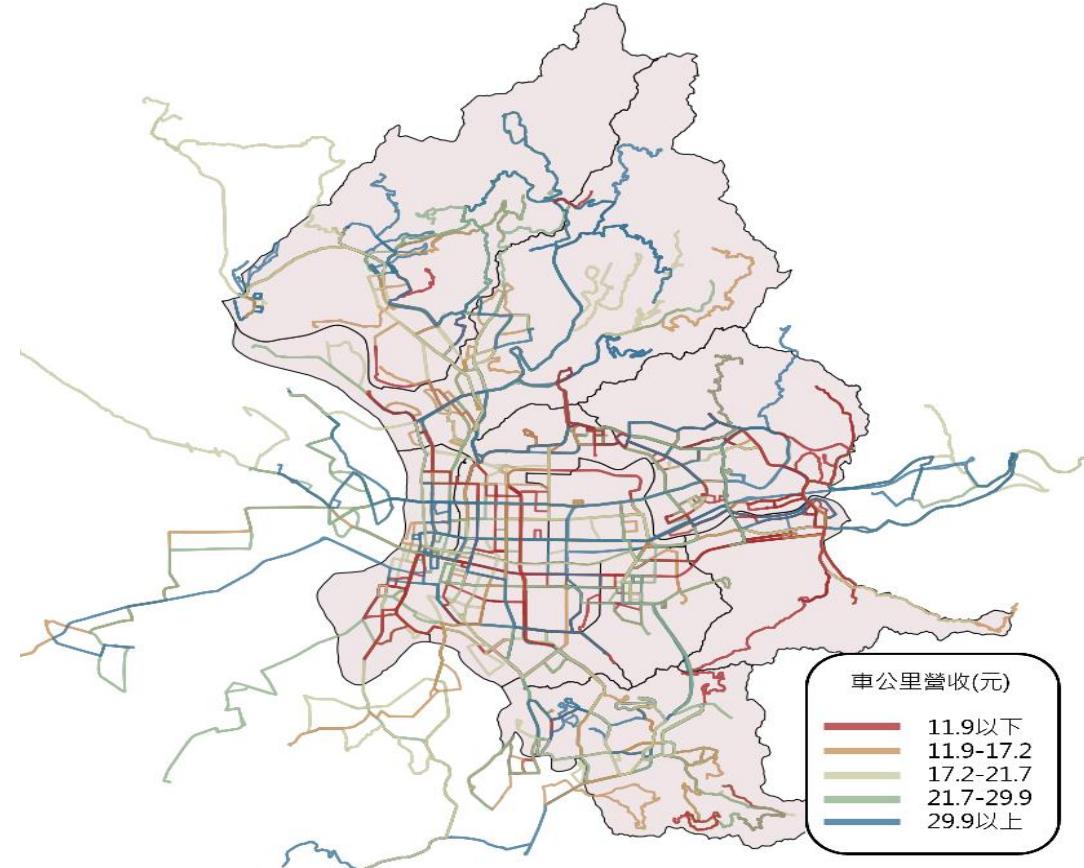


交大運管系



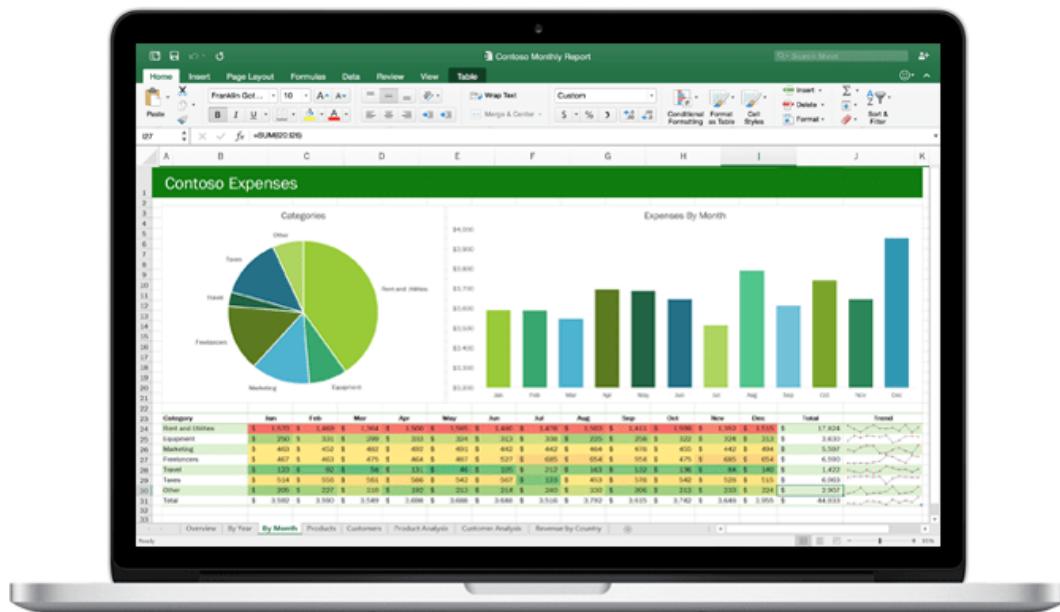
台大土木所

- 台北市虧損公車路線營運分析
  - 真的虧錢？
  - 為什麼虧錢？



# 技術能力

- Coding 經驗
  - 大一計算機概論 → 6 學分
- 資料分析技術
  - 基礎統計分析
  - ~~Machine learning~~
  - ~~Deep learning~~



# FAQ

---

Why

為什麼想要轉換跑道？

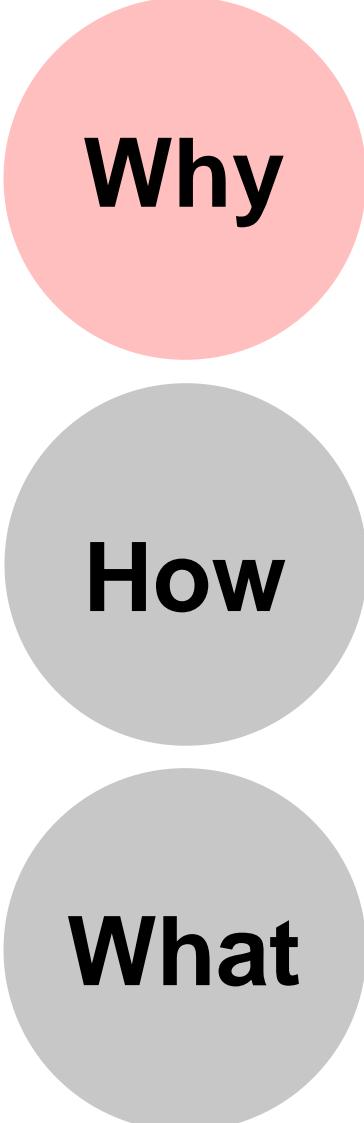
How

非本科系的要怎麼做？

What

想跨入這領域該做甚麼？

---



**Why**

為什麼想要轉換跑道？

**How**

非本科系的要怎麼做？

**What**

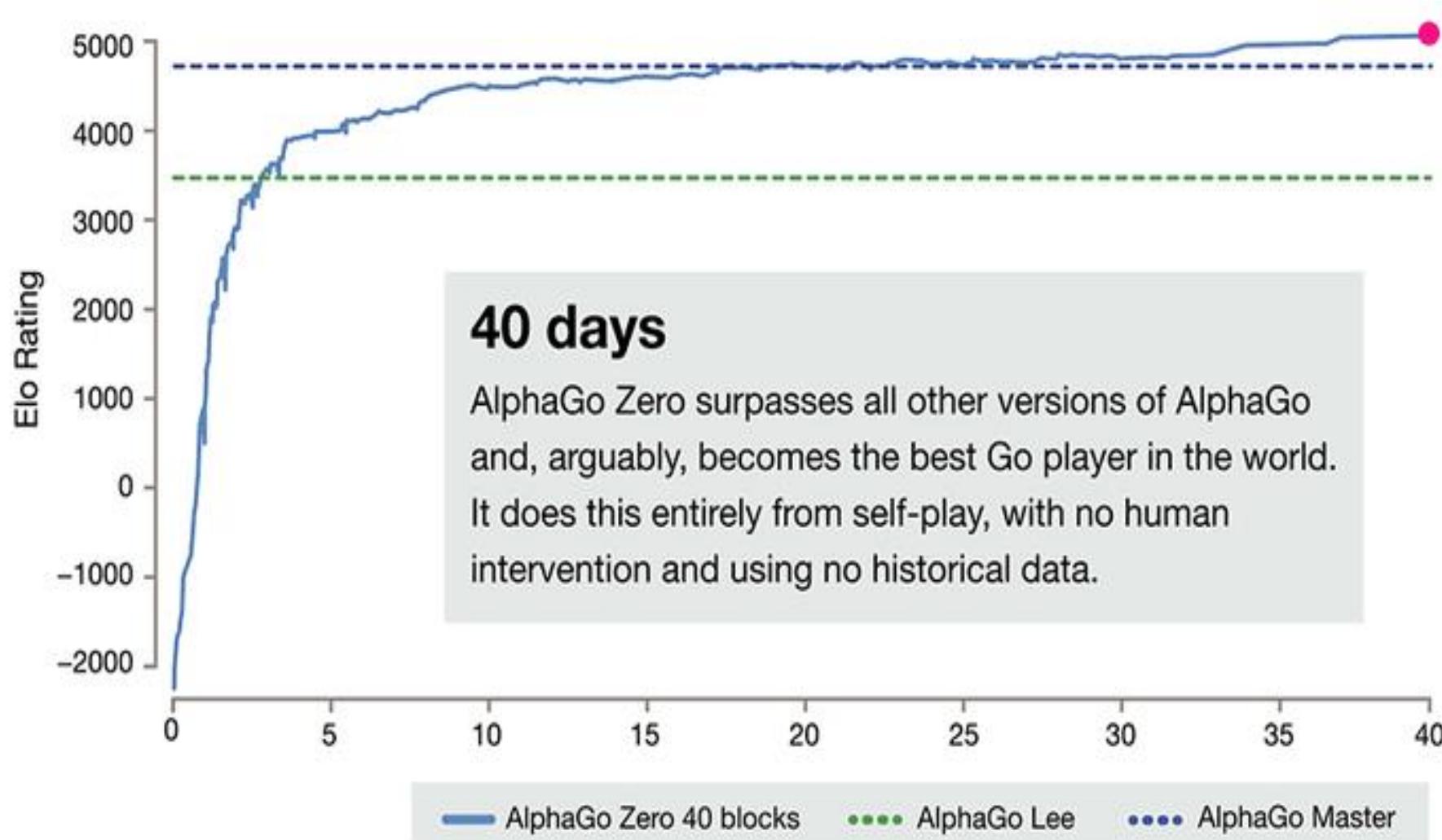
想跨入這領域該做甚麼？

# why data science?

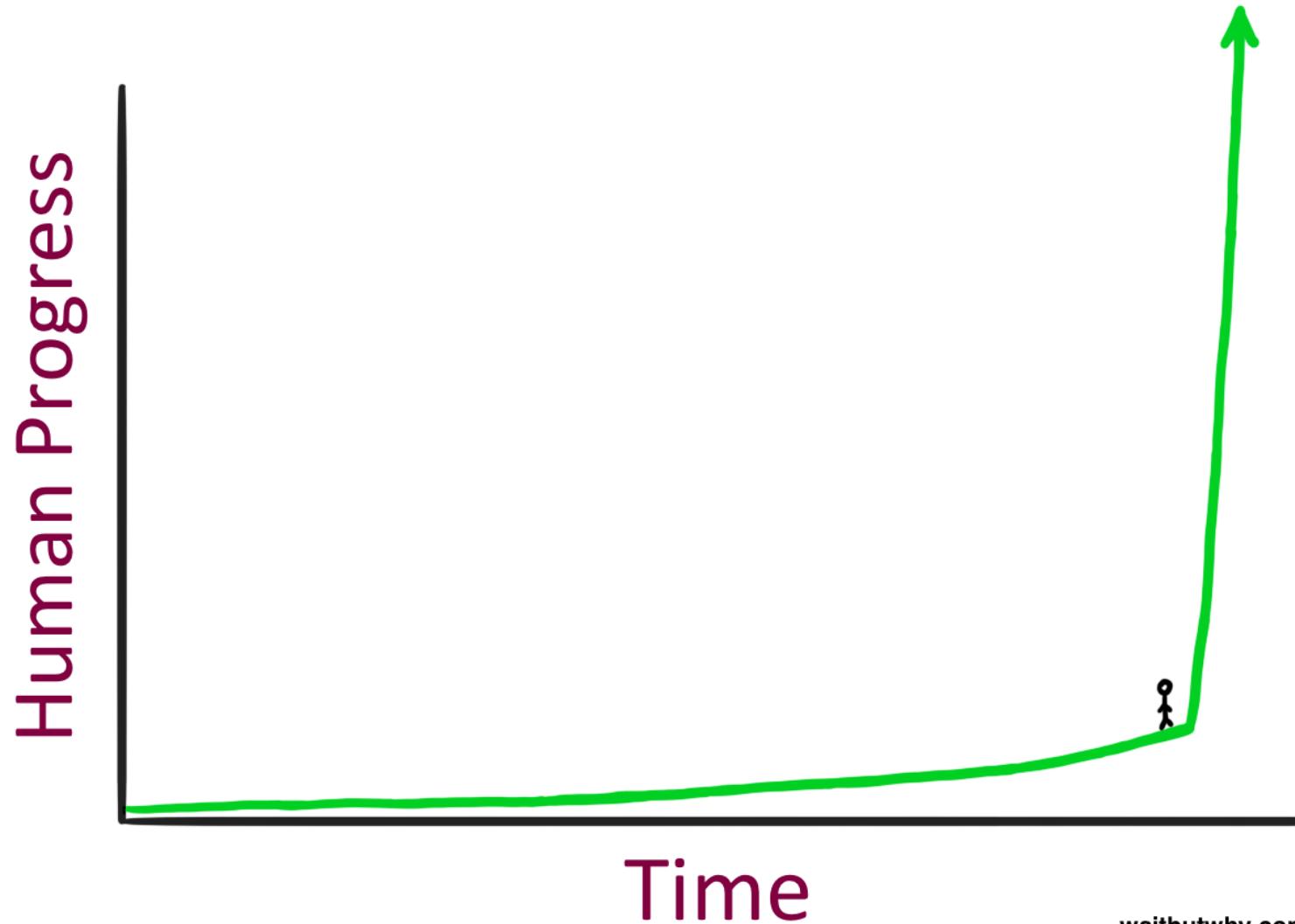


- 卷積神經網路 (Convolutional Neural Network, CNN) ?
- 深度強化學習 (Deep Q-Learning, DQN) ?
- 蒙地卡羅樹狀搜尋演算法 (Monte Carlo Tree Search) ?

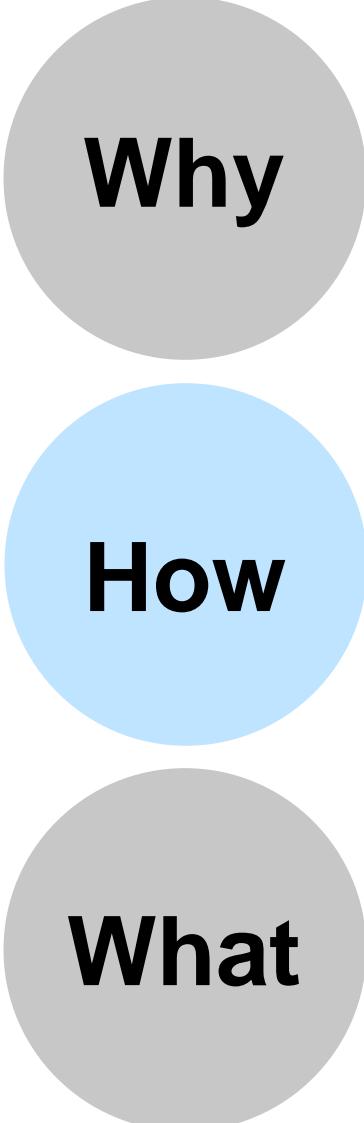
# Now? AlphaGo Zero!



