



徐唯哲  
Paul Wei-Che HSU

# 機器學習與多體學資料分析

Machine learning and integrative analysis of multi-omics data

分子與基因醫學研究所  
專案技術助研究員

2023/9/15

# 課程相關資料



[https://drive.google.com/drive/folders/1JwKLgLGlvscr9fVsAKc\\_AacJyYgkEfZ-?usp=sharing](https://drive.google.com/drive/folders/1JwKLgLGlvscr9fVsAKc_AacJyYgkEfZ-?usp=sharing)

蝗災

# 群聚 (Crowding)

引發的災難



# 蝗災的起源: CYP305M2

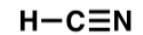
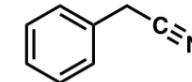
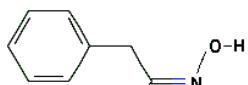
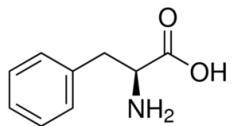
蚱蜢



Crowding

Active CYP305M2

蝗蟲



CYP305M2

CYP71

L-Phenylalanine (苯丙胺酸) → Z-Phenylacetaldoxime (苯乙醛肟) → Phenylacetonitrile, PAN (苯乙腈) → HCN (氰化氫)



# 細胞色素P450 CYP450 (Cytochrome P450)

CYP enzymes cover about **80%** of oxidative metabolism (Zhao et al., 2021)

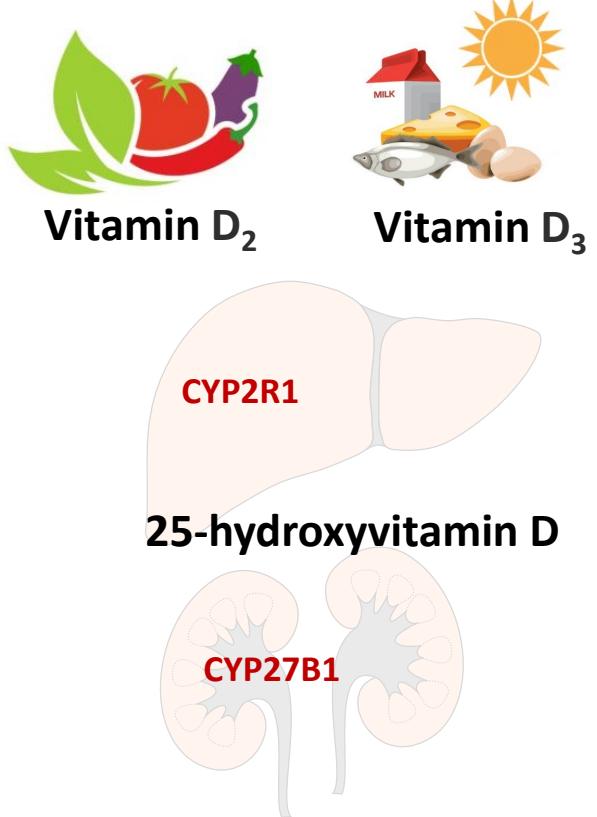
**70-80%** of clinical drugs are metabolized by CYP enzymes

(Zanger and Schwab., 2013)

藥物學名	適應症	相關基因
Clomipramine	憂鬱症	CYP2D6
Clopidogrel	中風，心肌梗塞，周邊動脈血管疾病	CYP2C19
Clozapine	精神分裂症	CYP2D6
Codeine	鎮咳、鎮痛	CYP2D6
Daclatasvir	慢性 C 型肝炎	IFNL3
Darifenacin	膀胱過動症	CYP2D6
Dasabuvir, Omibitasvir, Paritaprevir, and Ritonavir	慢性 C 型肝炎	IFNL3
Desipramine	憂鬱症	CYP2D6
Desvenlafaxine	台灣未上市	CYP2D6
Deutetetabenazine	遲發性運動障礙	CYP2D6
Dexlansoprazole	糜爛性食道炎	CYP2C19
Dextromethorphan & Quinidine	台灣未上市	CYP2D6
Diazepam	焦慮	CYP2C19
Donepezil	阿滋海默症	CYP2D6
Doxepin	憂鬱症，焦慮症，睡眠障礙，搔癢症	CYP2C19
Doxepin	憂鬱症，焦慮症，睡眠障礙，搔癢症	CYP2D6
Dronabinol	癌症化療止吐，刺激食慾	CYP2C9
Drospirenone & Ethynodiol Estradiol	避孕	CYP2C19
Duloxetine	憂鬱症，糖尿病週邊神經痛	CYP2D6
Elagolix	子宮內膜異位症	SLCO1B1

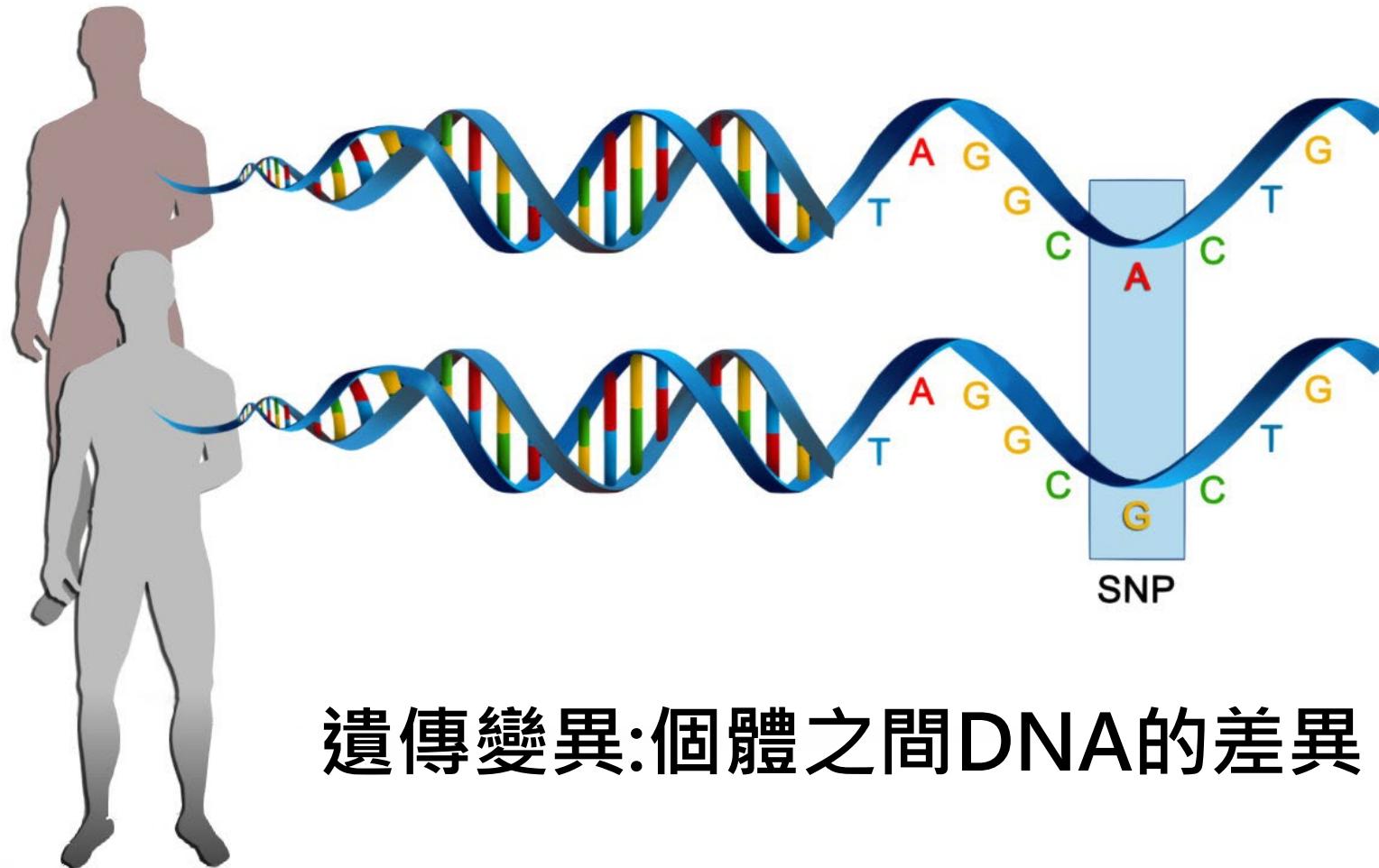
藥物學名	適應症	相關基因
Lesinurad	痛風	CYP2C9
Lofexidine	毒品戒斷	CYP2D6
Lusutrombopag	血小板減少症	F5
Meclizine	暈車，暈船，暈機	CYP2D6
Mercaptopurine	急性白血病，慢性骨髓性白血病	NUDT15
Mercaptopurine	急性白血病，慢性骨髓性白血病	TPMT
Metoclopramide	嘔吐，胃食道逆流，糖尿病引起之胃腸異常	CYP2D6
Metoprolol	高血壓，心絞痛，心衰竭，心肌梗塞，心律不整	CYP2D6
Mirabegron	膀胱過動症	CYP2D6
Modafinil	猝睡症	CYP2D6
Nebivolol	高血壓	CYP2D6
Nefazodone	憂鬱症	CYP2D6
Nortriptyline	憂鬱症	CYP2D6
Ombitasvir, Paritaprevir, and Ritonavir	慢性 C 型肝炎	IFNL3
Omeprazole	胃潰瘍，十二指腸潰瘍，胃食道逆流	CYP2C19
Ondansetron	癌症化療止吐	CYP2D6
Ospemifene	女性停經後骨質疏鬆	CYP2C9
Palonosetron	癌症化療止吐	CYP2D6
Pantoprazole	胃潰瘍，十二指腸潰瘍，胃食道逆流	CYP2C19
Paroxetine	憂鬱症，恐慌症，焦慮症，強迫症	CYP2D6
Peginterferon Alfa-2b	慢性 C 型肝炎	IFNL3
Perphenazine	鎮靜	CYP2D6

為什麼每個人獲得的維他命D的能力不一樣？



I, 25-dihydroxyvitamin D  
1,25(OH)<sub>2</sub>D

# Genetic Variants



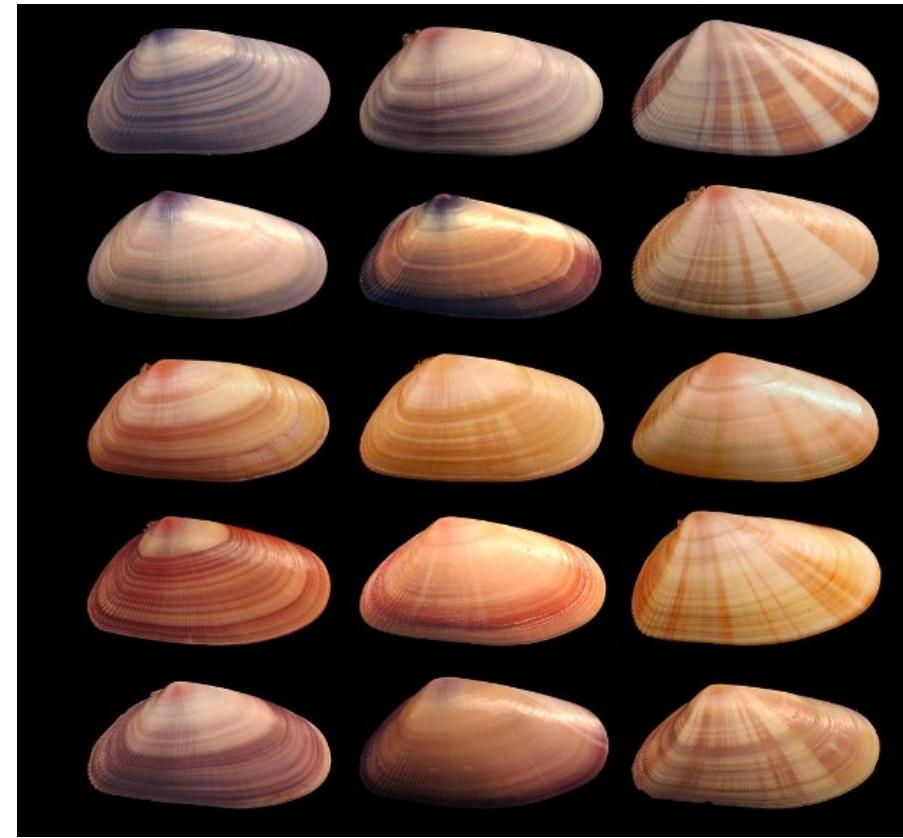
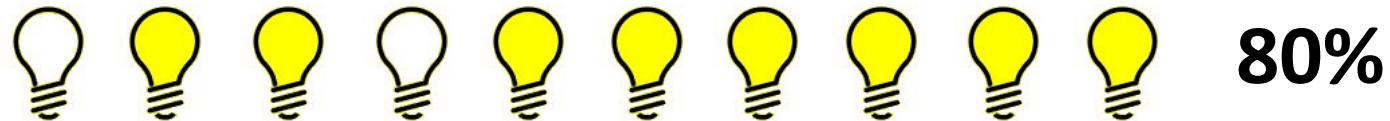
遺傳變異:個體之間DNA的差異

# Phenotype Variability

表型多變性

- 外顯率 (Penetrance)

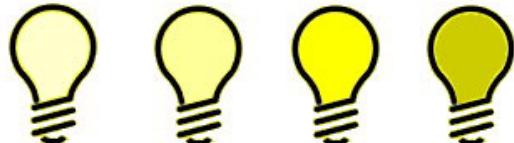
相同基因突變下，造成phenotype表現的比率



(<https://www.wikiwand.com/en/Phenotype>)

- 表現度 (Expressivity)