# 科技進展



---



**Why**

為什麼想要轉換跑道？

**How**

非本科系的要怎麼做？

**What**

想跨入這領域該做甚麼？

# MOOC!

---

Code



DATAQUEST

codecademy



DataCamp

Project

coursera

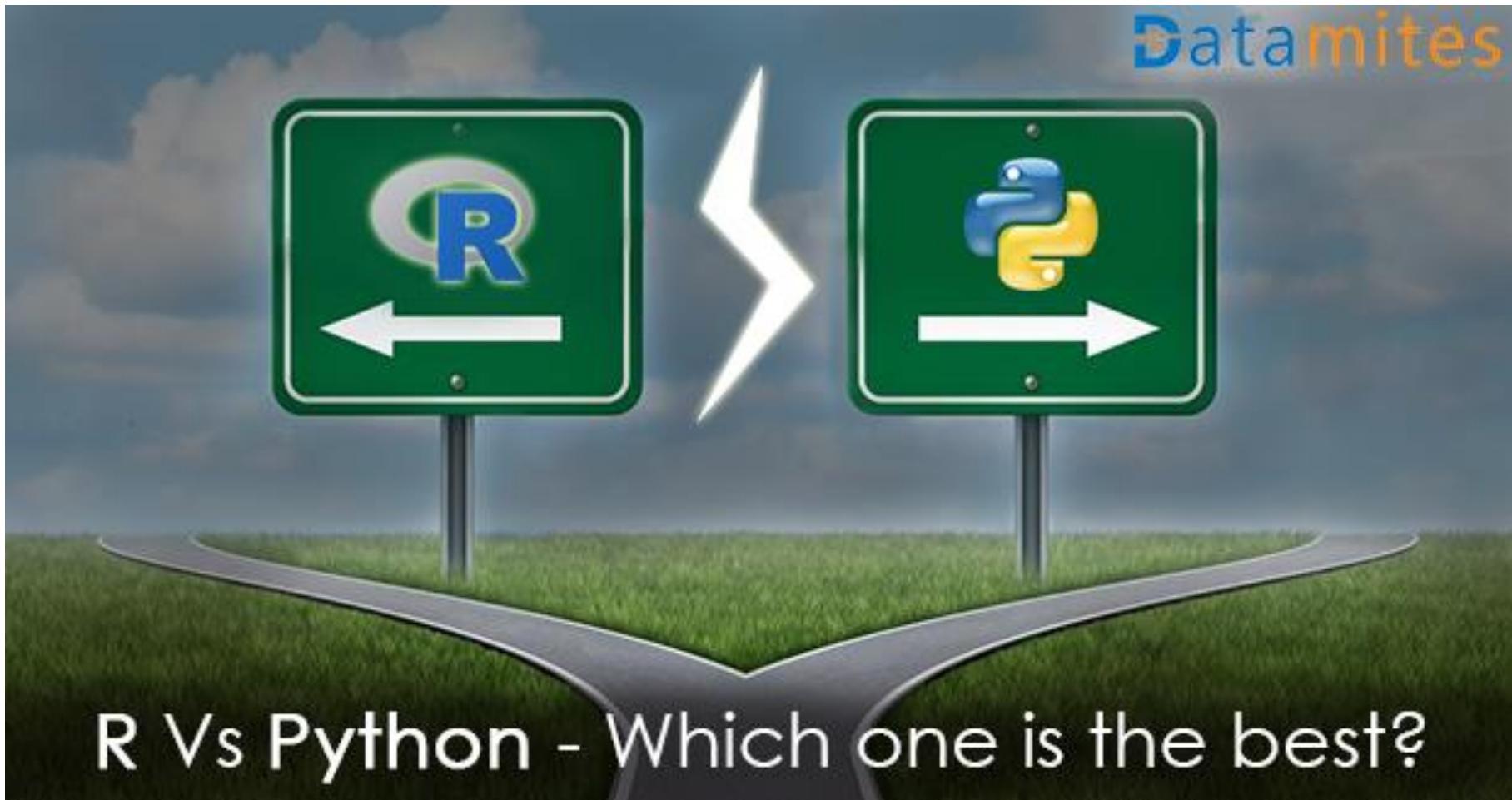


UDACITY

edX

# 該選擇甚麼語言？

- 世紀之爭: R vs. Python



# R 的優點



- 資料整理與清洗非常簡潔



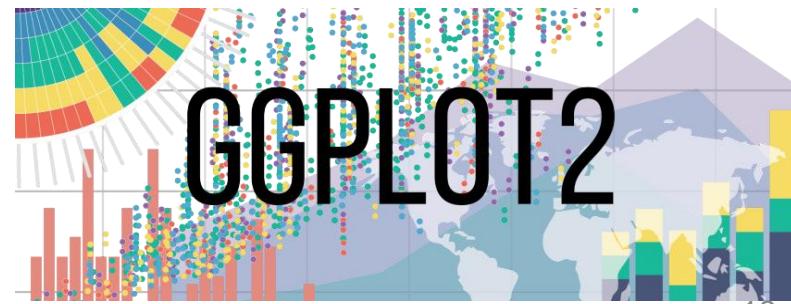
- 快速的統計建模

## Linear Regression in R!!!

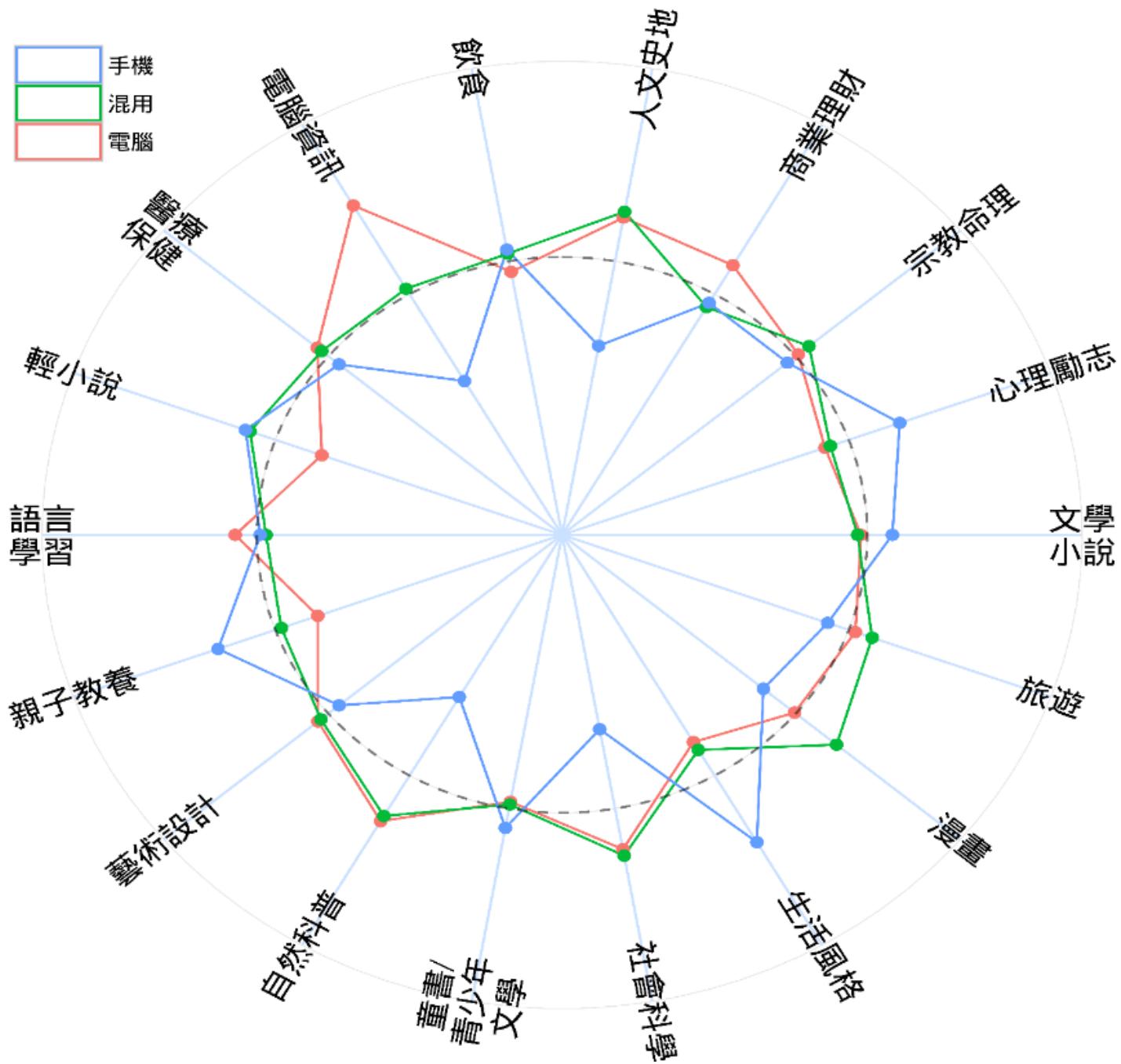
```
Call:  
lm(formula = size ~ weight, data = mouse.data)  
  
Residuals:  
    Min      1Q  Median      3Q     Max  
-1.5482 -0.8037  0.1186  0.6186  1.8852  
  
Coefficients:  
              Estimate Std. Error t value Pr(>|t|)  
(Intercept)  0.5813    0.9647   0.603   0.5658  
weight       0.7778    0.2334   3.332   0.0126 *
```

So Easy!!!!

- 強大的 ggplot2

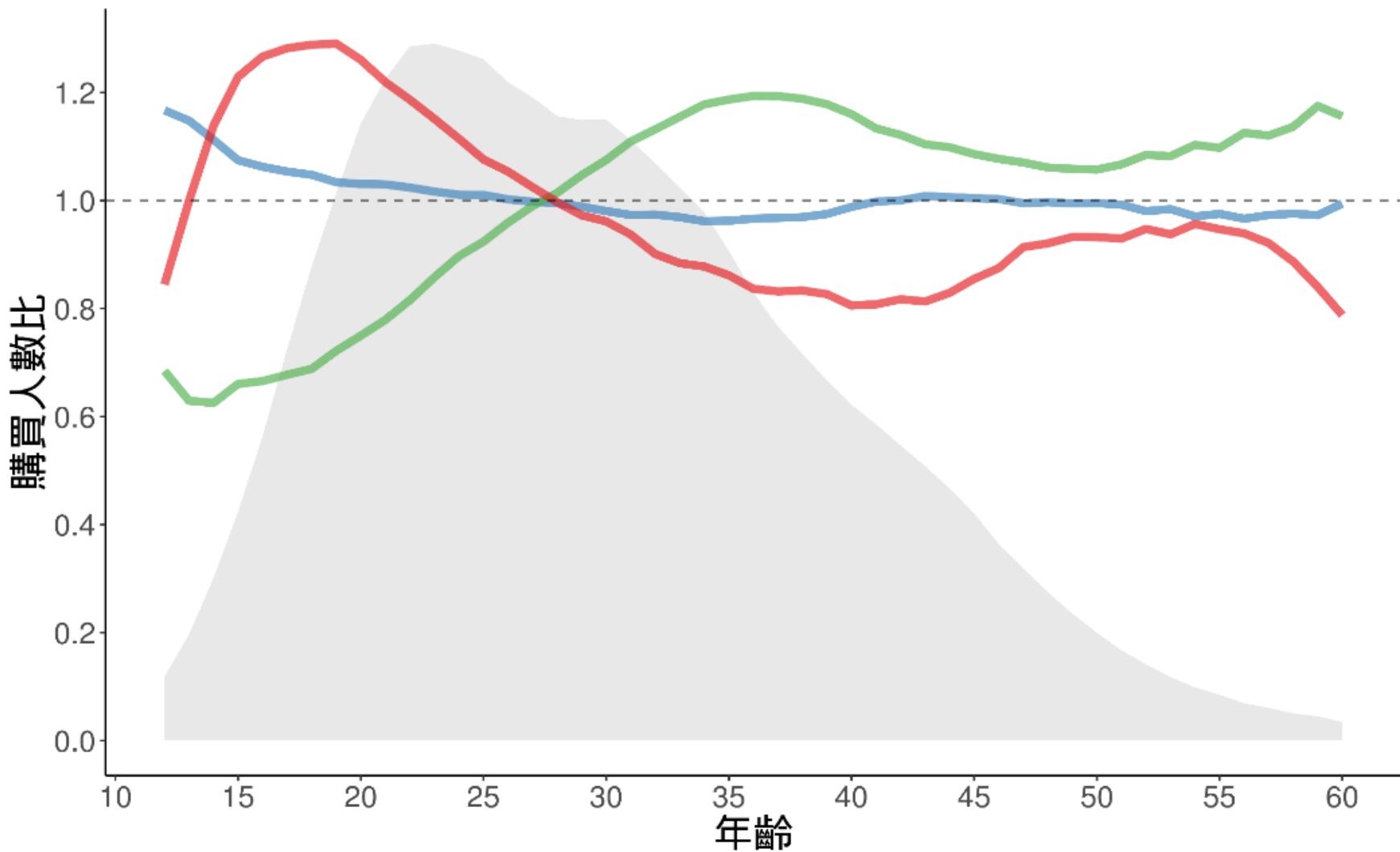


手機  
混用  
電腦



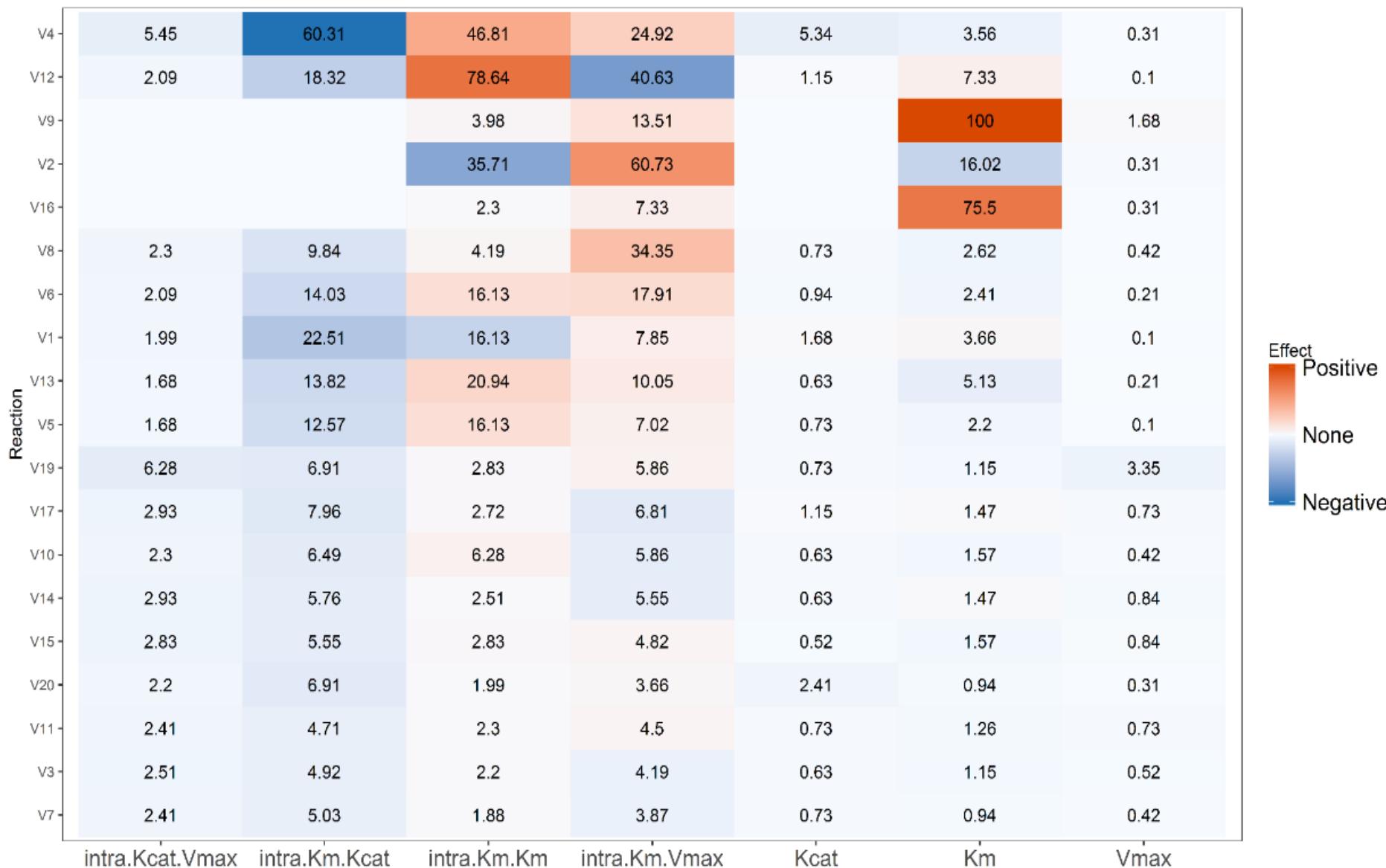
# 商業理財 - 成功法 (年齡)

■ 自我成長 ■ 致富 ■ 生涯規劃



百搭造型初學縫紉彩繪天然  
手帳禮物美日記甜物超可愛手工皂  
瘦腿身纏創意毛氈生活雜貨  
一個曲線簡單編織巴黎針法  
串珠上愛繩體我美立  
吊飾鉤針樣織貝光設計  
優雅鉤織黏士女孩療  
美甲飾品幸法式實用  
彩妝拼布圖案毛線  
妝包點ok玩偶韓國絲  
提袋娃娃刺繡拼布王學者  
提動物時尚人氣手縫卡片  
最愛時尚OK學會布包開心  
時光公斤包包手感

Feature importance in each reaction



# Python 的優點



- 許多套件支援連接資料庫
- 成熟的網路爬蟲套件
- 深度學習最常用的語言



# 如何選擇?

---

- 常需做探索式資料分析 (EDA)、資料視覺化與統計建模



- 網路爬蟲、串接資料庫或是做 deep learning



# 常被問的問題

---

**Why**

為什麼想要轉換跑道？

**How**

非本科系的要怎麼做？

**What**

想跨入這領域該做甚麼？

# 想跨入這領域該做甚麼？

- 做 Project!



- Scikit-learn toy datasets

鳶尾花資料集  
IRIS dataset

鐵達尼號資料集  
titanic dataset

波士頓房價資料集  
BOSTON dataset

手寫數字資料集  
MNIST dataset

嫌資料集太簡單？



- Kaggle 是一個 data science 競賽平台。企業和研究者可在其上發布資料集，讓全球專家可在其上進行競賽以產生最好的模型
  - Heritage Health Prize，獎金高達 300 萬美金
- **11/12(日) 11:20-12:05 我在 Kaggle 數海獅**
  - 使用深度學習方法協助 NOAA 生物學家監控阿留申群島的海獅數量
  - 勇奪世界第一！



# My project in 資料洞察實驗室

---



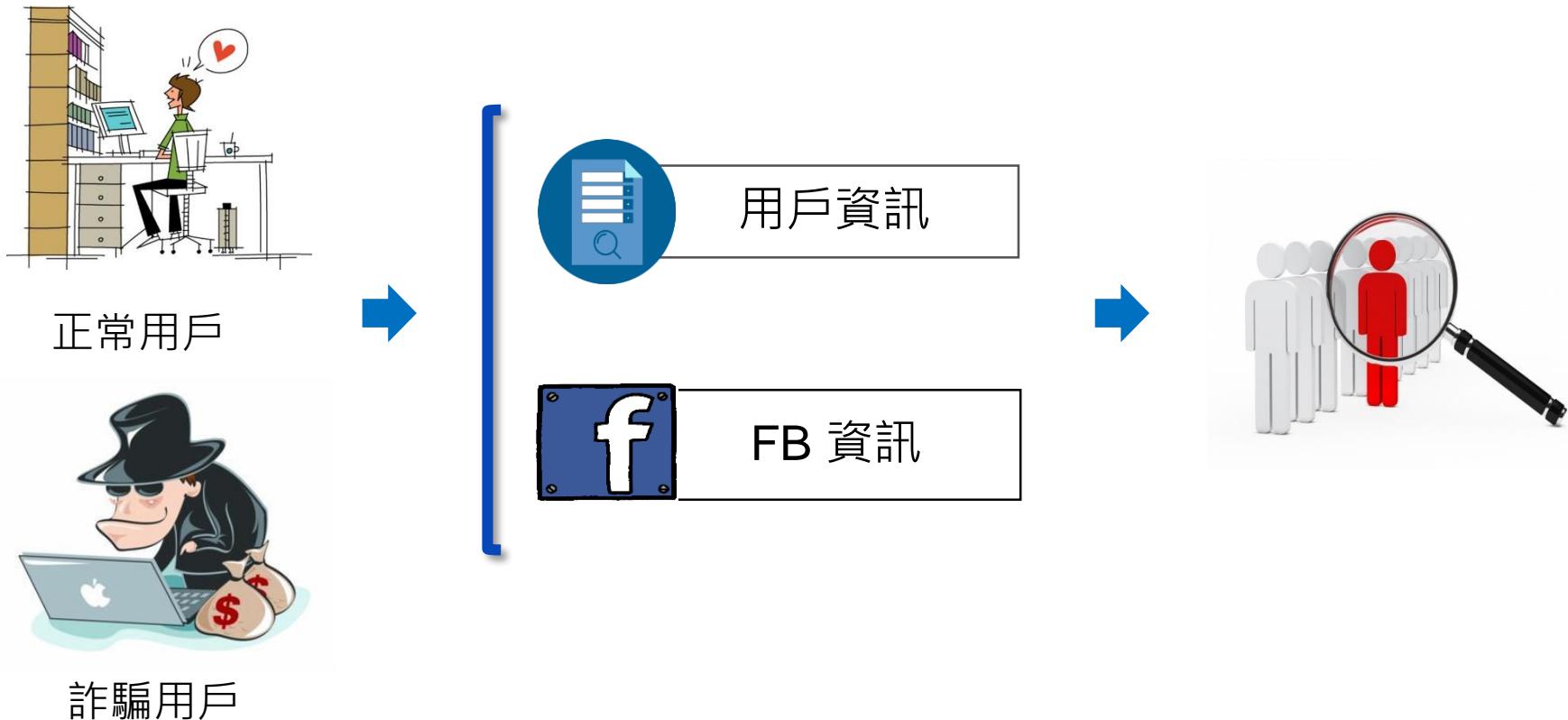
# 線上交友詐騙帳號分析

- 線上交友平台充斥詐騙帳號，如何透過機器學習技術，揪出這些斂財騙子，拯救悲情男女的荷包？



# 建立 features

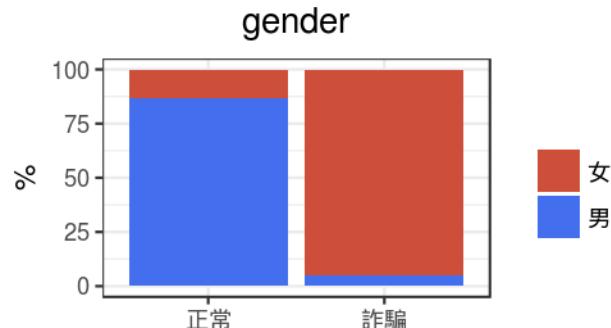
- 觀察詐騙用戶與正常用戶在我們建立的 features 上是否會有所不同。



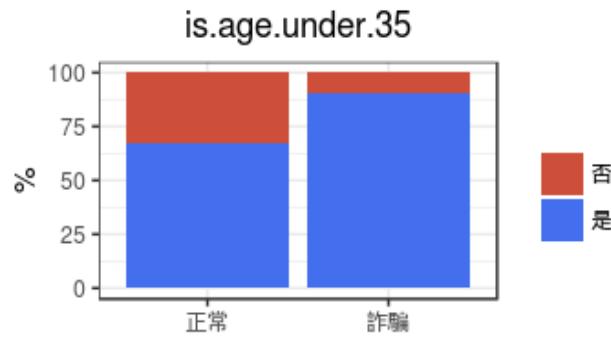
# 詐騙用戶與正常用戶之差異 - 用戶資訊



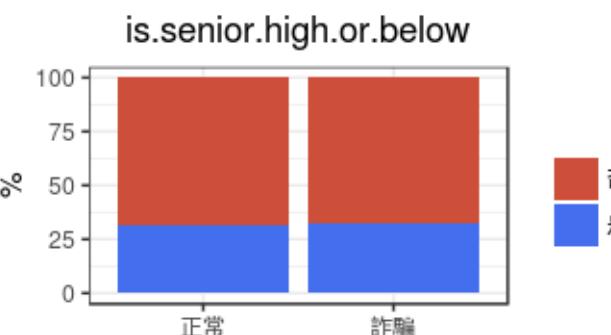
性別



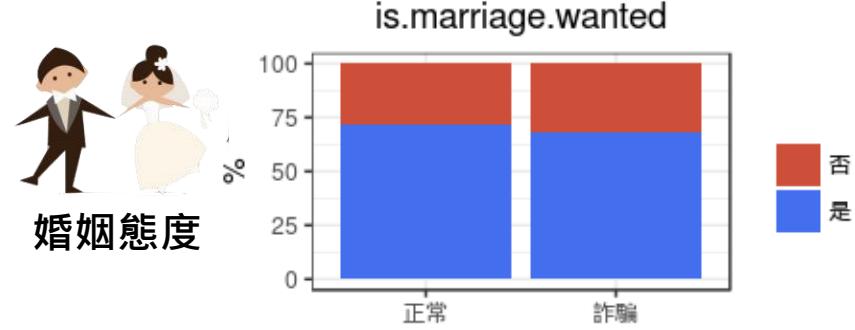
年齡



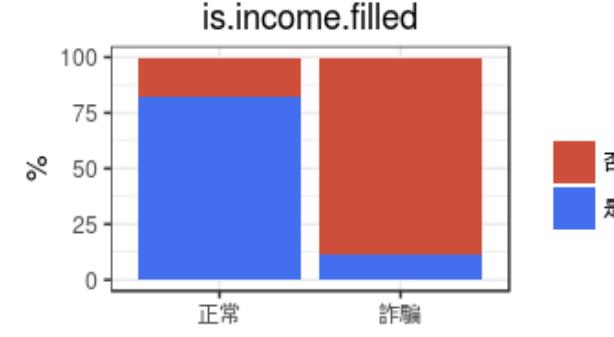
學歷



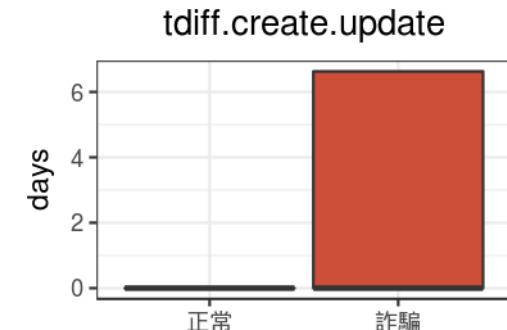
婚姻態度



年收入



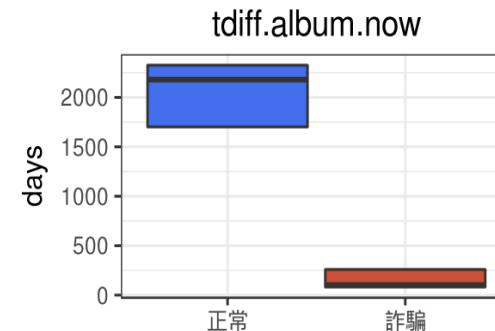
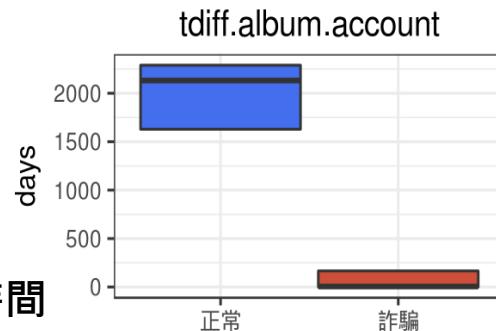
帳號資訊



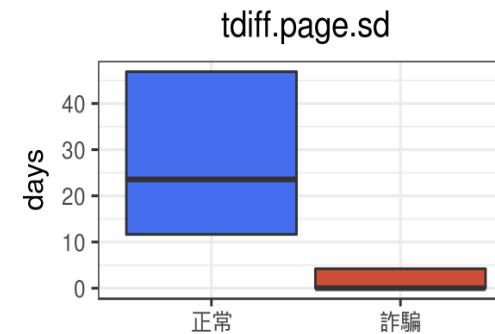
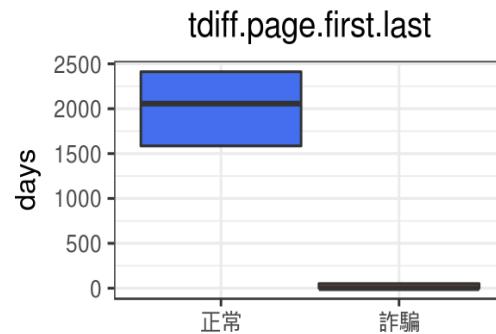
# 詐騙用戶與正常用戶之差異 – FB 資訊



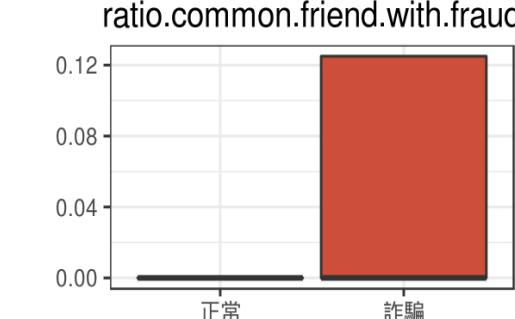
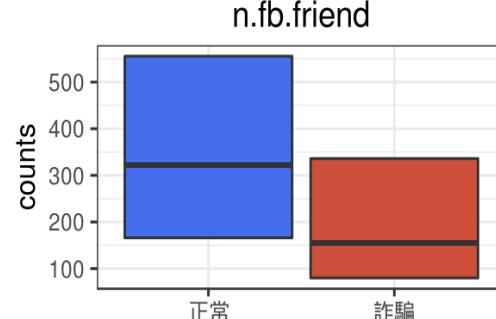
大頭貼建立時間