相同基因突變下，phenotype表現的程度



# CYP450, CYP (細胞色素P450)



<https://www.pharmvar.org/>

23 Pharmacogenes

21 CYP genes

## CYP1 family

CYP1A1  
CYP1A2  
CYP1B1

## CYP2 family

CYP2A6  
CYP2A13  
CYP2B6  
CYP2C8  
CYP2C9  
CYP2C19  
CYP2D6  
CYP2E1  
CYP2F1  
CYP2J2  
CYP2R1  
CYP2S1  
CYP2W1

## CYP3 family

CYP3A4  
CYP3A5  
CYP3A7  
CYP3A43

## CYP4 family

CYP4F2

## Non-Cytochrome P450

DPYD  
NUDT15



Clopidogrel  
(pro-drug) (保栓通,  
抗凝血藥物)

ABCB1

Small intestine

Intestinal absorption

<15% >85%

CYP2C19  
CYP2B6  
CYP1A2

CES1

Liver

CYP2C19  
CYP2C9

2-oxo Clopidogrel

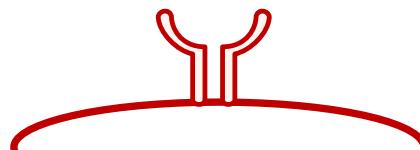
CYP2B6  
CYP3A4

Clopidogrel  
Inactive metabolite

CES1

Clopidogrel  
(active metabolite)

P2Y<sub>12</sub> receptor



降低血小板的活性  
抑制血栓

# Platelet reactivity unit (PRU)

血小板高反應性



Platelet



Activated Platelet

PRU

低

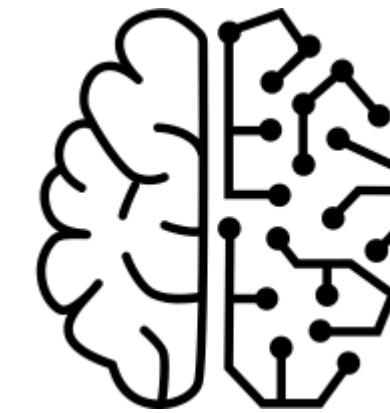
(抗凝血效果佳)



高  
(抗凝血效果差)

sample_id	PRU	CYP2C19	CYP2B6	CYP2C9	CYP3A4	CYP3A5
1	1	*1B/*1B	*1/*4	*1/*1	*1/*1	*3/*3
2	6	*1B/*1B	*1/*1	*1/*1	*1/*1	*1/*1
3	6	*2/*2	*1/*6	*1/*1	*1/*1	*1/*3
4	6	*1B/*2	*1/*6	*1/*1	*1/*1	*3/*3G
5	102	*1B/*1B	*1/*6	*1/*1	*1/*1	*3/*3
6	106	*1B/*1B	*1/*2	*1/*1	*1/*5	*1/*3
7	119	*1B/*1B	*1/*1	*1/*1	*1/*1	*3/*3
8	121	*1B/*2	*1/*1	*1/*1	*1/*18	*1/*3
9	122	*1B/*1B	*1/*1	*1/*1	*1/*1	*3/*3
10	131	*1B/*2	*1/*6	*1/*1	*1/*1	*1/*3
11	137	*1B/*3	*6/*6	*1/*1	*1/*1	*3/*3
12	142	*1B/*3	*1/*1	*1/*1	*1/*1	*3/*3
13	143	*1B/*2	*1/*1	*1/*1	*1/*1	*1/*3
14	147	*1B/*1B	*1/*1	*1/*1	*1/*1	*3/*3
15	154	*1B/*2	*1/*6	*1/*1	*1/*1	*3/*3
16	158	*1B/*2	*1/*4	*1/*1	*1/*1	*1/*1
17	160	*1B/*1B	*1/*4	*1/*1	*1/*1	*1/*1
18	171	*1B/*2	*1/*4	*1/*1	*1/*1	*3/*3
19	175	*1B/*2	*2/*26	*1/*1	*1/*1	*3/*3
20	179	*1B/*2	*1/*6	*1/*1	*1/*1	*3/*3
21	188	*1B/*1B	*1/*1	*1/*1	*1/*1	*3/*3
22	197	*2/*2	*1/*1	*1/*1	*1/*1	*1/*3
23	197	*2/*2	*2/*6	*1/*1	*1/*1	*3/*3
24	198	*1B/*1B	*1/*4	*1/*1	*1/*1	*3/*3
25	202	*1/*2	*1/*1	*3/*13	*1/*1	*1/*3
26	204	*1/*1B	*1/*6	*1/*3	*1/*1	*1/*3
27	208	*1B/*1B	*1/*4	*1/*1	*1/*1	*1/*1
28	211	*1/*1+rs4	*1/*4	*1/*3	*1/*1	*1/*3
29	213	*1B/*2	*1/*1	*1/*1	*1/*1	*1/*3
30	215	*1/*1B	*1/*1	*1/*1	*1/*1	*3/*3
31	218	*1B/*2	*1/*6	*1/*27	*1/*1	*1/*3
32	225	*1B/*3	*1/*4	*1/*1	*1/*1	*1/*3
35	240	*1B/*2	*1/*1	*1/*1	*1/*1	*1/*3
36	242	*1B/*2	*1/*1	*1/*1	*1/*1	*1/*3
37	249	*1B/*3	*4/*6	*1/*1	*1/*1	*3/*3
38	257	*3/*3	*1/*6	*1/*1	*1/*1	*3/*3
39	310	*1B/*1B	*4/*6	*1/*1	*1/*18	*1/*1
40	325	*1B/*3	*1/*4	*1/*1	*1/*1	*3/*3
41	356	*1B/*2	*1/*6	*1/*1	*1/*1	*1/*3
42	467	*2/*2	*1/*6	*1/*1	*1/*1	*1/*3

Pharmacogenes



Machine learning

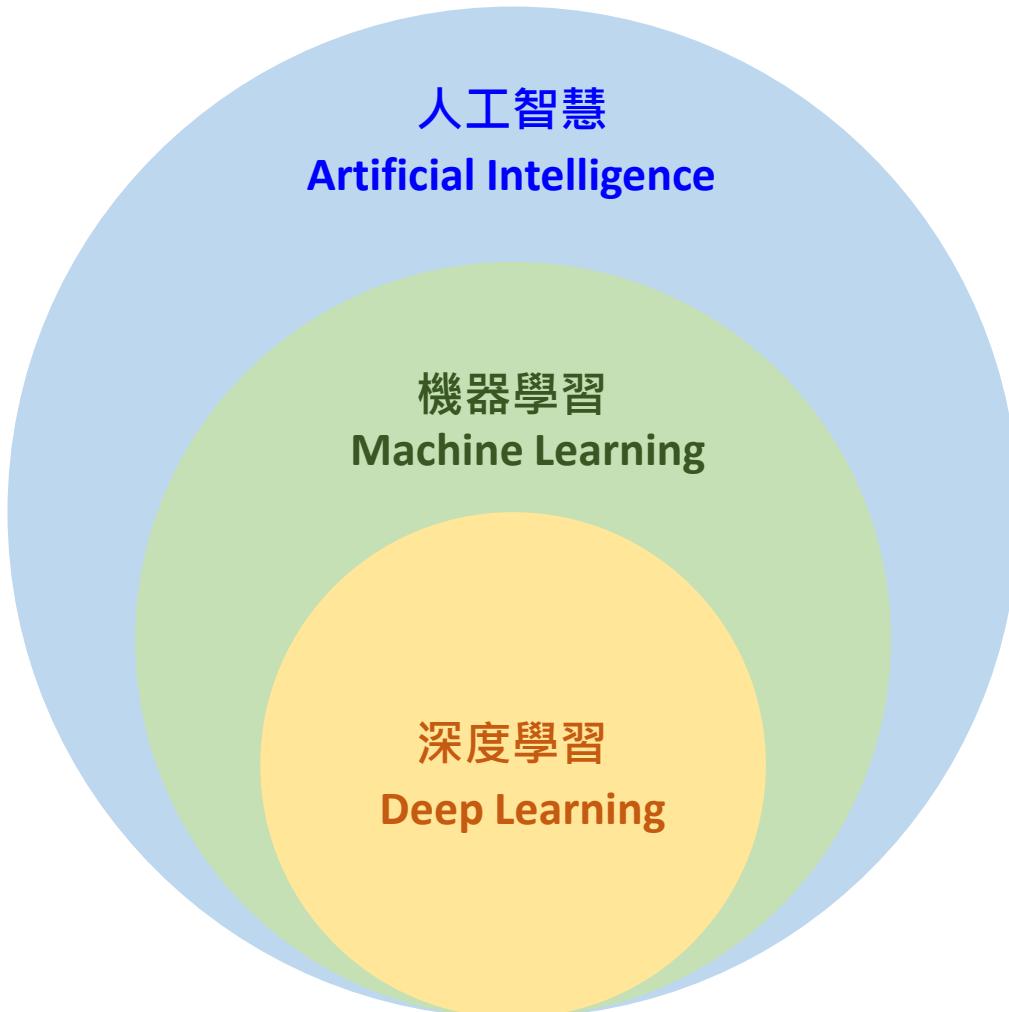
從基因型態預測藥效

## Machine learning



**Machine Learning:**  
**"The field of study that gives  
computers the ability to learn  
without being explicitly  
programmed."**