粉絲團

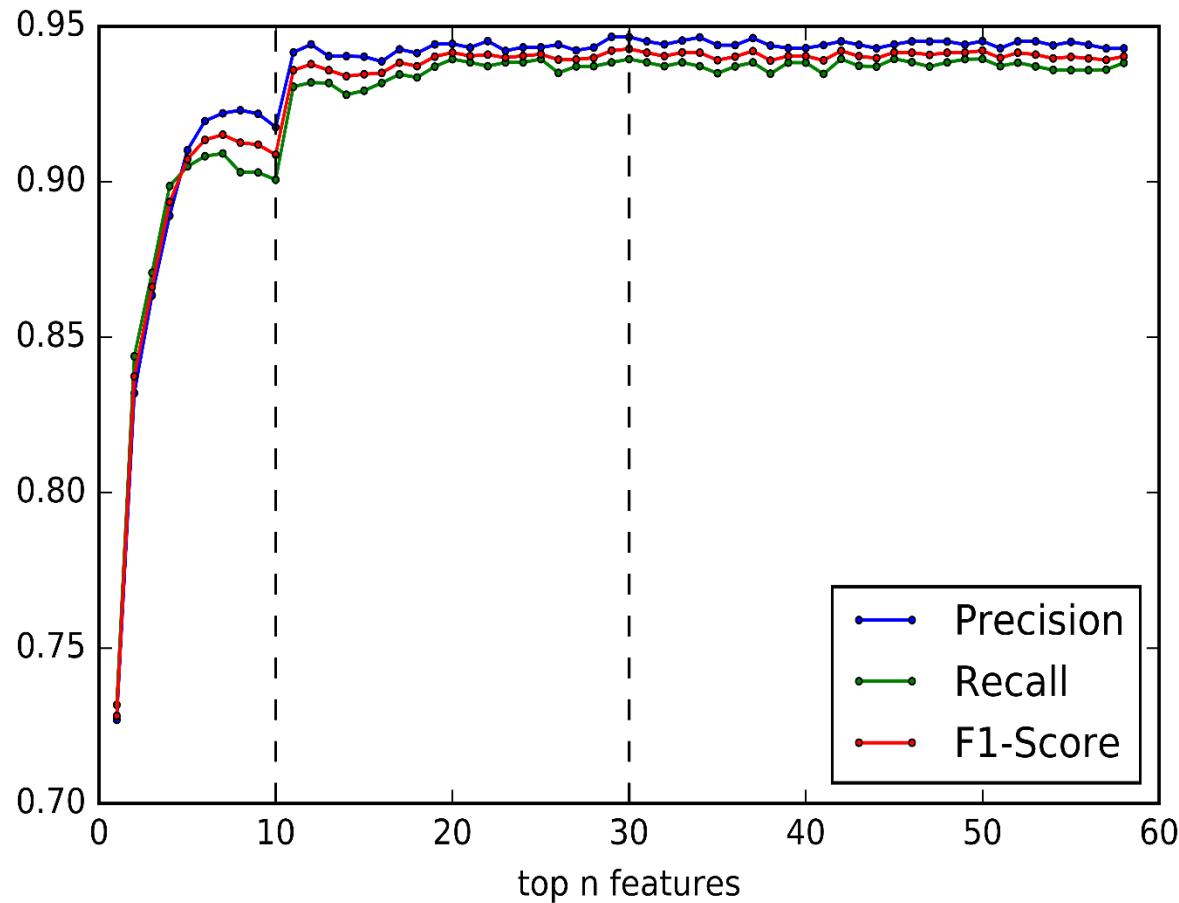


好友

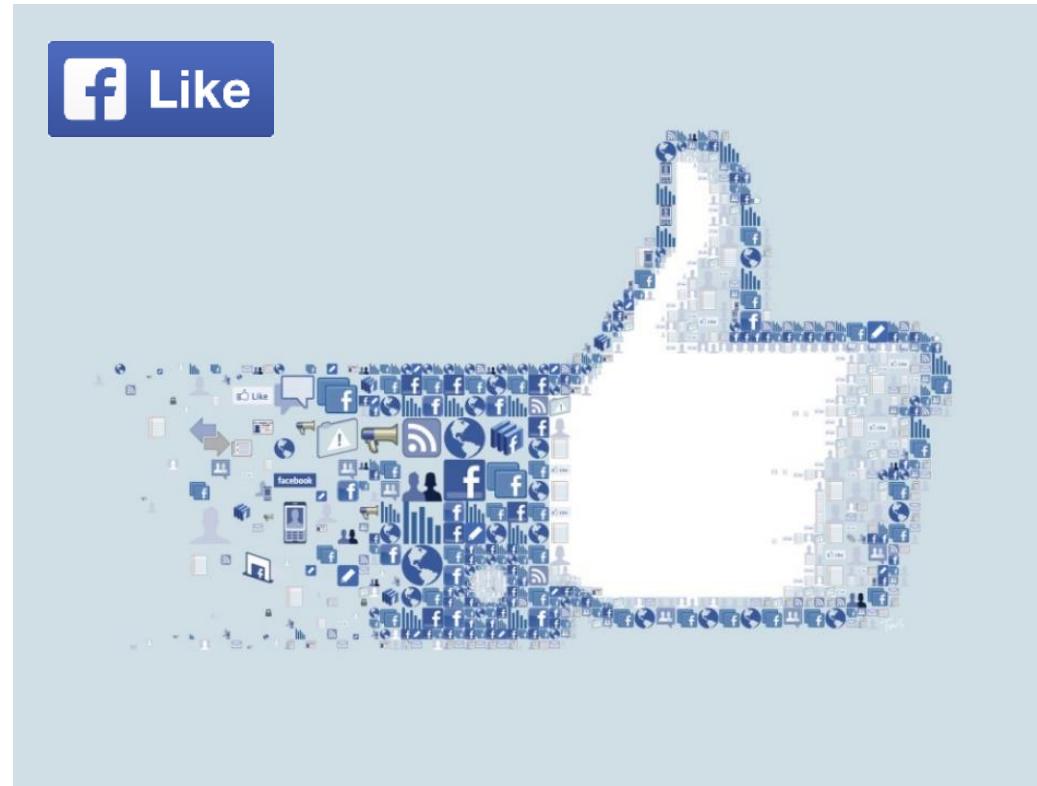
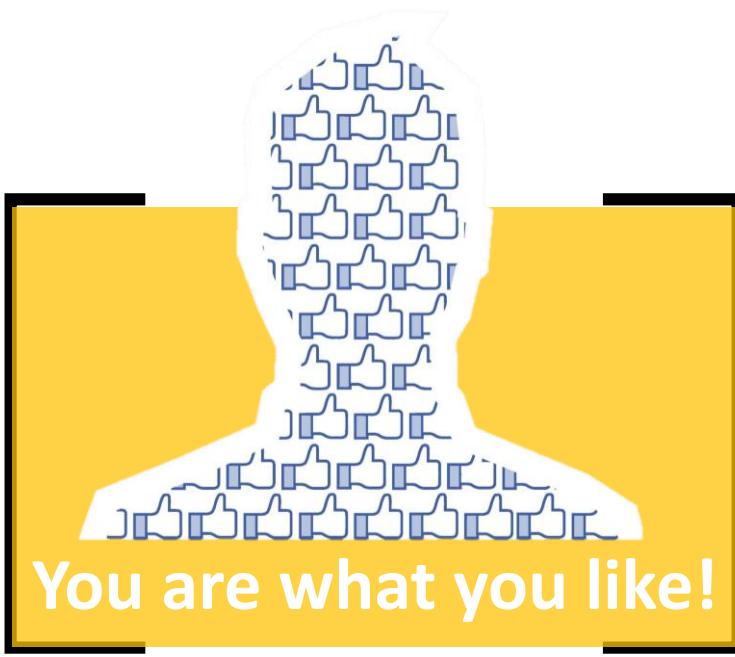


# modeling

- 用有效的 features · f1-score 達到 0.94

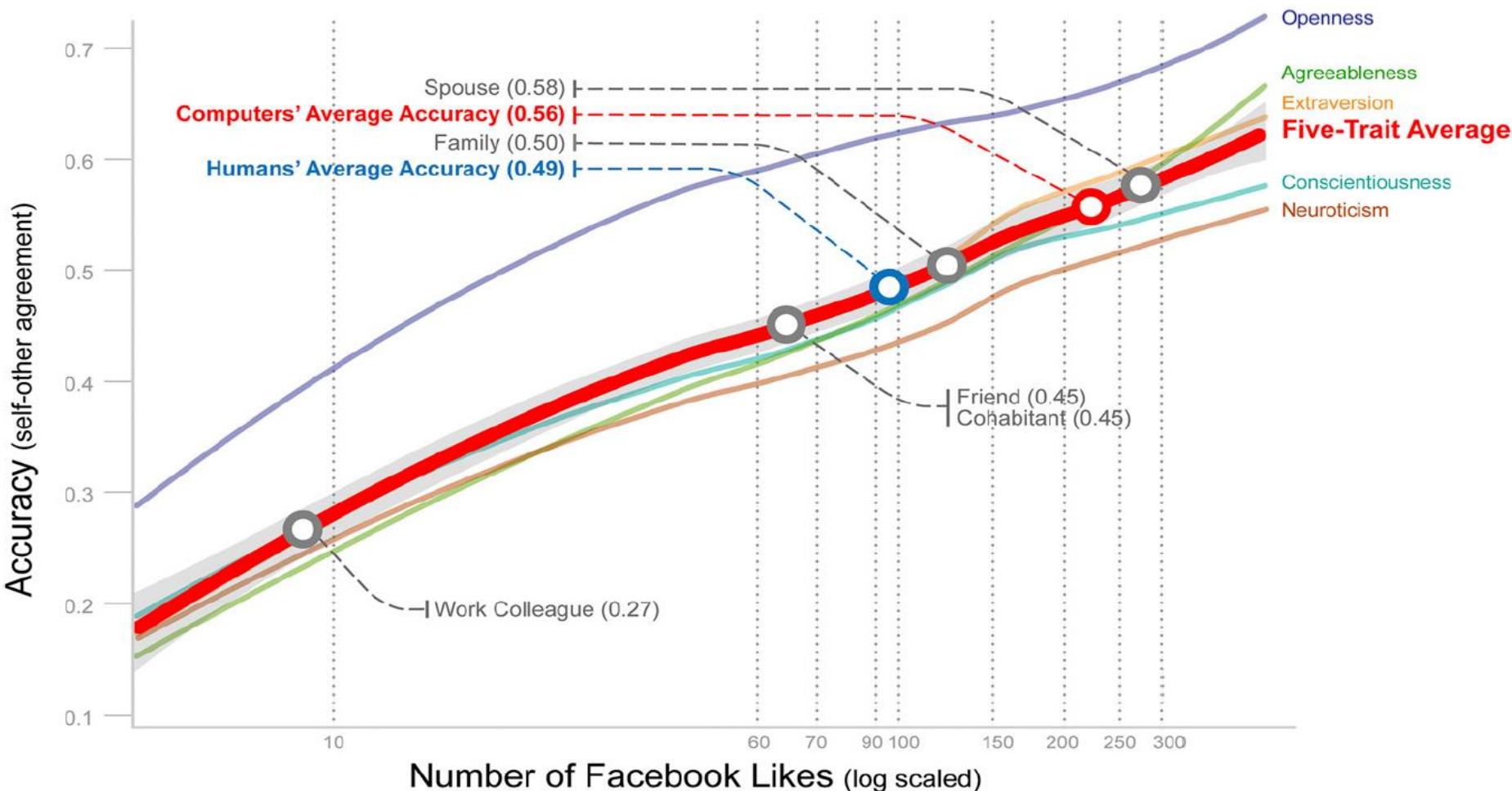


Random Forest	
Precision	0.947
Recall	0.940
F1-Score	0.943



# Computer vs. Humans

- 用你按過讚的粉絲團來猜你的個性



Computer-based personality judgments are more accurate than those made by humans.  
Proceedings of the National Academy of Sciences

# 誰按的粉絲團？

??

[小茉莉-陳瑀希](#)

[Garena 《英雄聯盟 LOL》](#)

[波多野結衣Hatano Yui](#)

[解婕翎](#)

[豆花妹 蔡黃汝](#)

[Nono\\_辜莞允](#)

[張景嵐](#)

[FHM Taiwan 男人幫國際中文版](#)

[潮物blog - 街頭潮流男著](#)

[張小筑 Ya Chu](#)

??

[Catworld小舖](#)

[EYESCREAM Inc.](#)

[LOVFEEL](#)

[OB嚴選](#)

[grace gift](#)

[BEVY C.](#)

[Joyceshopstyle](#)

[QUEEN FASHION SHOP](#)

[SweeSa水莎](#)

[Lulus](#)

# 性別

男

[小茉莉-陳瑀希](#)

[Garena 《英雄聯盟 LOL》](#)

[波多野結衣Hatano Yui](#)

[解婕翎](#)

[豆花妹 蔡黃汝](#)

[Nono\\_辜莞允](#)

[張景嵐](#)

[FHM Taiwan 男人幫國際中文版](#)

[潮物blog - 街頭潮流男著](#)

[張小筑 Ya Chu](#)

女

[Catworld小舖](#)

[EYESCREAM Inc.](#)

[LOVFEEL](#)

[OB嚴選](#)

[grace gift](#)

[BEVY C.](#)

[Joyceshopstyle](#)

[QUEEN FASHION SHOP](#)

[SweeSa水莎](#)

[Lulus](#)

# 誰按的粉絲團？

??

一休陪你一起愛瘦身

iFit 愛瘦身

小甜甜 張可昀

FB減肥達人 輕鬆教你瘦

美樂蒂 Melody

BEMAX

Woma

OB嚴選

鍾欣凌

杜詩梅 Tu Shih Mei

??

《HITO 本舖》

潮物部落格

PAZZO

Image

《OneBoy》

UNO STORE

RockSteady

SweeSa水莎

高高-流行服飾Store

Maxy

# 體型

胖

一休陪你一起愛瘦身

iFit 愛瘦身

小甜甜 張可昀

FB減肥達人 輕鬆教你瘦

美樂蒂 Melody

BEMAX

Woma

OB嚴選

鍾欣凌

杜詩梅 Tu Shih Mei

瘦

《HITO 本舖》

潮物部落格

PAZZO

Image

《OneBoy》

UNO STORE

RockSteady

SweeSa水莎

高高-流行服飾Store

Maxy

# 誰按的粉絲團？

高

TED

Technews 科技新報

批踢踢實業坊 (Ptt.cc)

PanSci 科學新聞網

商業周刊 (商周.com)

TEDxTaipei

The News Lens 關鍵評論網

背包客棧

VoiceTube 看影片學英語

Cheers : 快樂工作人

低

在不瘋狂就等死

羅志祥 SHOW

Mimi Dancing Club

小A辣

小三魔法棒 (小三美日)

真愛談戀愛 × 真愛橋

宛宛兒

連靜雯joanne lien

BY2

爆笑禁區

# 學歷

高

TED

Technews 科技新報

批踢踢實業坊 (Ptt.cc)

PanSci 科學新聞網

商業周刊 (商周.com)

TEDxTaipei

The News Lens 關鍵評論網

背包客棧

VoiceTube 看影片學英語

Cheers : 快樂工作人

低

在不瘋狂就等死

羅志祥 SHOW

Mimi Dancing Club

小A辣

小三魔法棒 (小三美日)

真愛談戀愛 × 真愛橋

宛宛兒

連靜雯joanne lien

BY2

爆笑禁區



研之有物

本站盼以具體的研究案例、真實的研究員生活，  
帶您前往數理科學、生命科學、人文社會三大  
領域研究現場，揭開中央研究院神秘的面紗

研之有物  
@research.sinica

研之有物 | 中央研究院



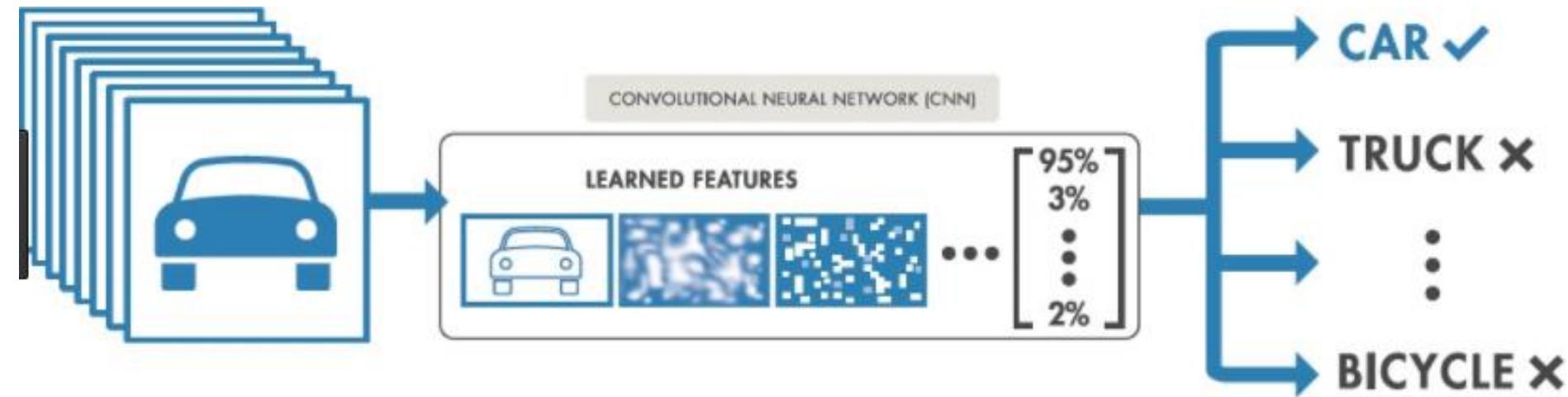
讚

追蹤

分享

...

來去逛逛

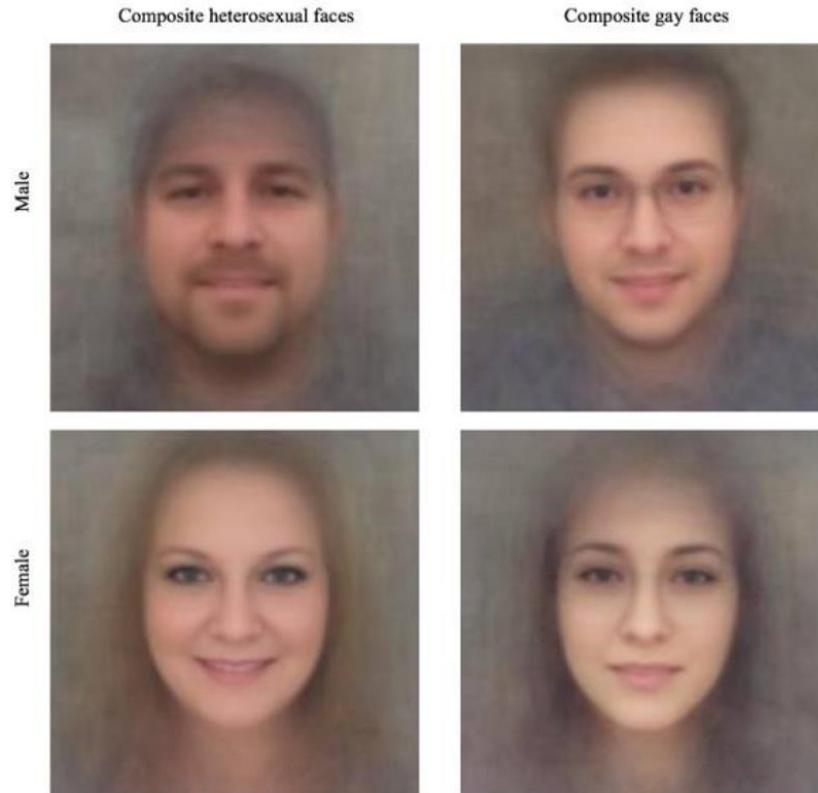


# 卷積神經網路

# Convolutional Neural Network

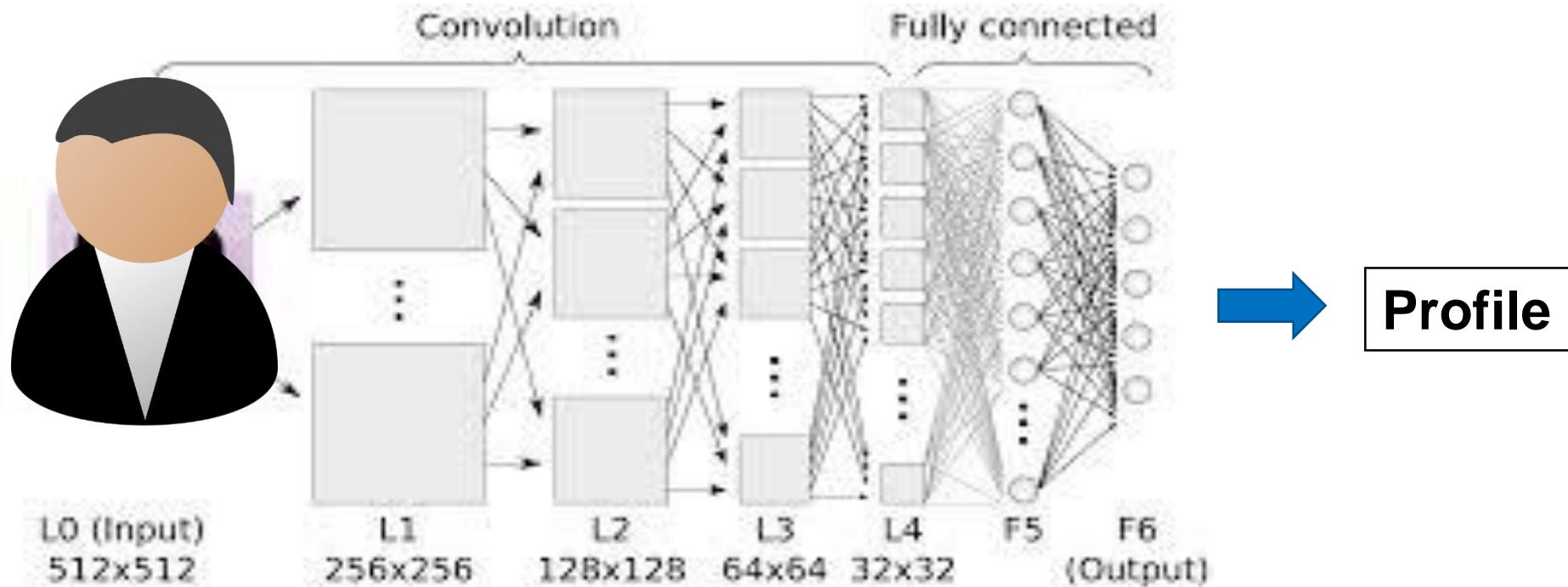
# AI 判斷同性戀？

- 史丹佛大學透過美國約會網站 3 萬 5000 多名男女會員的臉部照片為樣本，分辨出同男和異男的準確度可達 81%



# 用照片預測基本資料

- 運用 CNN 提取圖片特徵，進行分類



# 預測結果

- 純粹用照片來預測



會不會被喜歡 (女生) 0.9

會不會被喜歡 (男生) 0.78

FB 朋友數 0.67

是否抽菸喝酒 0.65

# 預測誰是基督徒?

---



基督教徒的機率: 0.9999

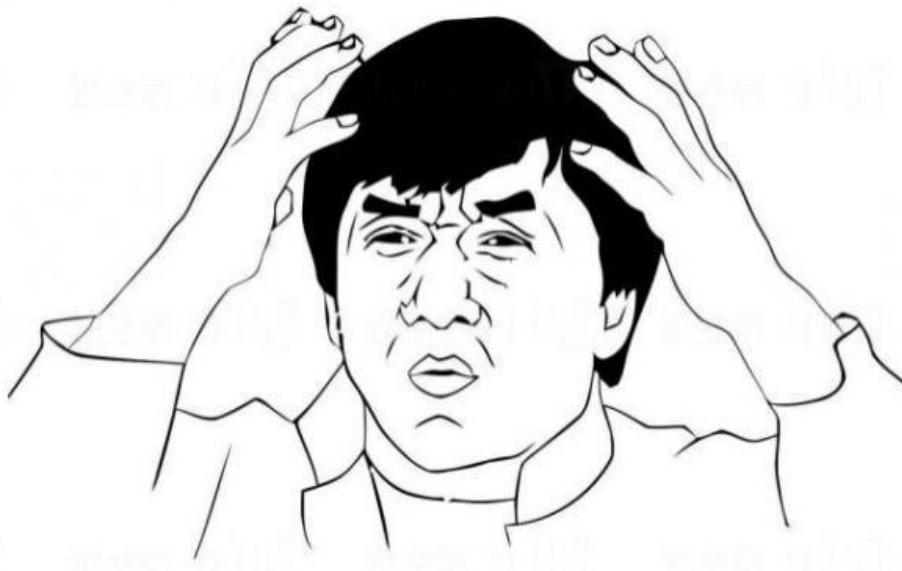


基督教徒的機率: 0

# 三冠王模型的故事

---

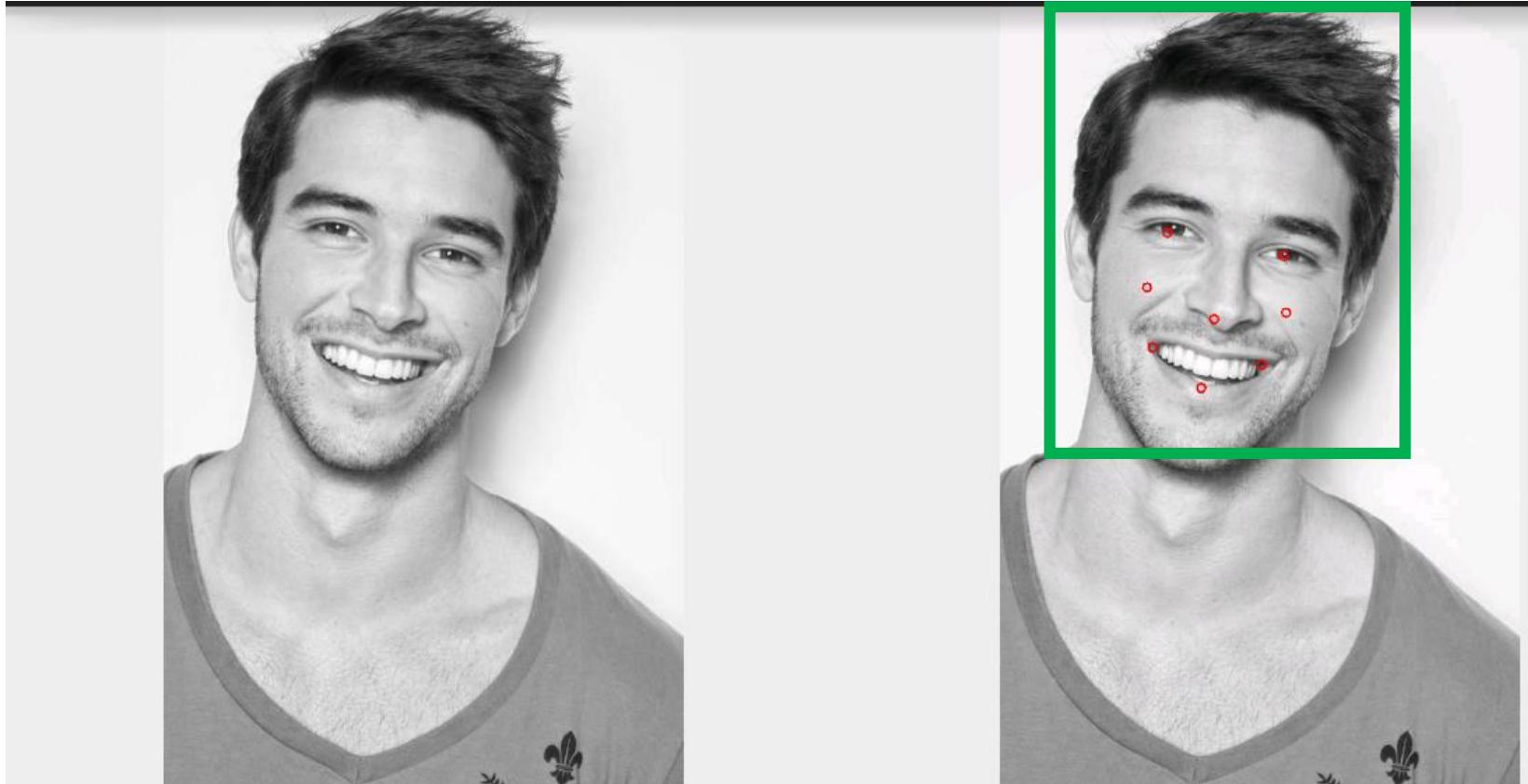
- 預測是否為基督徒
  - Accuracy: 1.00
  - AUC: 1.00
  - F1-score: 1.00
- 確定是在 testing data 上面的表現



# 三冠王模型的故事

---

- 資料前處理 → 用 dlib 確保每張圖片都有包含臉，且全臉佔圖的比例不能太小



# 三冠王模型的故事

- 忘記把所有基督徒照片的綠框刪除



# 資料 Recheck

---



基督徒



非基督徒

- 我訓練好的模型不是一個基督徒分類器
- 而是一個**“綠框”**分類器!