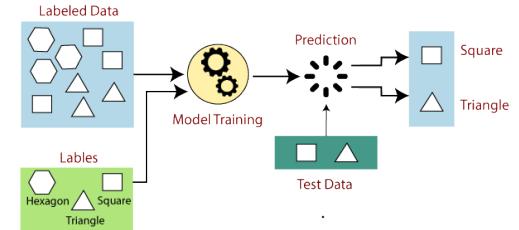
# 人工智能，機器學習與深度學習



# Machine learning

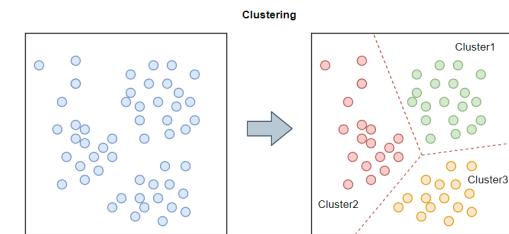
Supervised Learning

監督學習



Unsupervised Learning

無監督學習

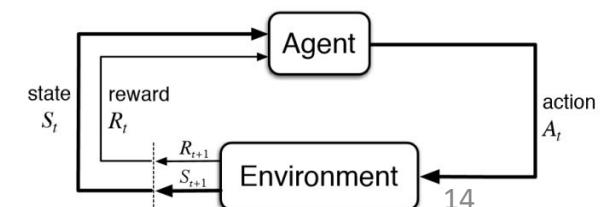


Semi-supervised learning

半監督學習

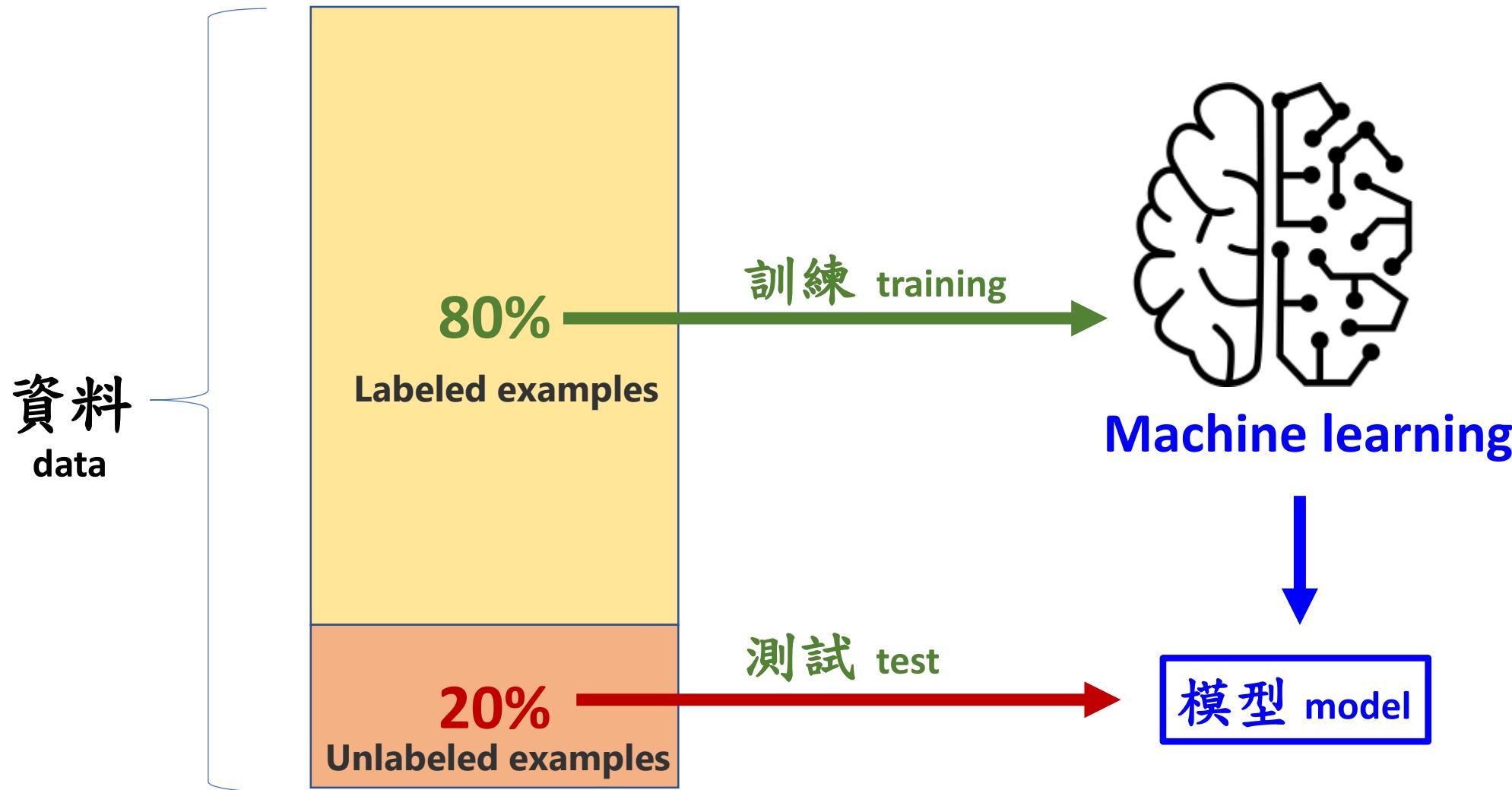
Reinforcement Learning

強化學習



# Supervised Learning

監督學習



# Supervised Learning

監督學習

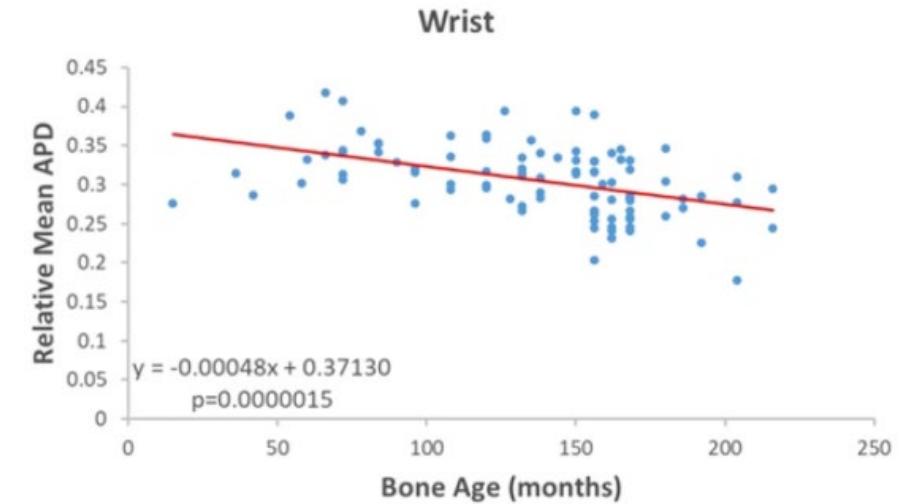
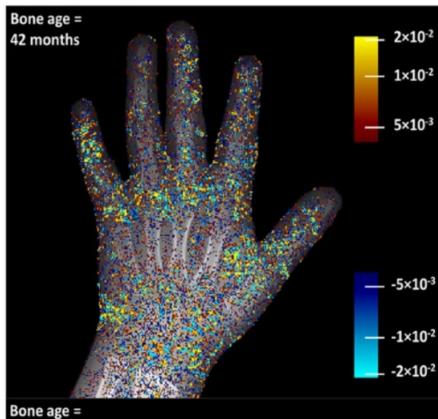
## Regression(回歸)

骨齡預測

房價評估

降雨機率

**predicting a quantity**



## Classification(分類)

良性腫瘤/惡性腫瘤

圖片：食物/狗

垃圾郵件/非垃圾郵件

**predicting a label**

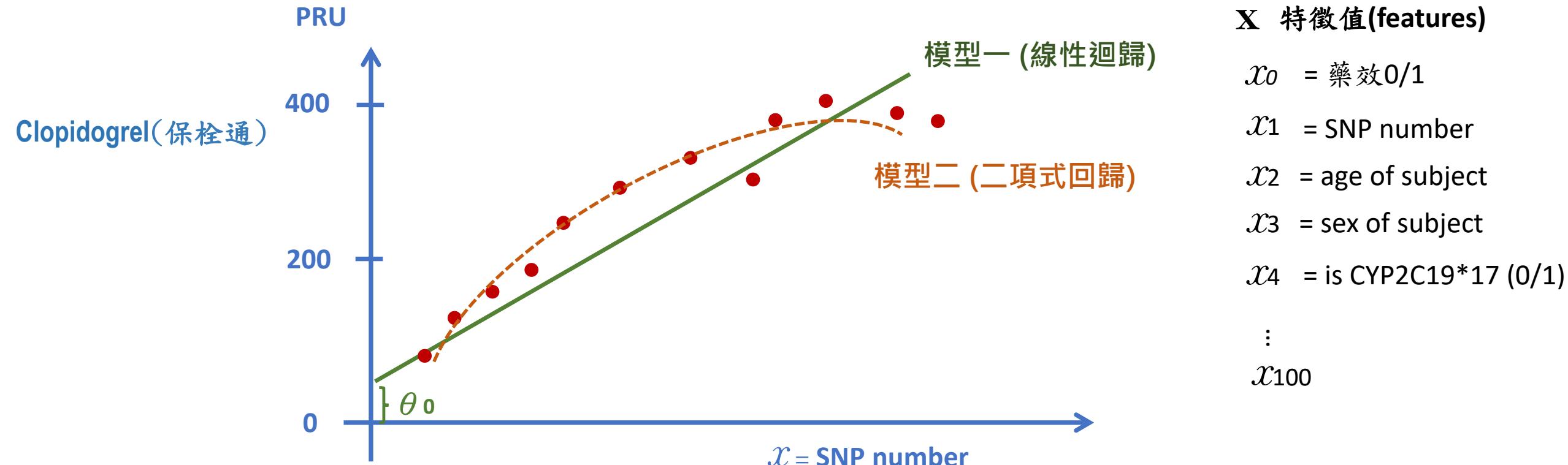


(a)



(b)

(by Karen Zack )



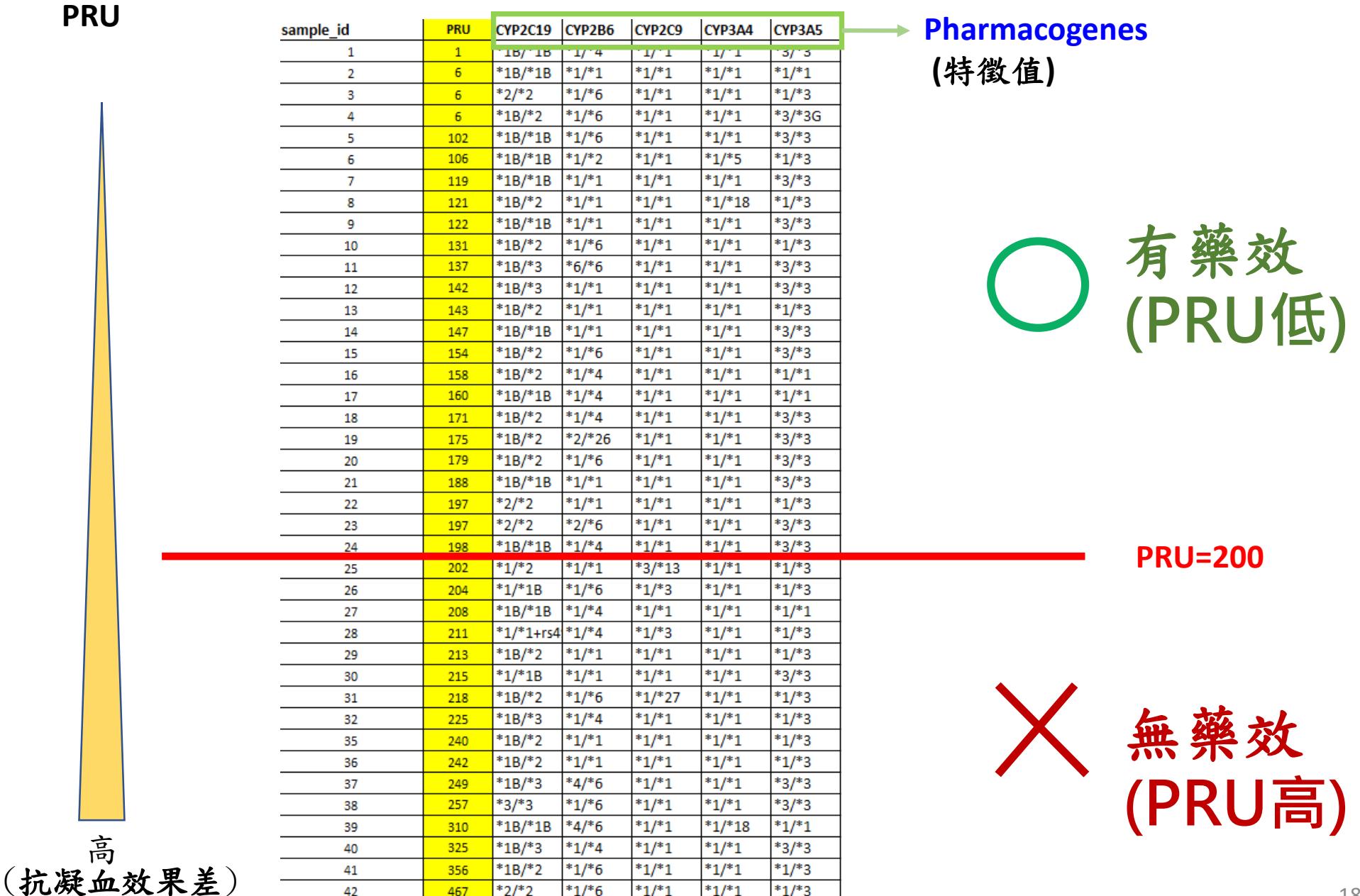
$$h_\theta(x) = \theta_0 + \theta_1 x$$

模型一 (線性迴歸)

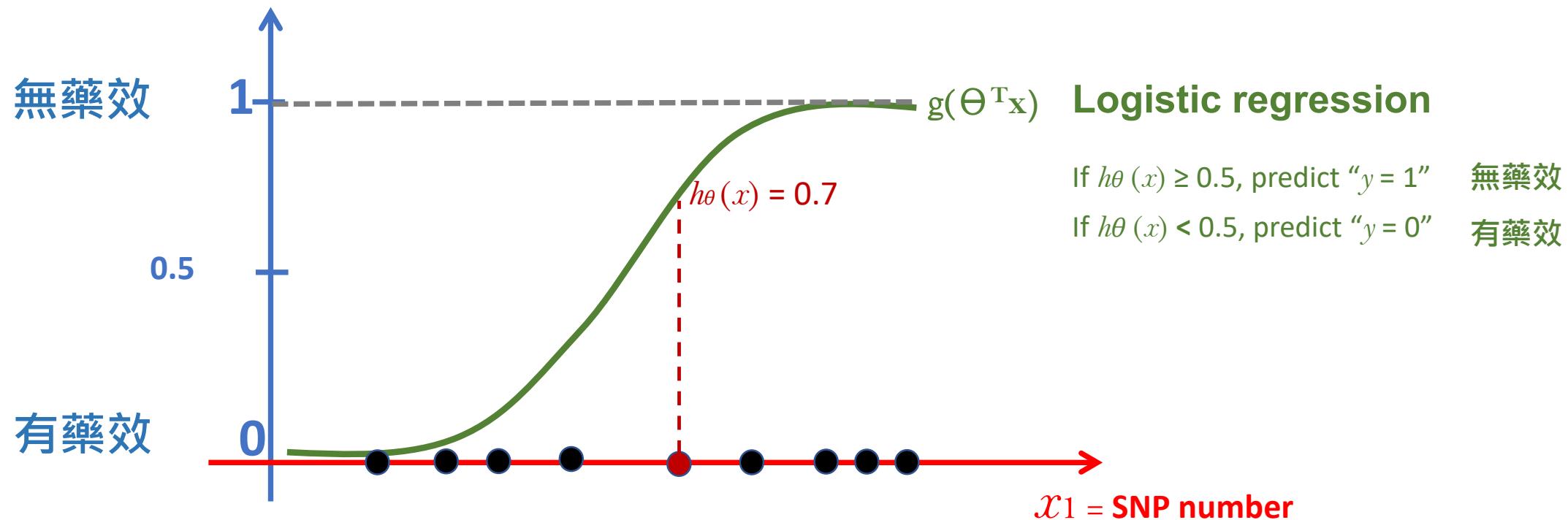
$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

模型二 (二項式回歸)