# 預測誰是基督徒?

---



基督教徒的機率: 0.9999



基督教徒的機率: 0

# ~~預測誰是基督徒？哪張照片有綠框？~~

---



~~基督教徒的機率: 0.9999~~

有綠框的機率: 0.9999



~~基督教徒的機率: 0~~

有綠框的機率: 0

# 深度學習在現實世界應用中的 各種趣味陷阱

---

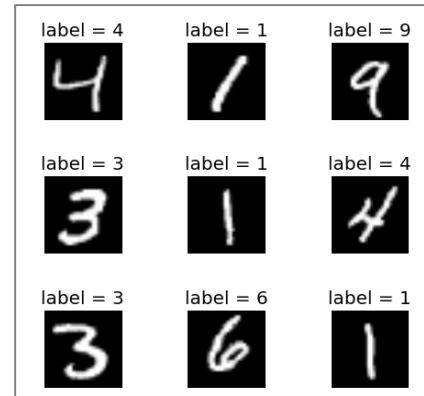
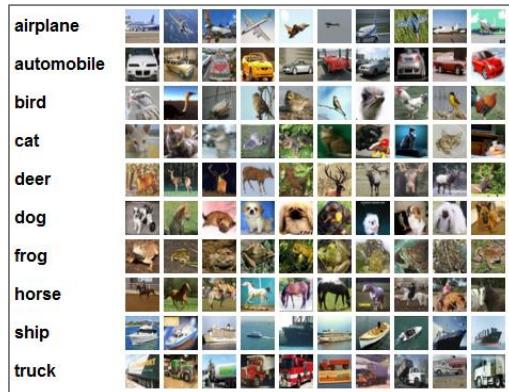
# 深度學習的美好想像與迷思

---

- 萬物皆可 Train
  - FizzBuzz 的故事

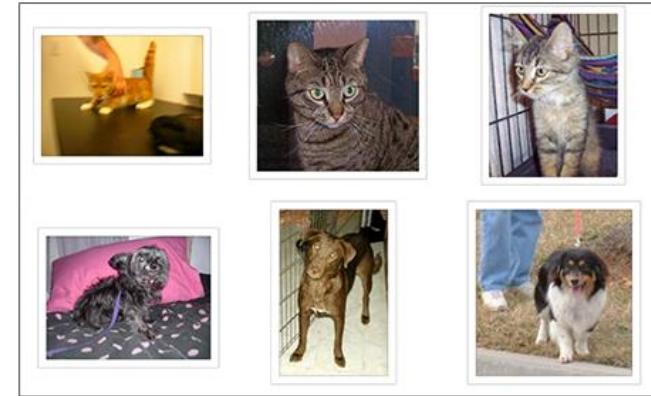
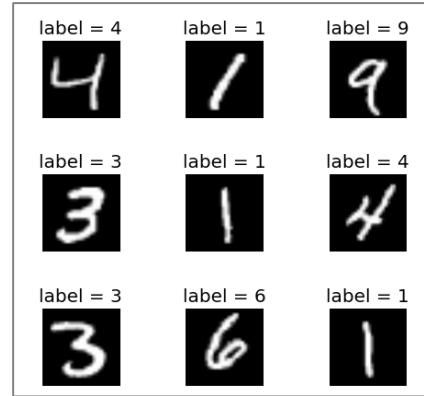
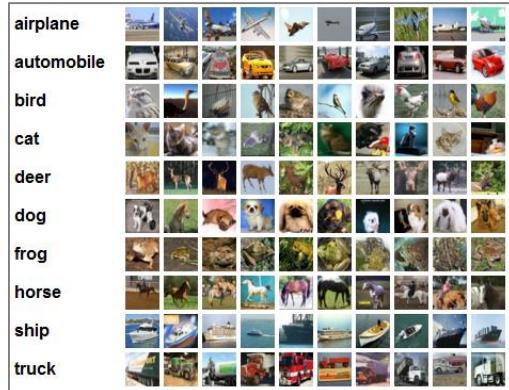
# 深度學習的美好想像與迷思

- 萬物皆可 Train
  - FizzBuzz 的故事
  - 網路上的資源很多，教學也很多，直接套不難吧？



# 深度學習的美好想像與迷思

- 萬物皆可 Train
  - FizzBuzz 的故事
  - 網路上的資源很多，教學也很多，直接套不難吧？



- 深度學習一定強！
  - 很多時候傳統的演算法反而更好
  - Occam's razor

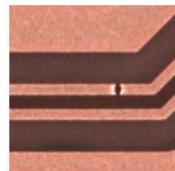




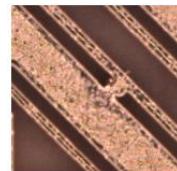
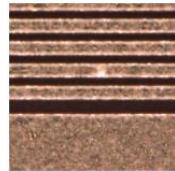
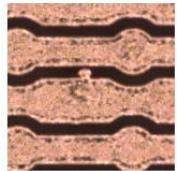
ForGIFs.com



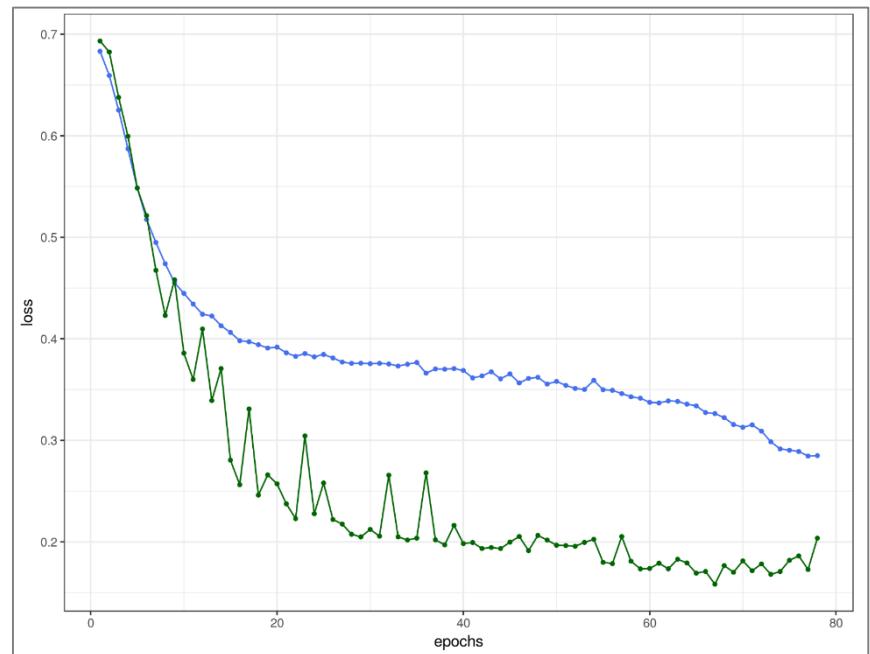
# 1. 異常分類 - 不就是個分類問題嗎?



Fail



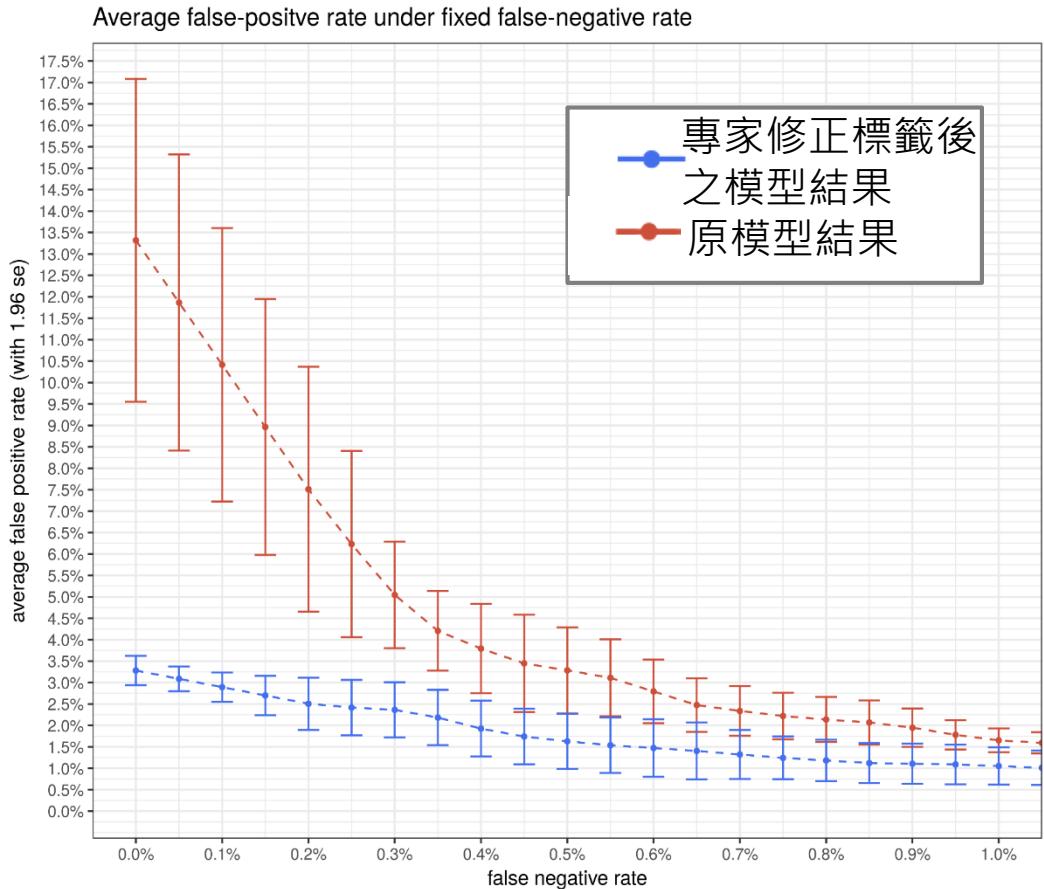
Pass



# 異常分類 - 不就是個分類問題嗎?

原本人工標記的 label 真的是乾淨的?

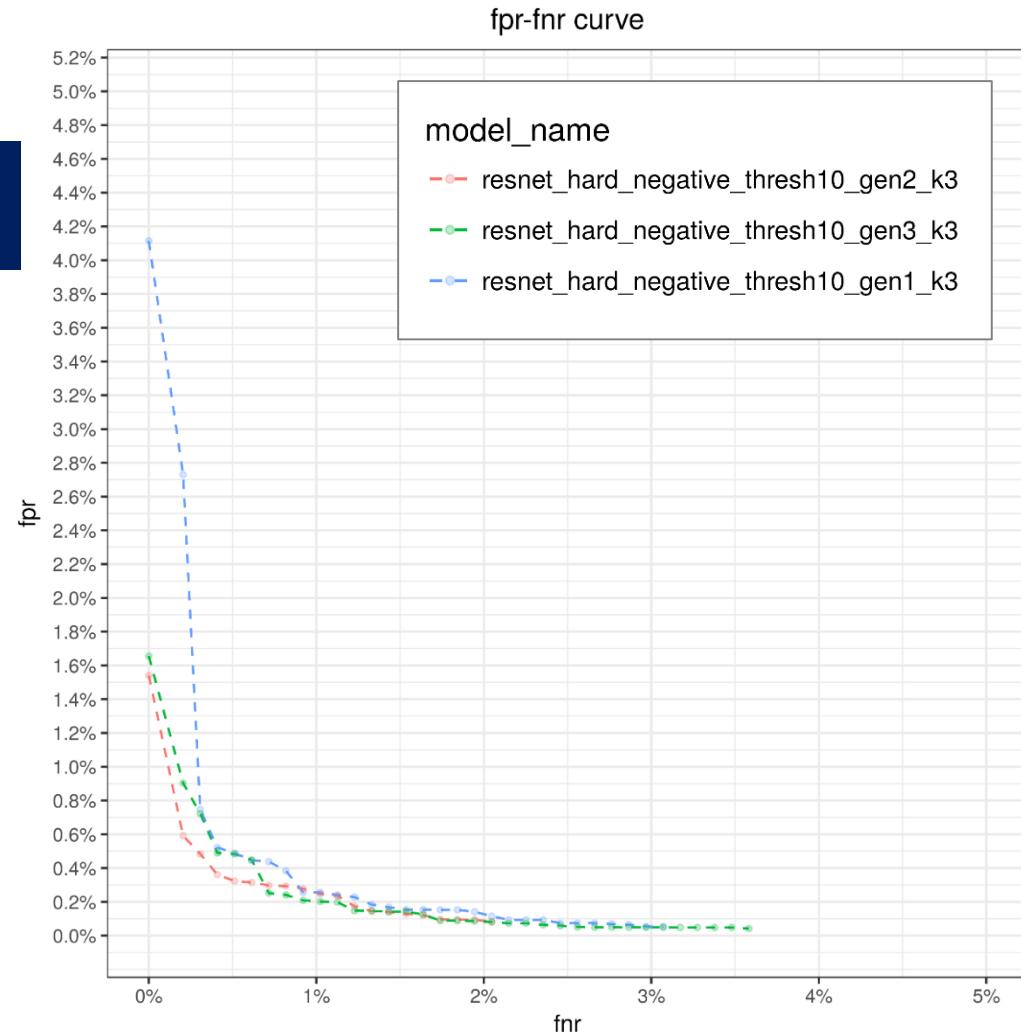
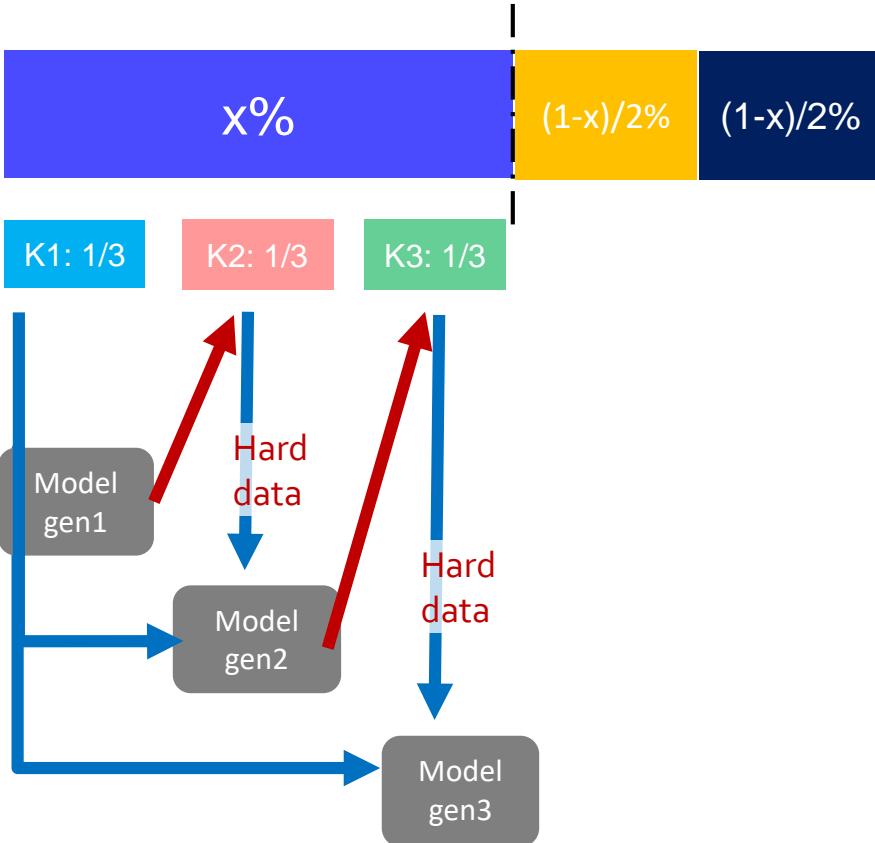
- 專家複判 Web UI



# 異常分類 - 不就是個分類問題嗎?

訓練一次不好, 你有試過訓練兩次嗎?