# Classification(分類)



# 從受試者的特徵值去預測Clopidogrel(保栓通)的藥效



$h_\theta(x)$  = estimated probability that  $y = 1$  on input  $x$

$$\text{If } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{SNP number} \end{bmatrix}$$

$$h_\theta(x) = 0.7$$

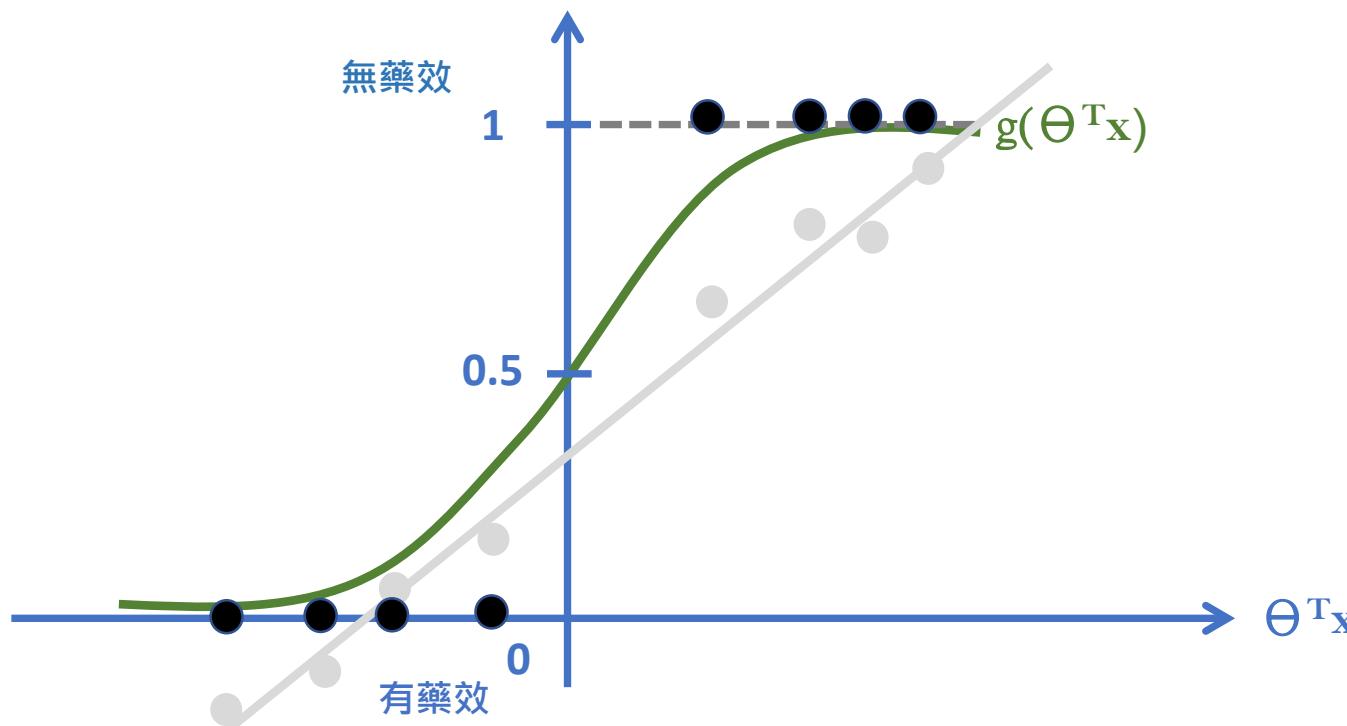
該名患者服藥後70%無效

**X 特徵值(features)**

$x_0$  = 藥效0/1

$x_1$  = SNP number

# 從受試者的特徵值去預測Clopidogrel(保栓通)的藥效



**X 特徵值(features)**

$x_0$  = 藥效 0/1

$x_1$  = SNP number

$x_2$  = age of subject

$x_3$  = sex of subject

$x_4$  = is CYP2C19\*17 (0/1)

:

$x_{100}$

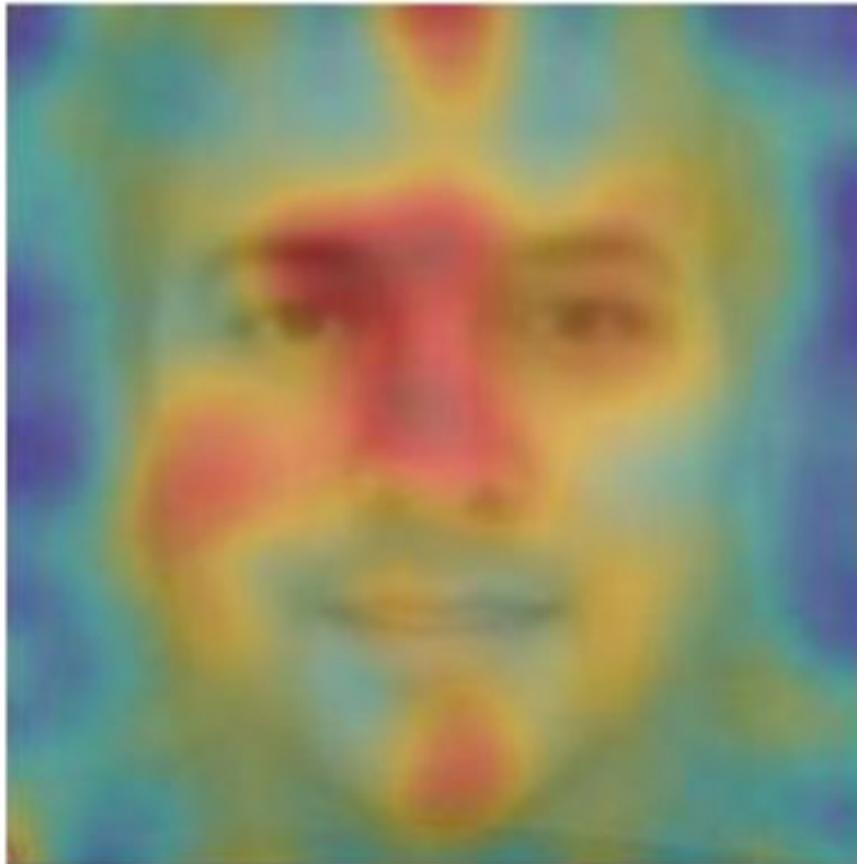
## Logistic regression

$$h_\theta(x) = g(\Theta^T x)$$

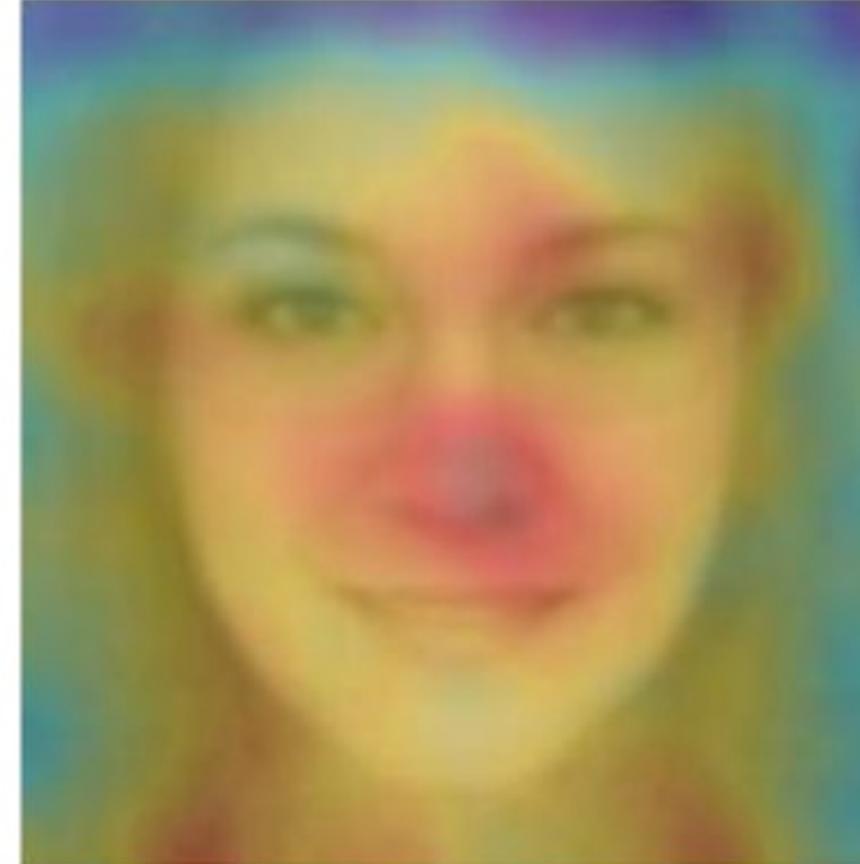
$$= \frac{1}{1 + e^{-\Theta^T x}}$$

Sigmoid function

找出重要的特徵值(features)  $\mathbf{X}$



Male



Female

# 如何診斷代謝症候群？

台灣  
目前的評估標準



2007.1.18公告實施

## ATP-III

代謝症候群的診斷標準  
符合五項指標任三項者



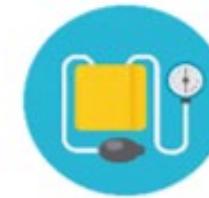
### 腰圍

男 $\geq 90\text{cm}$   
女 $\geq 80\text{cm}$



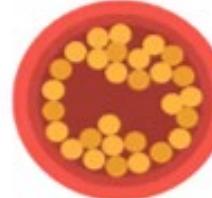
### 血糖

空腹  
 $\geq 100\text{mg}$   
或已用藥



### 血壓

$\geq 130/85$   
mg  
或已用藥



### TG

$\geq 150$   
mg/dL  
或已用藥



### HDL

男 $<40$   
女 $<50$   
mg/dL  
或已用藥

# 如何判斷指數的嚴重性？



n=217,196

1

Scores Based

建立各項指標風險程度的評分

根據不同性別年齡/性別，標準化各項指標 (Z-score)

2

Weighting

分析各項指標的權重

統計模型

Pearson Correlation

Regression

LASSO

Confirmatory factor analysis (CFA)

# 從台灣人體資料庫來看不同性別與年齡，各項指數的重要性也不同

年齡	20-39	40-54	55-64	≥65	
男性	TG	0.593812	0.564152	0.54498	0.529211
	HDL	0.32498	0.323892	0.286285	0.265515
	腰圍	0.894987	0.891831	0.902409	0.938052
	血糖	0.778805	0.683606	0.641712	0.63396
	血壓	0.423434	0.335106	0.261107	0.215299

**重要特徵排名**  
feature scores analysis

- |        | 1. 腰圍  | 1. 腰圍  | 1. 腰圍  | 1. 腰圍 |
|--------|--------|--------|--------|-------|
| 2. 血糖  | 2. 血糖  | 2. 血糖  | 2. 血糖  | 2. 血糖 |
| 3. TG  | 3. TG  | 3. TG  | 3. TG  | 3. TG |
| 4. 血壓  | 4. 血壓  | 4. HDL | 4. HDL | 4. 血壓 |
| 5. HDL | 5. HDL | 5. 血壓  | 5. 血壓  | 5. 血壓 |

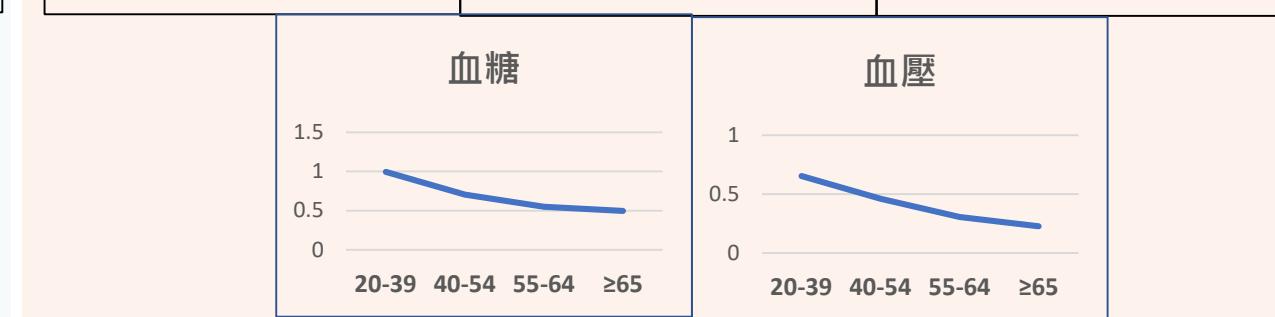
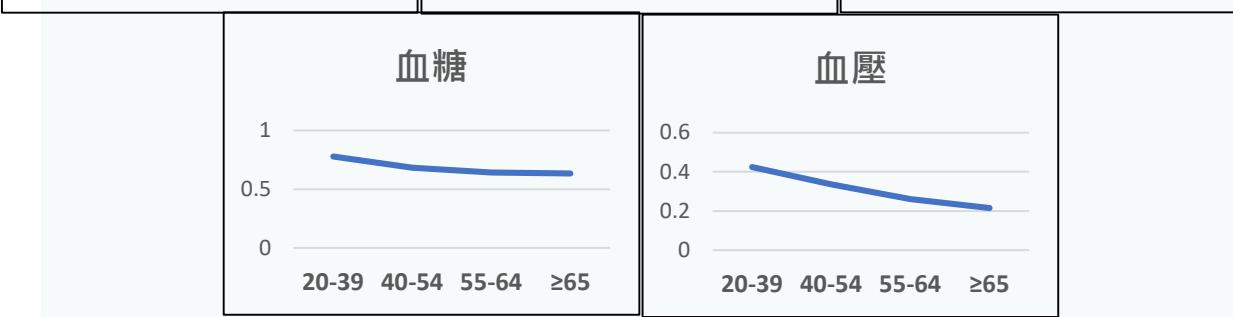
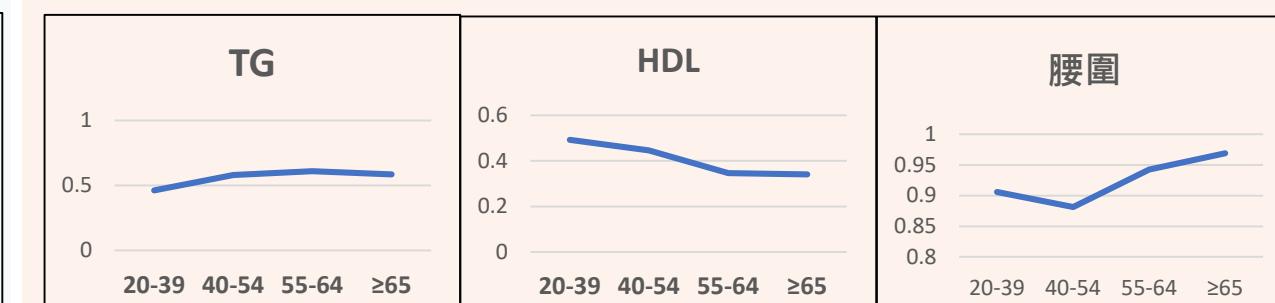
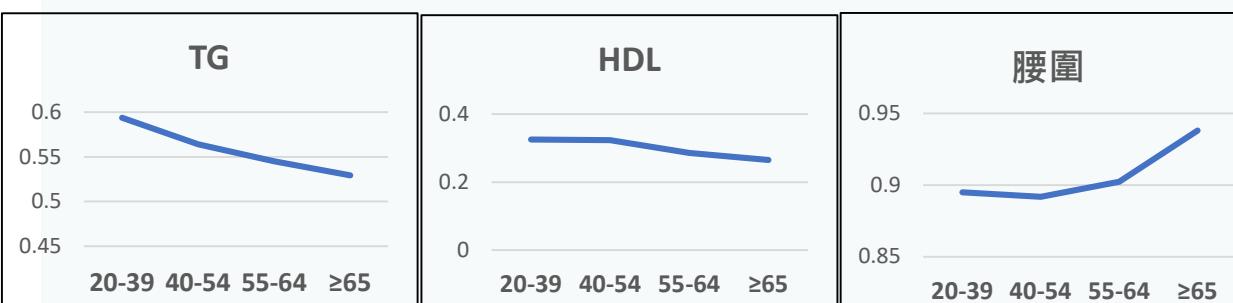
年齡	20-39	40-54	55-64	≥65	
女性	TG	0.460873	0.580948	0.608649	0.584282
	HDL	0.493218	0.446388	0.346302	0.341003
	腰圍	0.905682	0.881332	0.942321	0.969141
	血糖	0.995502	0.706625	0.54936	0.497753
	血壓	0.651754	0.458852	0.306021	0.224916

**重要特徵排名**  
feature scores analysis

- |        | 1. 血糖  | 1. 腰圍  | 1. 腰圍  | 1. 腰圍 |
|--------|--------|--------|--------|-------|
| 2. 腰圍  | 2. 腰圍  | 2. 血糖  | 2. TG  | 2. TG |
| 3. 血壓  | 3. TG  | 3. 血糖  | 3. 血糖  | 3. 血糖 |
| 4. HDL | 4. 血壓  | 4. HDL | 4. HDL | 4. 血壓 |
| 5. TG  | 5. HDL | 5. 血壓  | 5. 血壓  | 5. 血壓 |



隨著年齡的增加，腰圍的重要性逐漸上升。

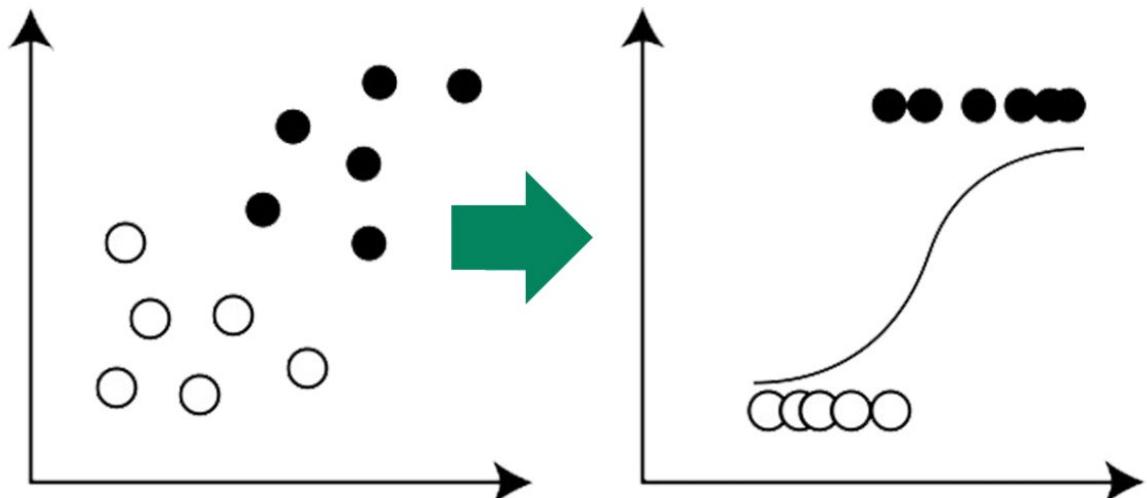


# 監督學習常用的演算法

如何從細胞外觀判定惡性與良性腫瘤

## 羅吉斯迴歸 (Logistic Regression)

分類

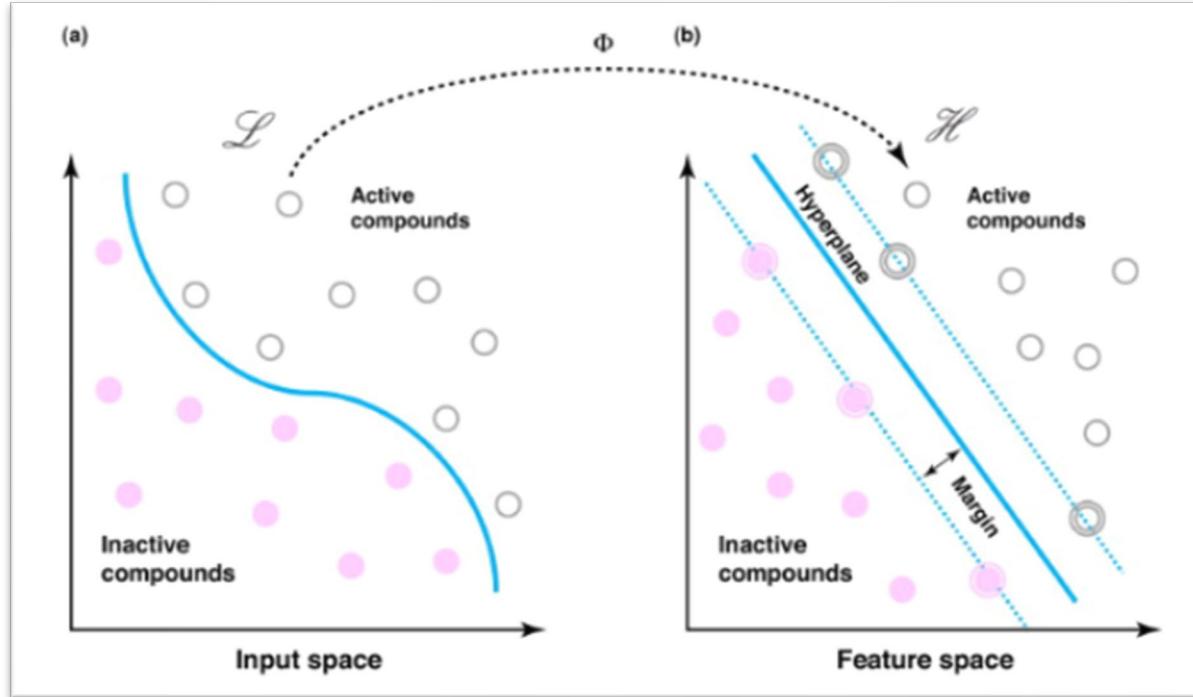


將細胞特徵綜合評分，畫出惡性與良性腫瘤的分界

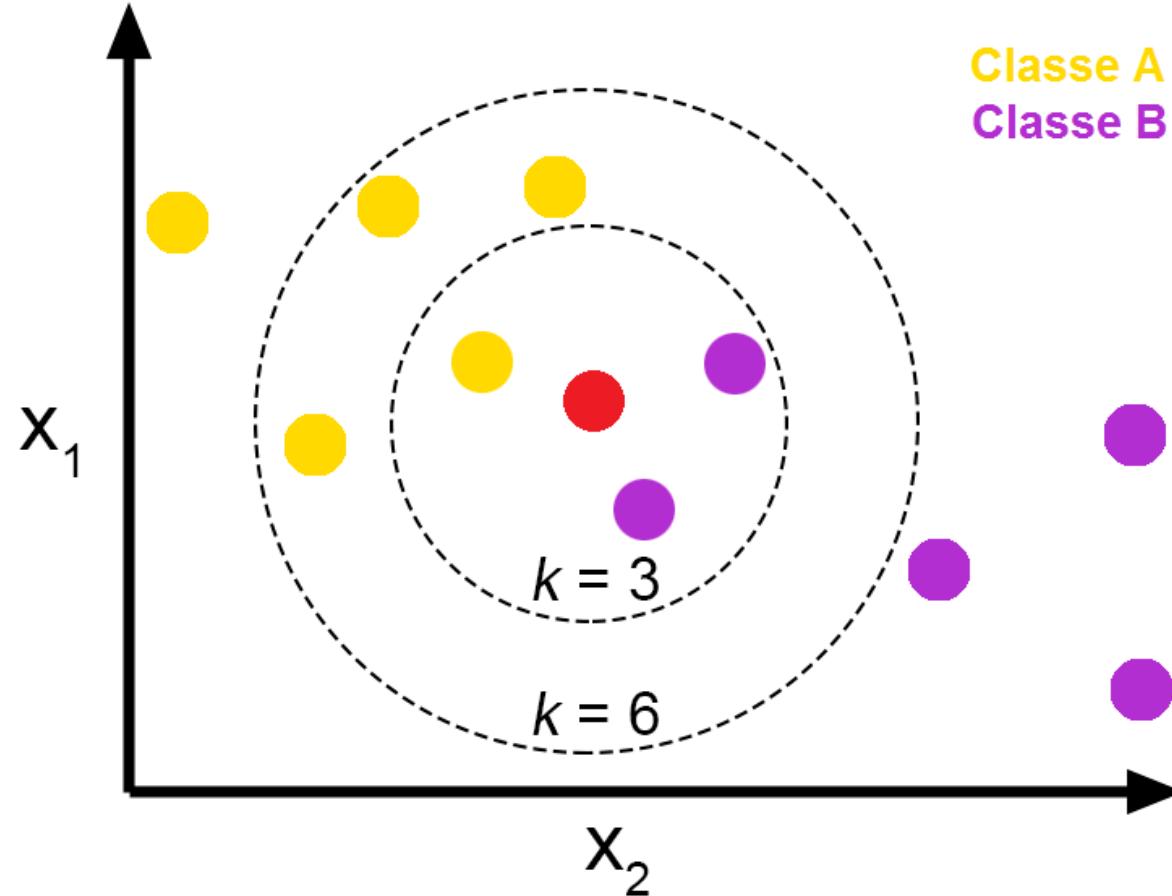
The diagram illustrates the components of the Naive Bayes formula. At the top, 'Likelihood' points to  $P(x|c)$ , 'Class Prior Probability' points to  $P(c)$ , and 'Posterior Probability' points to  $P(c|x)$ . Below this, 'Predictor Prior Probability' points to  $P(x)$ . The formula itself is shown as  $P(c|x) = \frac{P(x|c)P(c)}{P(x)}$ . At the bottom, the formula is expanded to  $P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$ .

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$
$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

根據經驗，腫瘤越大，惡性的機率越大



把二維變成三維，找出超平面分開惡性與良性腫瘤



Classe A 惡性腫瘤  
Classe B 良性腫瘤

腫瘤是良性或惡性，由周圍的細胞表決

# 決策樹 ( Decision Tree )

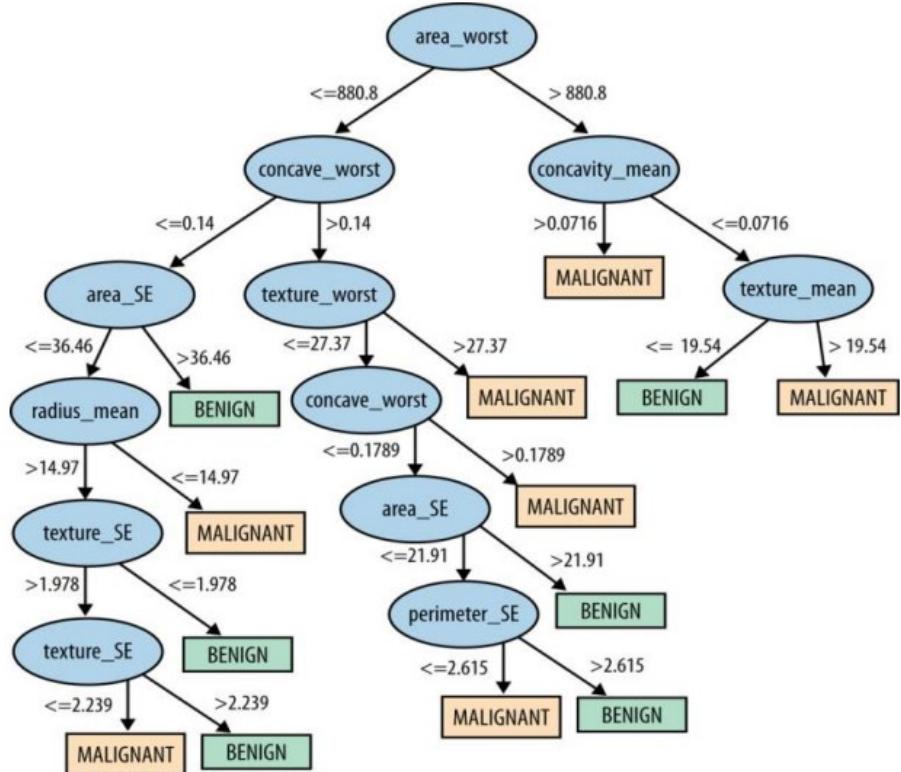
# 分類和迴歸

隨機森林 ( Random Forest )

極端梯度提升樹 (XGBoost)