- Data boosting



## 2. 原料配色 – 神秘的那道光?



顏色差異（英語：Color difference），亦稱距離，是色彩學上的一個關注點。它量化了一個概念。在未量化之前，人們只能用形容詞來大概描述這個概念，這使得對顏色要求嚴格的工作者們很不方便。顏色差異可以通過色彩空間內的歐氏距離簡單計算得出，也可以使用CIE較為複雜、均勻的人類知覺公式計算。

## 2. 原料配色 – 神秘的那道光?

$$\Delta E = \sqrt{\left(\frac{\Delta L}{lS_L}\right)^2 + \left(\frac{\Delta C}{cS_C}\right)^2 + \left(\frac{\Delta H}{S_H}\right)^2}$$

where

$$\Delta C = C_1 - C_2$$

$$C_1 = \sqrt{a_1^2 + b_1^2}$$

$$C_2 = \sqrt{a_2^2 + b_2^2}$$

$$\Delta H = \sqrt{\Delta a^2 + \Delta b^2 - \Delta C^2}$$

$$\Delta L = L_1 - L_2$$

$$\Delta a = a_1 - a_2$$

$$\Delta b = b_1 - b_2$$

$$S_L = \begin{cases} 0.511 & \text{if } L_1 < 16 \\ \frac{0.040975L_1}{1+0.01765L_1} & \text{if } L_1 \geq 16 \end{cases}$$

$$S_C = \frac{0.0638C_1}{1 + 0.0131C_1} + 0.638$$

$$S_H = S_C(FT + 1 - F)$$

$$T = \begin{cases} 0.56 + |0.2 \cos(H_1 + 168^\circ)| & \text{if } 164^\circ \leq H_1 \leq 345^\circ \\ 0.36 + |0.4 \cos(H_1 + 35^\circ)| & \text{otherwise} \end{cases}$$

$$F = \sqrt{\frac{C_1^4}{C_1^4 + 1900}}$$

$$H = \arctan\left(\frac{b_1}{a_1}\right)$$

$$H_1 = \begin{cases} H & \text{if } H \geq 0 \\ H + 360^\circ & \text{otherwise} \end{cases}$$

## 2. 原料配色 – 神秘的那道光?

$$\Delta E = \sqrt{\left(\frac{\Delta L}{lS_L}\right)^2 + \left(\frac{\Delta C}{lS_C}\right)^2 + \left(\frac{\Delta H}{lS_H}\right)^2}$$

where

$$\Delta C = C_1 - C_2$$

$$C_1 = \sqrt{a_1^2 + b_1^2}$$

$$C_2 = \sqrt{a_2^2 + b_2^2}$$

$$\Delta H = \sqrt{\Delta a^2 + \Delta b^2}$$

$$\Delta L = L_1 - L_2$$

$$\Delta a = a_1 - a_2$$

$$\Delta b = b_1 - b_2$$

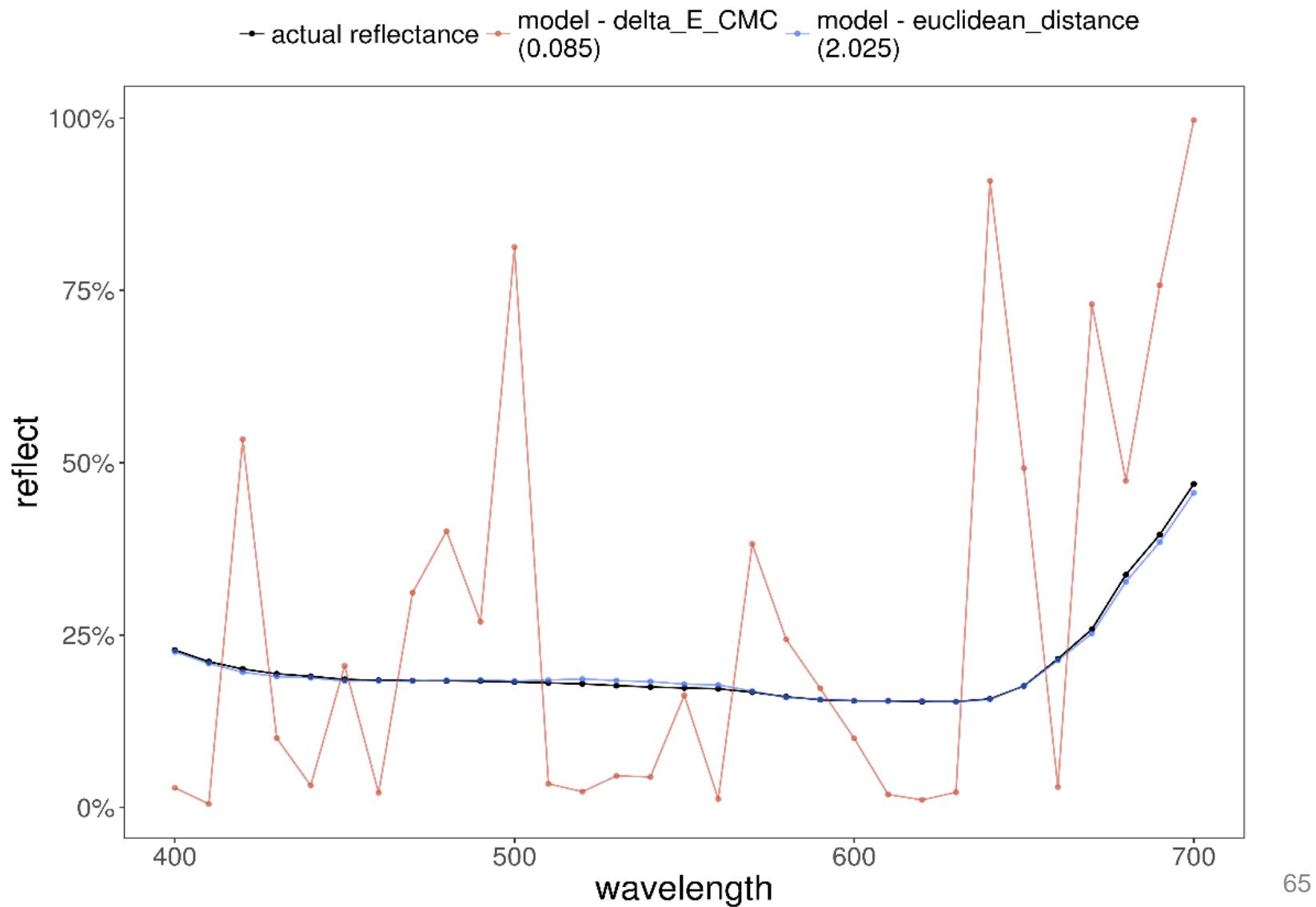


if  $164^\circ \leq H_1 \leq 345^\circ$   
otherwise

$$H_1 = \begin{cases} H & \text{if } H \geq 0 \\ H + 360^\circ & \text{otherwise} \end{cases}$$

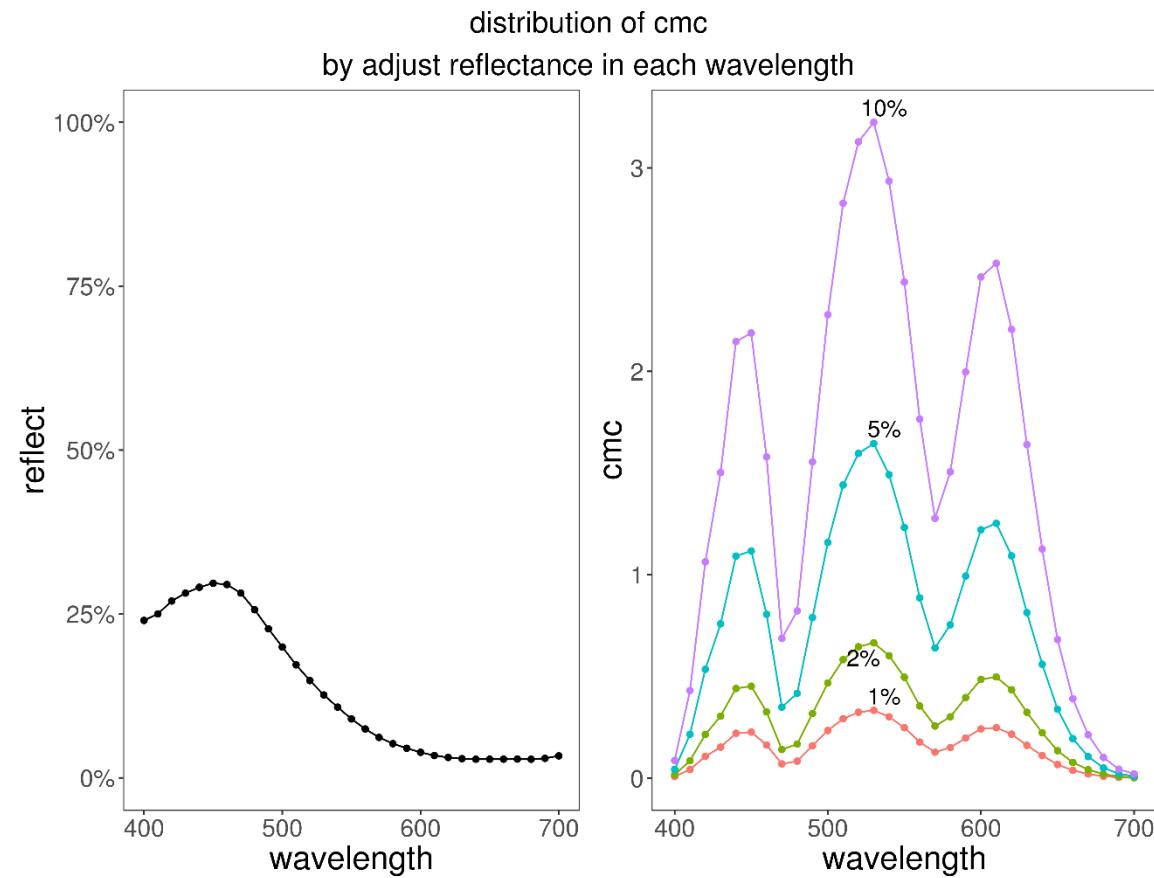
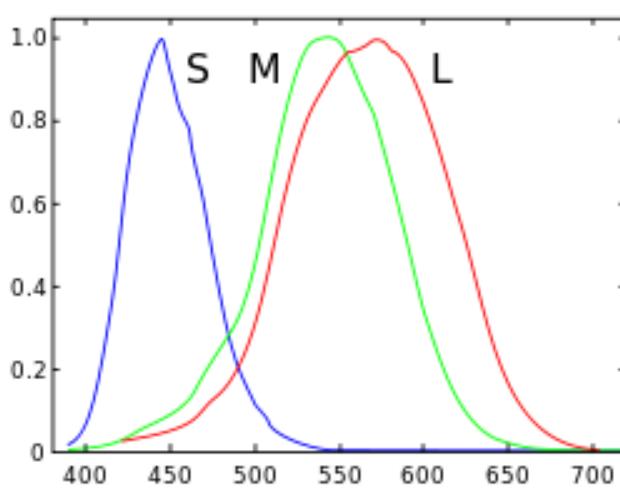
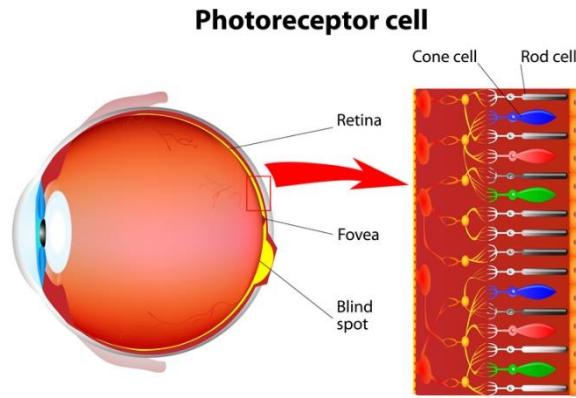
# 原料配色 – 神秘的那道光?

直觀的 Loss function = 詭異的 fitting result? : dECMC vs. eud.loss



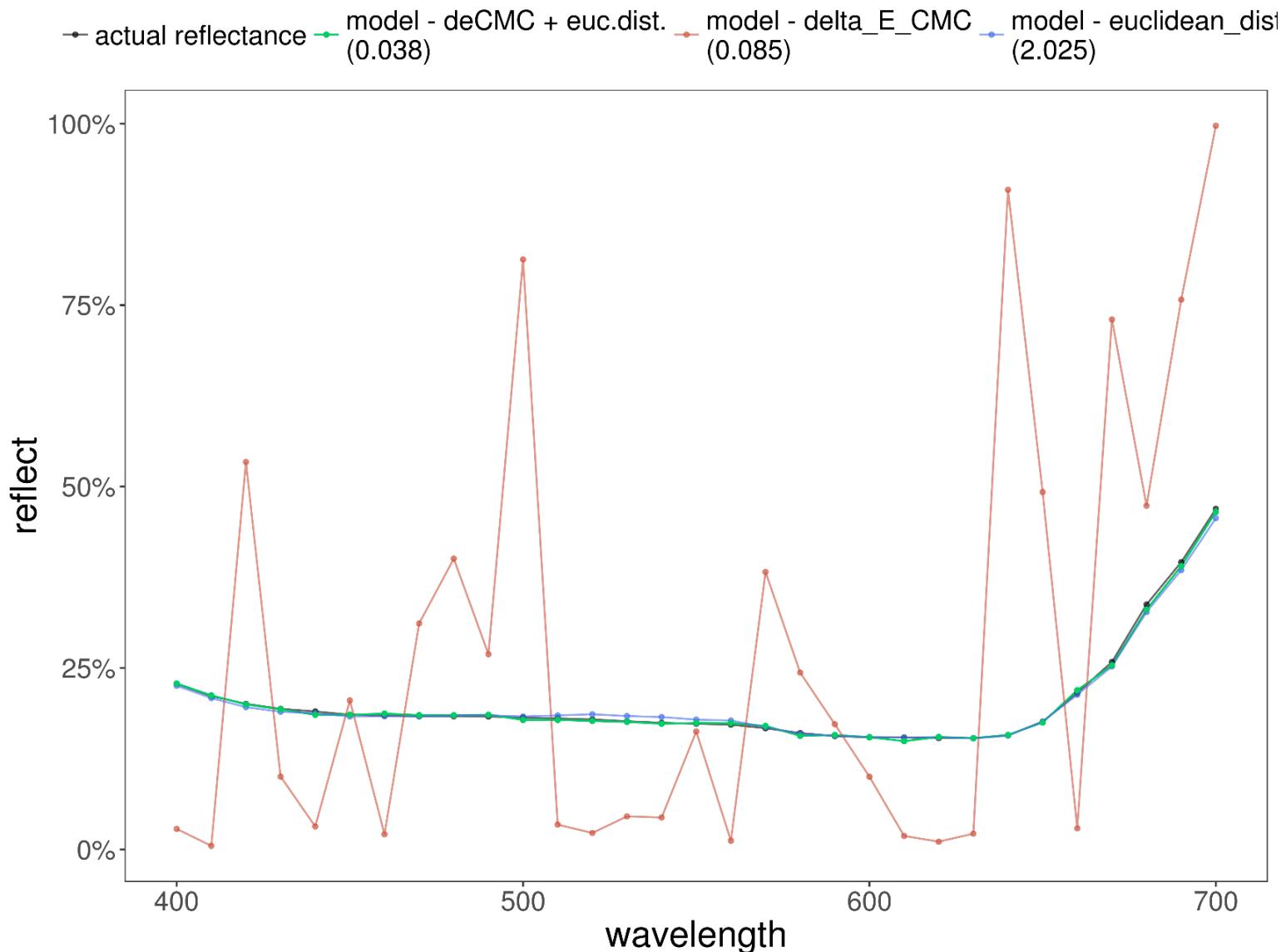
# 原料配色 – 神秘的那道光?

為何出現詭異的結果?