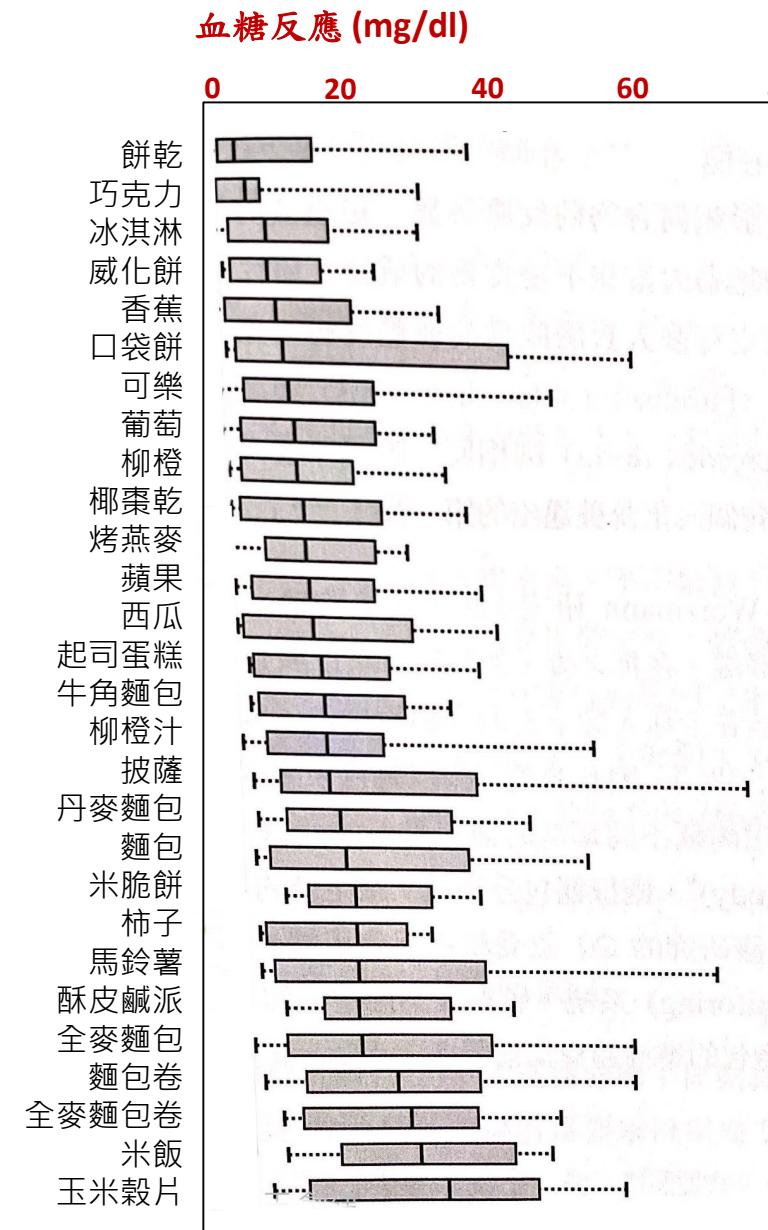
輕量級的高效梯度提升樹(lightGBM)

梯度提升樹 (GBDT)

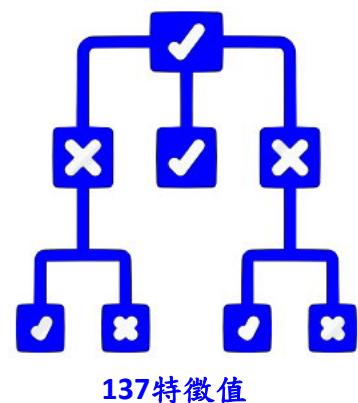


先看腫瘤大小，再看形狀，然後看硬度

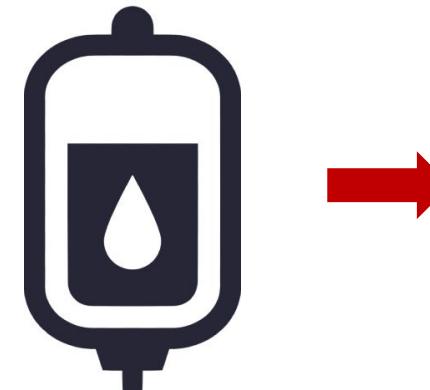
# 由決策樹 ( Decision Tree )來預測個人對特定食物的血糖反應



決策樹 Decision Tree



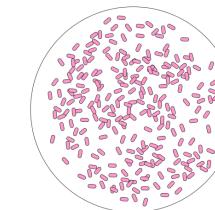
血糖預測模型



腸道菌才是關鍵

狄氏副擬桿菌

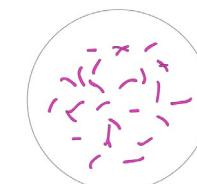
*Parabacteroides distasonis*



↑ 血糖反應

多氏擬桿菌

*Bacteroides dorei*



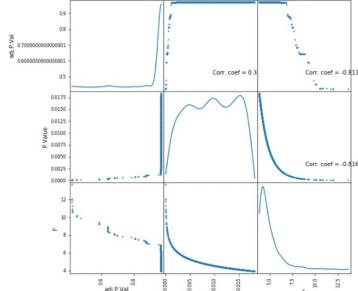
↓ 血糖反應

# 如何入門？

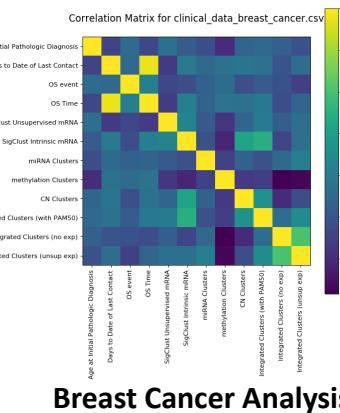


全球最大的資料科學社群

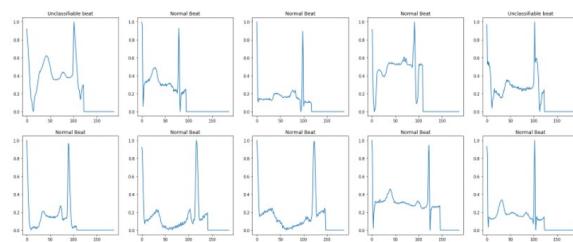
Inside Kaggle you'll find all the code & data you need to do your data science work. Use over 50,000 public datasets and 400,000 public notebooks to conquer any analysis in no time.



Alzheimer Microarray Analysis



Breast Cancer Analysis



ECG Heartbeat Categorization

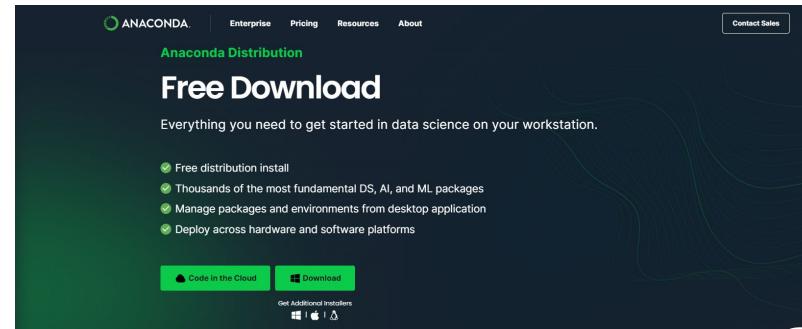
#CHR	POS	REF	ALT	SNP	all_meta_N	all_inv_var_meta_beta	all_inv_var_meta_sebeta	all_inv_var_meta
1	858952	G	A	1:858952::G:A	13	-1.675e-02	6.1199e-02	7.843e-01
1	898883	C	T	1:898883::C:T	16	-9.0672e-02	9.1222e-02	3.283e-01
1	982173	G	T	1:982173::G:T,G	14	-1.521e-01	9.9293e-01	1.250e-01
1	983175	C	A	1:983175::C:A	13	1.8374e-01	4.1175e-02	1.176e-02
1	983285	T	C	1:983285::T:C	16	-1.2357e-01	7.4183e-02	9.577e-02
1	983352	G	A	1:983352::G:A	16	-9.7613e-02	8.8222e-02	2.265e-01
1	983510	A	G	1:983510::A:G	16	-1.3837e-01	7.5025e-02	8.226e-02
1	983536	A	T	1:983536::A:T	16	-1.2747e-01	7.7711e-02	1.099e-01
1	983551	A	C	1:983551::A:C	13	1.0780e-02	4.9737e-02	8.284e-01
1	983723	A	G	1:983723::A:G	13	7.4167e-02	4.8956e-02	7.082e-02
1	984081	T	C	1:984081::T:C	16	-1.3112e-01	7.4747e-02	7.939e-02
1	984115	G	T	1:984115::G:T	16	-1.2833e-01	7.6793e-02	9.470e-02
1	984148	A	G	1:984148::A:G	16	-1.4801e-01	7.4546e-02	5.982e-02
1	984149	T	G	1:984149::T:G	16	-1.3968e-01	7.4557e-02	6.101e-02
1	985373	T	C	1:985373::T:C	17	2.8055e-02	2.5669e-02	2.577e-01
1	985769	C	G	1:985769::C:G	13	7.6952e-02	4.4086e-02	8.098e-02
1	986633	T	G	1:986633::T:G	13	7.3357e-02	4.5711e-02	1.086e-01
1	98894	G	T	1:98894::G:T	14	8.1297e-02	3.9340e-02	3.878e-02
1	916255	C	T	1:916255::C:T	15	1.1488e-01	4.1834e-02	6.033e-03
1	916558	G	A	1:916558::G:A	15	8.4447e-02	3.8855e-02	2.648e-01

COVID19-hg GWAS round 5 release (GRCh38)

# 快速入門

上課資料

1. 下載 ANACONDA: <https://www.anaconda.com>

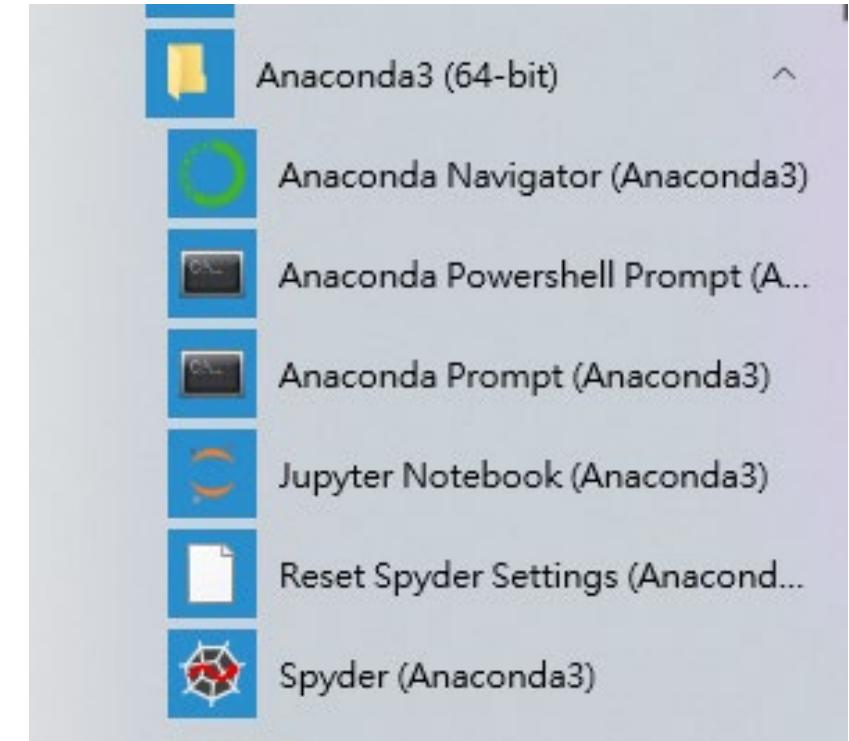


2. 安裝

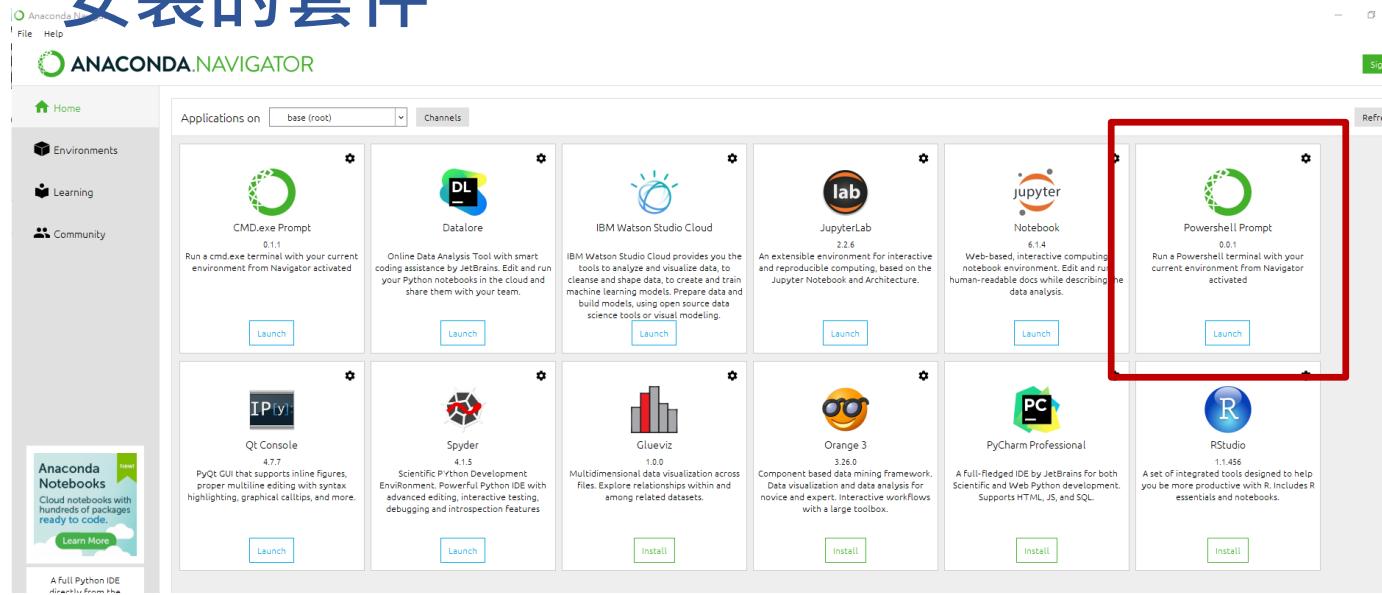
3. 打開ANACONDA

4. 準備工作: 安裝套件 ← 只需要做一次

5. 跑程式 (Jupyter)



# 安裝的套件



打開“安裝frameworks.txt”，一行一行輸入需要安裝的套件

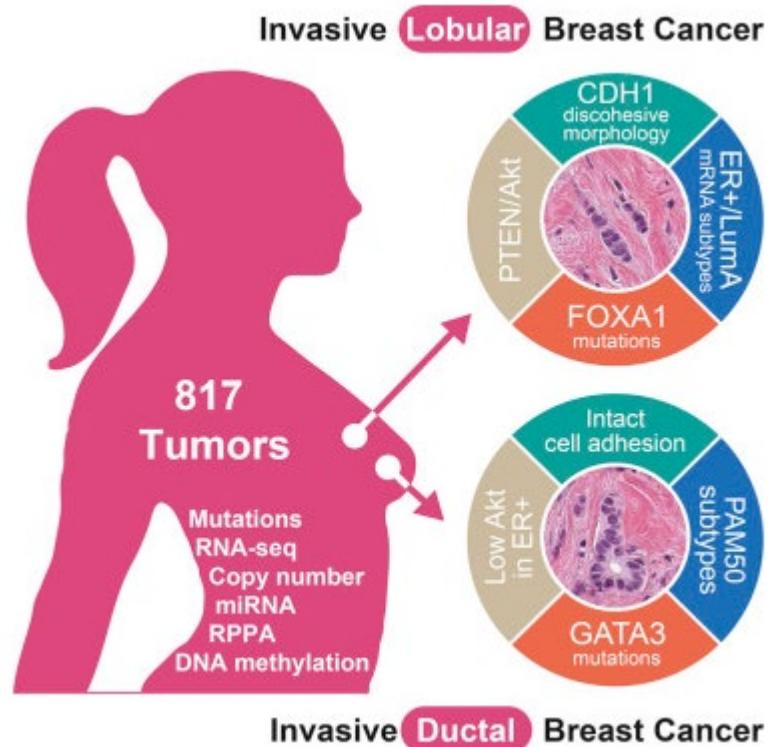
```
安装frameworks - 記事本
檔案(F) 編輯(E) 格式(O) 檢視(V) 說明
pip install dataframe
pip install IPython
pip install sklearn
pip install seaborn
pip install missingno
pip install featexp
pip install mlxtend
pip install pydotplus
pip install xgboost
pip install joblib
conda install numpy
conda install matplotlib
conda install pandas
conda install ipython notebook
```

系統管理員: C:\Windows\System32\WindowsPowerShell\v1.0\powershell.exe  
(base) PS C:\Users\090503> pip install dataframe

# 實戰演練

檔案: <https://www.kaggle.com/datasets/samdemharter/brca-multiomics-tcga>

paper: [https://www.cell.com/cell/fulltext/S0092-8674\(15\)01195-2](https://www.cell.com/cell/fulltext/S0092-8674(15)01195-2)



## Multi-Omics

Copy Number Variations (860)  
Somatic Mutations (249)  
Gene Expression (604)  
Protein Expression (223)  
Total: 1936 features

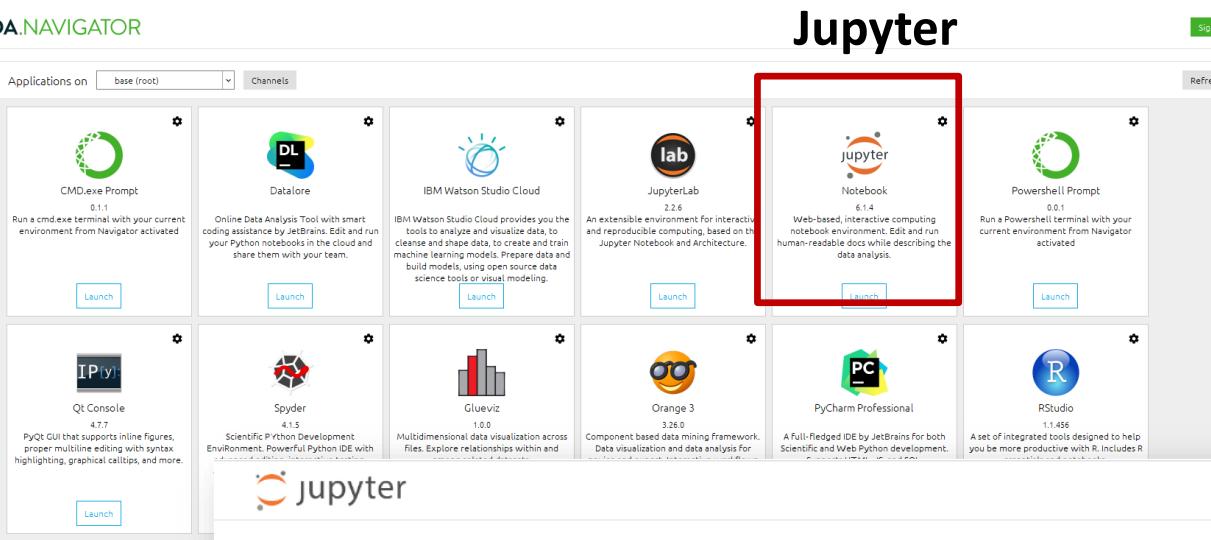
**Total: 1936 features**

n=705 (611名存活 · 94名死亡)

705 patients

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	BVL	BVM
1	rs_CLEC3	rs_CPB1	rs_SCGB2	rs_SCGB11	rs_TFF1	rs_MUCL1	rs_GSTM1	rs_PIP	rsADIPO	rs_ADH1B	rs_S100A7	rs_HMGCS	rs_CYP2B1	rs_ANKRD1	rs_PRAME	lpp_p90RSI	vital.status
2	0.892818	6.580103	14.12367	10.6065	13.18924	6.649466	10.52033	10.33849	10.24838	10.22997	0	7.904609	8.754801	9.681176	9.363258	-0.20726	0
3	0	3.691311	17.11609	15.51723	9.867616	9.691667	8.179522	7.911723	1.289598	1.818891	0.444879	1.289598	5.698288	7.885523	1.818891	0.26753	0
4	3.74815	4.375255	9.658123	5.326983	12.10954	11.64431	10.51733	5.114925	11.97535	11.91144	4.126056	9.383379	8.630913	9.61031	4.733674	-0.19852	0
5	0	18.23552	18.53548	14.53358	14.07899	8.91376	10.55746	13.30443	8.205059	9.211476	0	4.347992	7.836265	9.732374	9.069201	-0.0469	0
6	0	4.583724	15.71186	12.80452	8.881669	8.430028	12.96461	6.806517	4.294341	5.385714	0.841893	6.11157	7.701901	10.7592	1.370053	0.037473	0
7	0.602837	1.026658	9.248794	5.095713	0	6.085525	7.773748	11.08078	6.916635	7.935318	8.671148	5.846375	4.665484	0	11.48724	-0.04847	1
8	5.654312	2.842074	8.231955	7.062477	7.285706	8.490481	0	6.823051	5.907011	6.793263	6.287818	2.762306	2.390172	4.593306	8.209576	0.67704	0
9	2.370945	0.430499	0	2.635987	0.761711	2.788999	5.45192	2.552943	7.388579	6.464041	0.430499	1.625925	1.453491	0	11.82103	-0.10164	0
10	0	0.603597	2.21313	1.027933	1.847075	0	5.949234	1.622087	1.027933	0	0	1.847075	1.847075	0	8.118521	-0.10292	0
11	0	8.187287	5.250765	3.058923	10.29397	0	0.493032	0	9.300252	9.766249	0	0.859811	10.32087	10.69306	1.94557	-0.40892	1
12	0	1.878333	5.078392	2.741467	14.30011	14.61199	2.58938	10.27948	10.80134	10.55123	1.878333	8.667065	14.1371	10.77706	1.418352	-0.34292	0
13	0	8.152976	11.38638	8.111426	9.813765	6.457277	1.594644	7.26423	11.55646	11.6801	0	8.013315	11.8963	9.243775	1.594644	0.090501	0
14	12.08081	1.078542	4.287783	1.311852	9.451286	4.367665	9.611614	8.252329	11.80453	11.73297	0.454808	0	10.24658	2.344005	3.725883	-0.09927	0
15	0.54745	3.648235	0.54745	0.54745	0	7.824536	6.152814	0.54745	5.614657	4.474436	6.440336	0	0.54745	0	1.253747	-0.09538	1
16	10.22	10.13828	14.25146	10.51583	8.020775	12.02788	0.740021	10.50855	8.096916	9.961419	1.880098	8.743535	8.933501	9.745235	0	-0.419	0
17	0	5.03562	6.944828	6.301805	13.35564	5.467071	12.55429	17.15364	5.082013	5.478677	2.193141	7.069034	16.35991	11.34807	0	-0.17159	0
18	1.225028	7.223793	15.18962	12.26103	11.92687	13.07427	2.942646	12.181	8.253096	8.42015	2.764792	8.853269	5.837464	7.359244	4.269115	0.429118	0
19	0	12.18525	10.23754	8.880304	9.218761	9.478781	9.841524	15.88566	0	0	5.641419	3.537768	10.56868	9.668254	9.091519	0.241038	1
20	5.207627	6.388463	7.419122	4.816287	3.759348	4.83754	11.86855	3.714103	4.367986	4.899475	3.408087	4.338481	8.238421	0	2.11919	-0.38334	0
21	0	0.818605	1.337597	0.818605	8.2703	11.07211	8.985609	1.87688	11.8184	10.92404	0	1.101314	8.133799	8.391719	10.60697	0.024336	1
22	4.376457	4.162299	11.51267	10.10907	7.922094	11.97035	10.18525	12.85716	4.041988	7.48328	1.587701	7.839681	6.939862	9.395491	1.443182	-0.38335	0
23	5.272363	5.58896	5.424925	2.927915	12.06274	9.128513	12.65949	9.305378	2.814612	1.590626	12.65178	3.082754	2.408984	7.026013	9.361294	-0.04657	0
24	12.6757	1.38648	10.08193	4.788534	7.36861	14.90865	1.935497	9.742503	3.181389	0	11.00965	10.61633	1.774671	0	4.458566	0.016026	0
25	0	0	0.497638	0.497638	0	0	11.38833	0	0	0	11.44249	0.497638	0.497638	0	9.782566	-0.50348	0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
706	1.687374	0.530171	15.89707	12.01882	13.90623	14.52932	11.95344	8.899572	9.476798	9.451396	3.649374	9.933222	10.60729	9.804425	1.687374	-0.37637	0

## Jupyter



The screenshot shows the Jupyter Notebook interface with a file list. The 'data' folder contains the following files:

Name	Last Modified	File size
stage1_solution_filtered.csv	6年前	8.46 kB
stage2_sample_submission.csv	6年前	21.6 kB
stage2_test_text.csv	6年前	60.6 MB
stage2_test_variants.csv	6年前	16.3 kB
stage_2_private_solution.csv	6年前	2.8 kB
test_text	5年前	24.8 MB
test_text.zip	5年前	104 MB
test_variants	5年前	7.14 kB
test_variants.zip	5年前	48.6 kB
training_text	5年前	212 MB
<input checked="" type="checkbox"/> training_text.zip	5年前	62.8 MB
training_variants	5年前	66.7 kB
training_variants.zip	5年前	24.8 kB

開啟brca.ipynb



# Prediction of Breast Cancer Severity with Multi-Omics Data

導入必要的套件

```
In [1]: import pandas as pd
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split, StratifiedKFold, cross_val_score
from sklearn.metrics import roc_curve, roc_auc_score
import matplotlib.pyplot as plt
from joblib import dump
```

步驟1：載入數據

```
In [2]: data_path = 'data.csv'
df = pd.read_csv(data_path)
```

步驟2：將資料分為特徵和目標變量

```
In [3]: X = df.drop("vital.status", axis=1)
y = df["vital.status"]
```

步驟3：根據目標變量進行層次化分割，將資料分為訓練集和測試集

```
In [4]: X_train, X_test, y_train, y_test = train_test_split(X, y, stratify=y, train_size=0.85, random_state=40)
```



#### 步驟4: 訓練隨機森林分類器

```
In [5]: rf = RandomForestClassifier(random_state=30)  
rf.fit(X_train, y_train)
```

```
Out[5]: RandomForestClassifier  
RandomForestClassifier(random_state=30)
```

#### 步驟5: 在測試集上進行預測

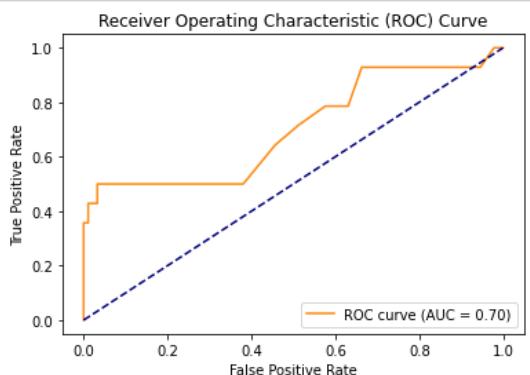
```
In [6]: y_pred_prob = rf.predict_proba(X_test)[:, 1]
```

#### 步驟6: 計算 ROC 曲線和 AUC (曲線下面積)

```
In [7]: fpr, tpr, _ = roc_curve(y_test, y_pred_prob)  
roc_auc = roc_auc_score(y_test, y_pred_prob)
```

#### 步驟7: 繪製 ROC 曲線

```
In [8]: plt.figure()  
plt.plot(fpr, tpr, color='darkorange', label=f'ROC curve (AUC = {roc_auc:.2f})')  
plt.plot([0, 1], [0, 1], color='navy', linestyle='--')  
plt.xlabel('False Positive Rate')  
plt.ylabel('True Positive Rate')  
plt.title('Receiver Operating Characteristic (ROC) Curve')  
plt.legend(loc="lower right")  
plt.show()
```