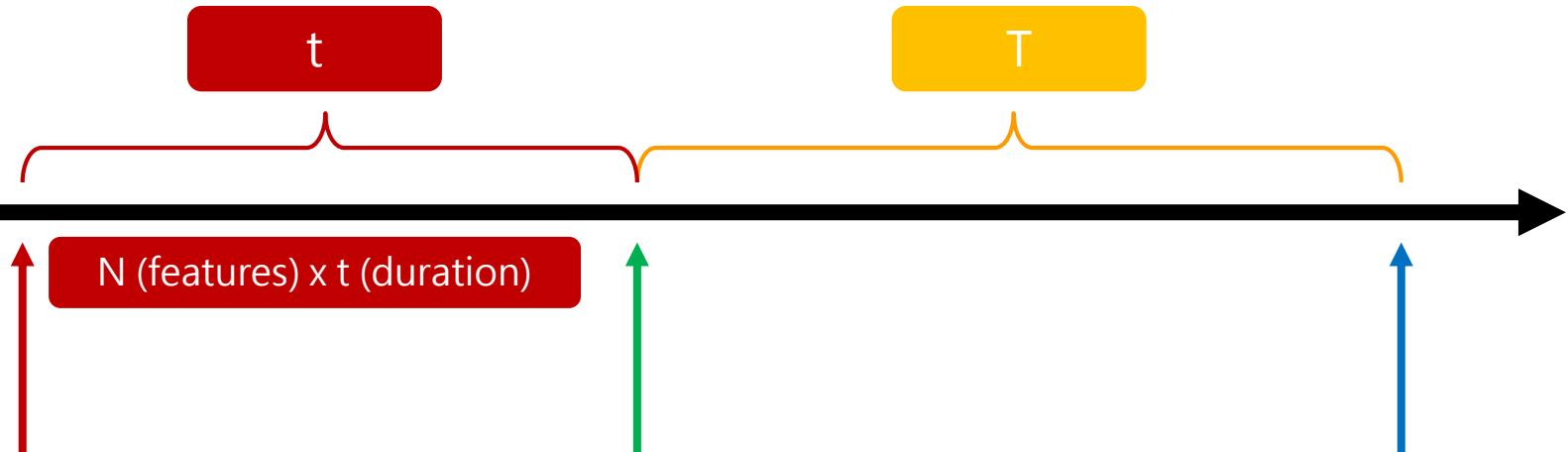
# 原料配色 – 神秘的那道光?

改成兩階段的 model fitting



### 3. 預測性維護 – 又是分類問題?

- 標準問題定義



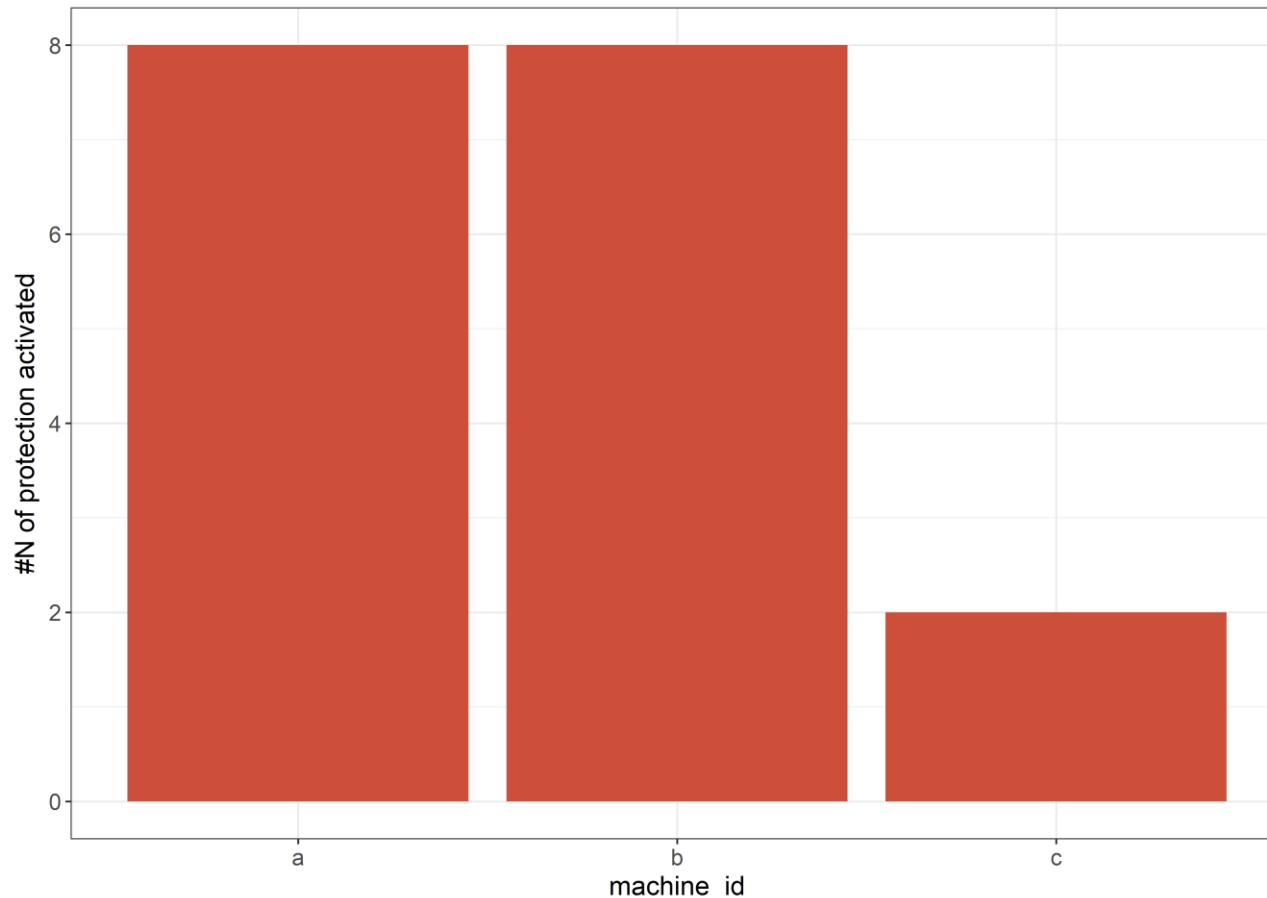
:(  
:(

Your PC ran into a problem and needs to restart. We're just collecting some error info, and then we'll restart for you.

# 預測性維護 – 又是分類問題?

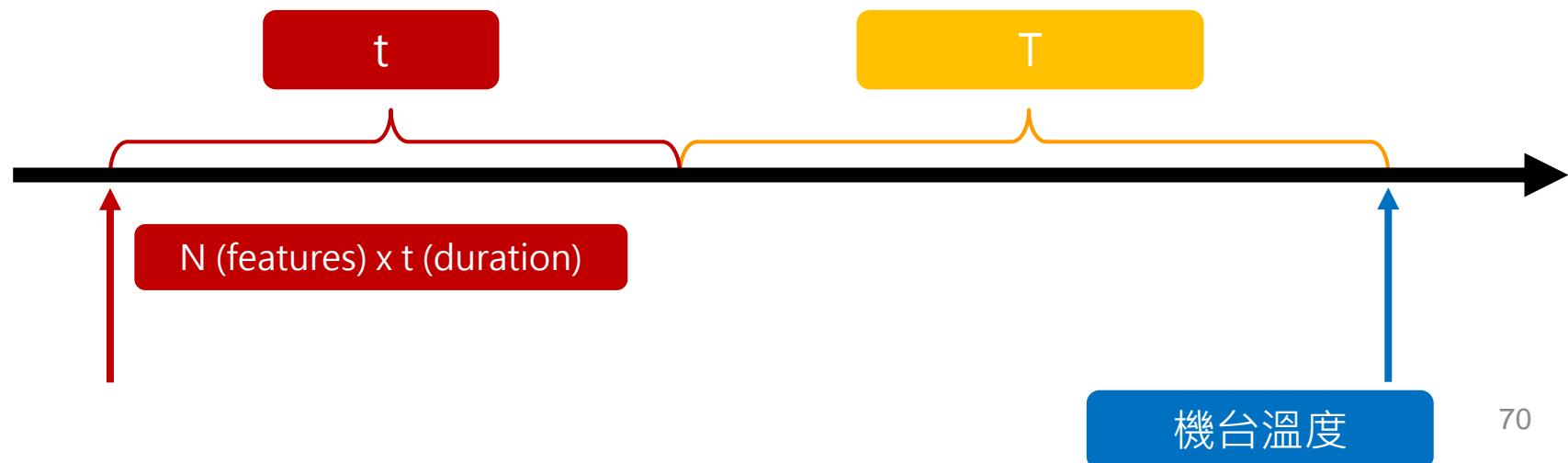
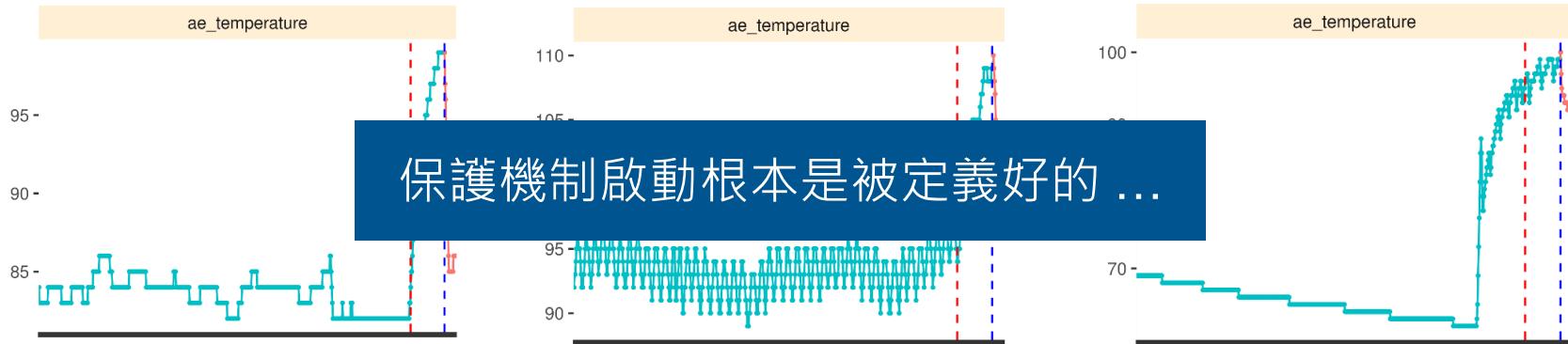
案情並不單純

- 八個月來... 最多也只出現八次異常?!
  - 每 5 秒一筆資料, 八個月就是 ...五百多萬筆 ...
- “異常” 不是損壞, 只是保護機制 ...



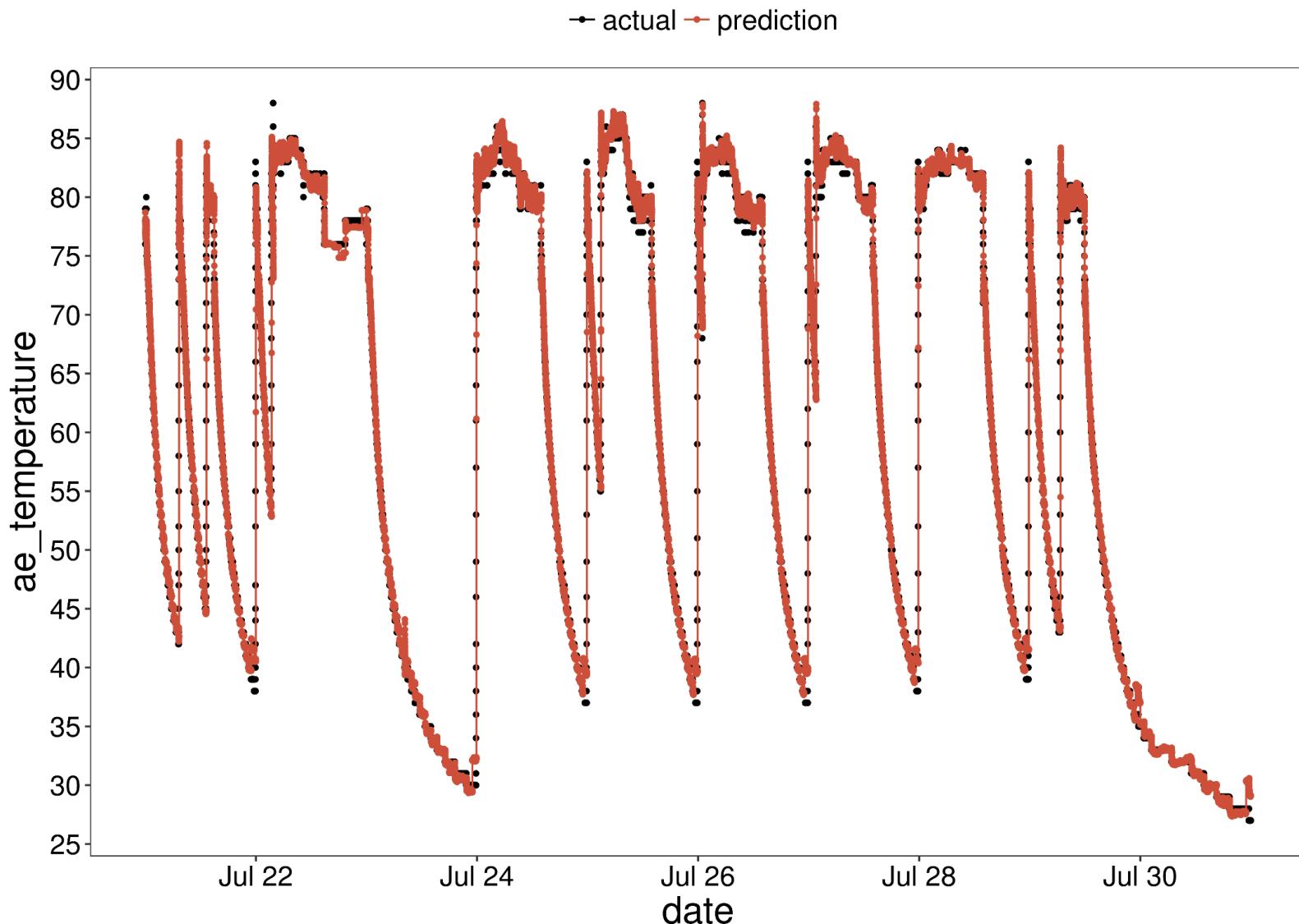
# 預測性維護 – 又是分類問題?

重新定義問題



# 預測性維護 – 又是分類問題?

轉換目標：變成迴歸問題



**借模型很容易，但是 ...**  
**問對問題、了解的模型限制、山不轉路轉**  
**才能讓你在陷阱重重的現實世界存活**

---

# Take home message

---

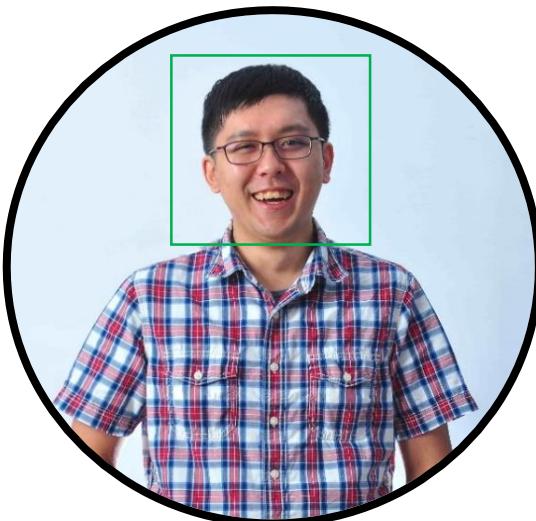
- 對於 ”跨領域”
  - 資料分析不單是資工、電機的領域
    - Lab: 會計、應用數學、風險管理、土木、醫學、光電, ...
  - 能帶來更多不同的思維，激盪出新的火花
  - 要踏入之前，可以先多上點線上課程或多聽演講，確認你真的有興趣
    - 資料分析中 90% 的時間是很枯燥乏味的，但當成功找到有意義的結果時，會覺得努力都是值得的



# Take home message

---

- “資料分析” 跟 “深度學習” 這件事
  - 別衝動，好好先檢查資料
  - 思考問題的本質是什麼
  - 了解模型的限制、改變模型的架構
    - 窮則變，變則通
- **結果太好的時候，不要高興得太早！**
  - 莫忘綠框



# Thank You!