## 步驟8: 使用 K-Fold 交叉驗證進行模型驗證

```
In [9]: cv = StratifiedKFold(n_splits=10)
scores = cross_val_score(RandomForestClassifier(random_state=42), X, y, cv=cv, scoring="accuracy")
print(f"Accuracy: {scores.mean():.4f} (+/- {scores.std() * 2:.4f})")

Accuracy: 0.8809 (+/- 0.0563)
```

## 步驟9: 重要特徵排序

```
In [10]: importance = pd.Series(data=rf.feature_importances_, index=X_train.columns)
importance = importance.sort_values(ascending=False)
print(importance)

rs_KIAA0408      0.010571
rs_MMRN1         0.010038
rs_ADIPQO        0.005953
rs_SLC19A3        0.005876
rs_KCNIP2         0.005841
...
cn_AMZ1           0.000000
cn_FRMD1          0.000000
cn_SYT5            0.000000
cn_SIRPG           0.000000
cn_SLC1A2          0.000000
Length: 1936, dtype: float64
```

## Step 10: 輸出訓練好的模型

```
In [11]: dump(rf, 'random_forest_model.joblib')

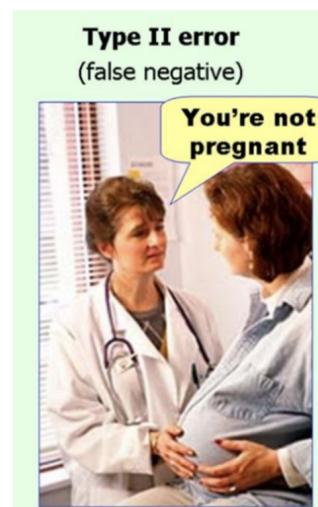
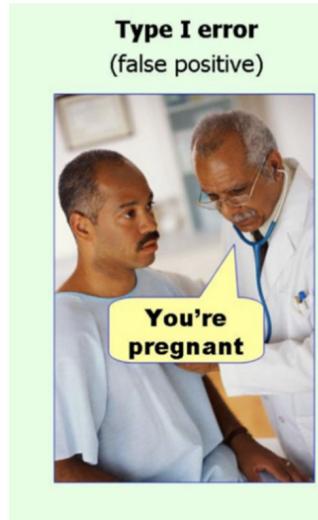
Out[11]: ['random_forest_model.joblib']
```

# 效能評估



# confusion matrix 混淆矩陣

		預測
		Positive      Negative
實際	Positive	True Positive (TP)
	Negative	False Negative (FN) 偽陰性
Negative	False Positive (FP)	True Negative (TN) 偽陽性



偽陽性率，沒病被診斷有病的比率

false positive rate (FPR)

$$\frac{FP}{N} = \frac{FP}{FP + TN}$$

偽陰性率，有病被診斷沒病的比率

false negative rate (FNR)

$$\frac{FN}{P} = \frac{FN}{FN + TP}$$

**Recall/Sensitivity** 
$$\frac{TP}{P} = \frac{TP}{TP + FN}$$
 所有得病者被正確診斷出得病的比率

**Specificity** 
$$\frac{TN}{N} = \frac{TN}{TN + FP}$$
 所有非得病者被正確診斷非得病的比率

**Precision** 
$$\frac{TP}{TP + FP}$$
 診斷出得病者中，真正得病的比率

**False Discovery Rate (FDR)** 
$$\frac{FP}{FP + TP} = 1 - Precision$$

**Accuracy** 
$$\frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$$
 所有受試者被正確診斷為有病和沒病的比率

# Harmonic Precision-Recall Mean (F1 Score)

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

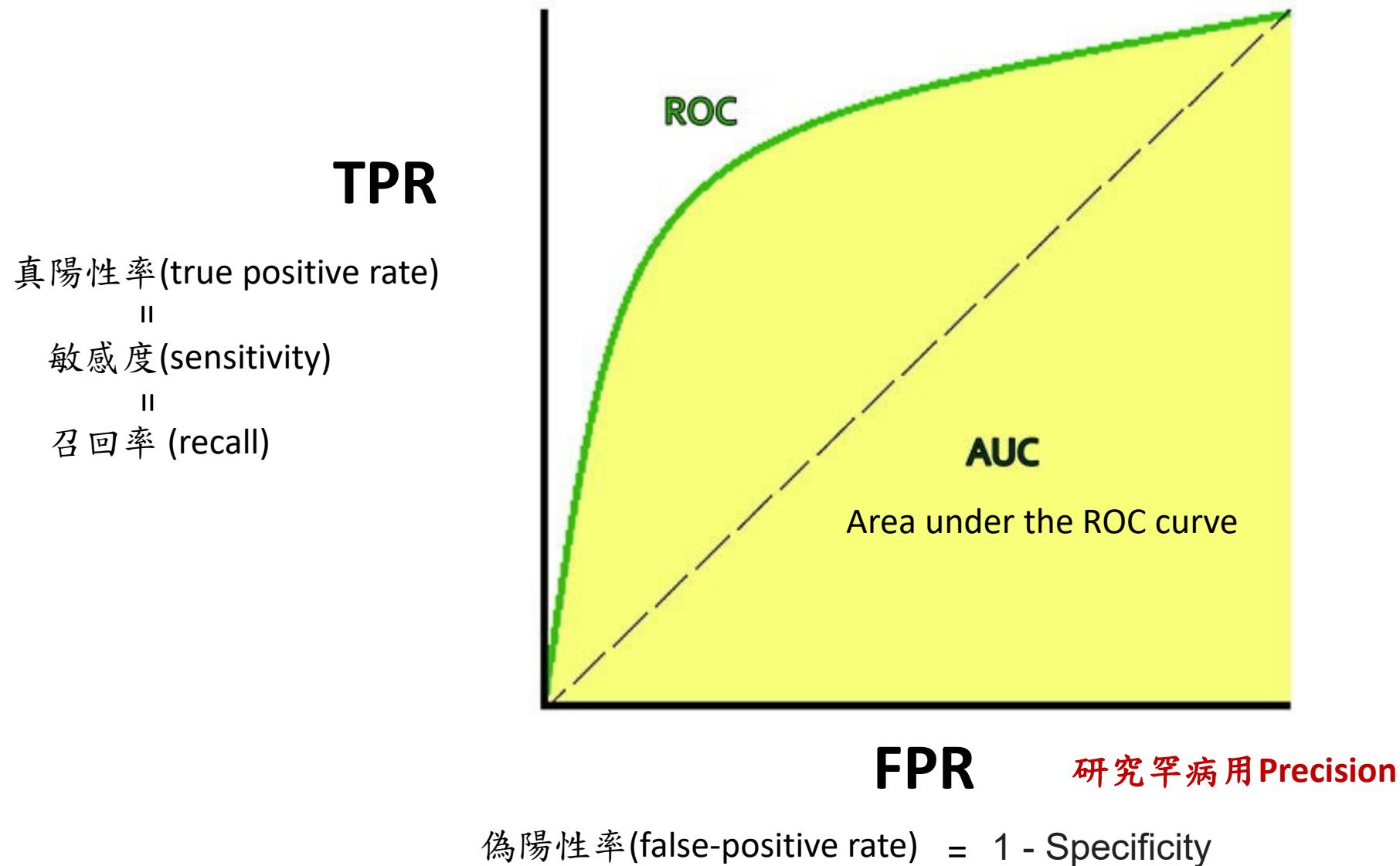
診斷出得病者中，真正得病的比率

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

所有得病者被正確診斷出得病的比率

$$F_1 = 2 * \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

# Receiver operating characteristic curve, ROC



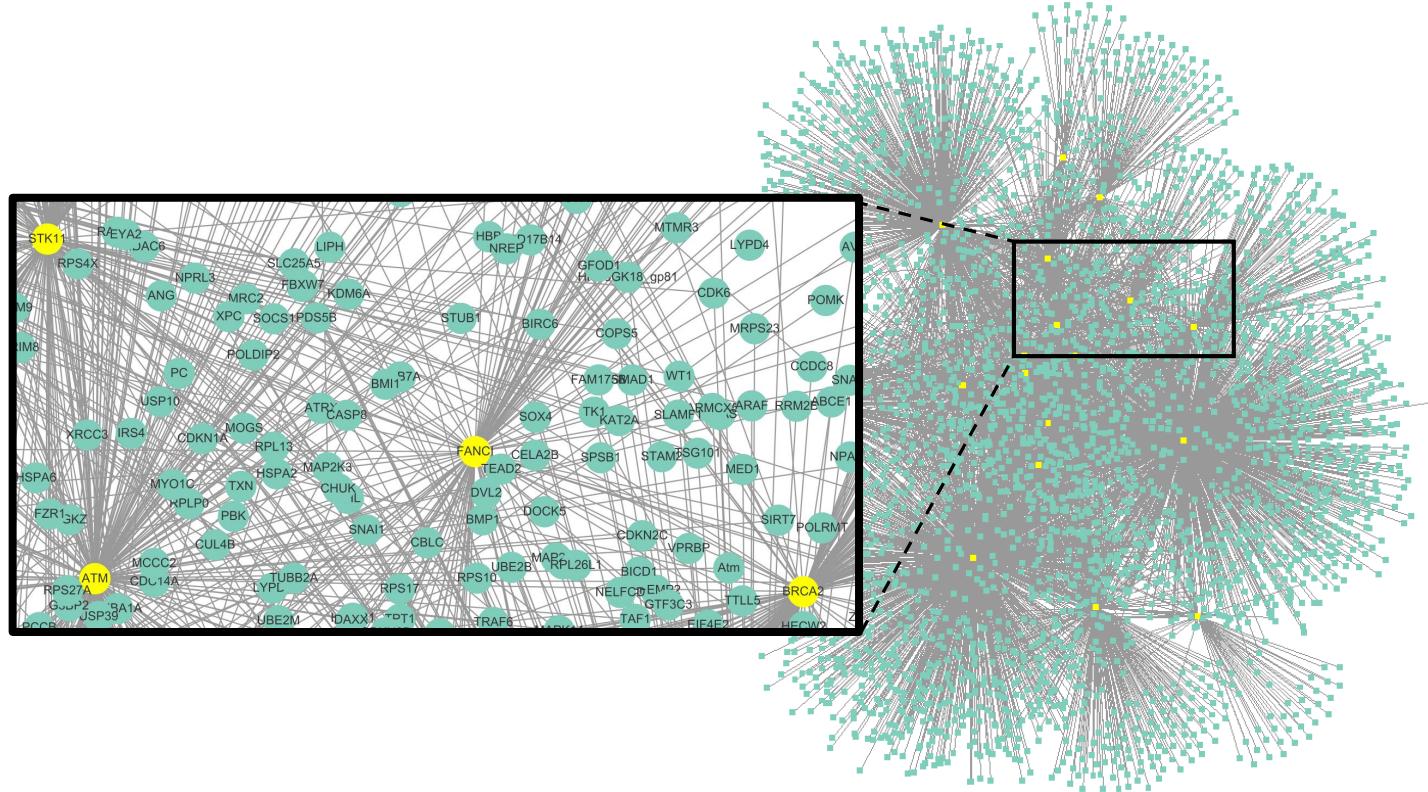
# 與遺傳性乳癌相關的基因

14%遺傳性乳癌患者找不到相關基因

---

<b>ATM</b>	ATM serine/threonine kinase
<b>BARD1</b>	BRCA1 associated RING domain 1
<b>BRCA1</b>	BRCA1 DNA repair associated
<b>BRCA2</b>	BRCA2 DNA repair associated
<b>BRIP1</b>	BRCA1 interacting helicase 1
<b>CDH1</b>	cadherin 1
<b>CHEK2</b>	checkpoint kinase 2
<b>FANCA</b>	FA complementation group A
<b>FANCI</b>	FA complementation group I
<b>FANCL</b>	FA complementation group L
<b>NBN</b>	nibrin
<b>NF1</b>	neurofibromin 1
<b>PALB2</b>	partner and localizer of BRCA2
<b>PIK3CA</b>	phosphatidylinositol-4,5-bisphosphate 3-kinase catalytic subunit alpha
<b>PMS2</b>	homolog 2, mismatch repair system component
<b>PPM1D</b>	protein phosphatase, Mg <sup>2+</sup> /Mn <sup>2+</sup> dependent 1D
<b>PTEN</b>	phosphatase and tensin homolog
<b>RAD51C</b>	RAD51 paralog C
<b>STK11</b>	serine/threonine kinase 11
<b>TP53</b>	tumor protein p53

---

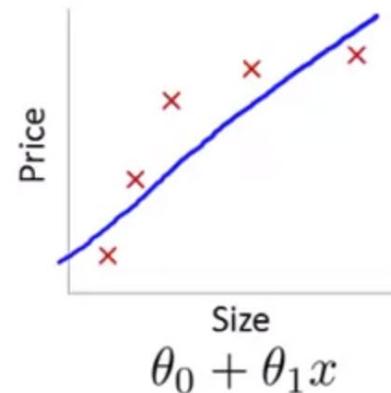


僅有10%乳癌患者有以上基因突變

5%乳癌患者有BRCA1/2突變

# 欠擬合(Underfitting) & 過擬合(overfitting)

欠擬合

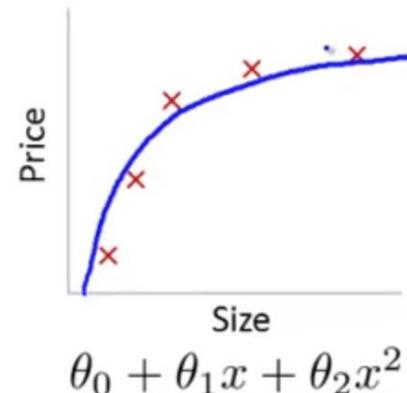


High bias  
(underfit)

特徵數/參數太少

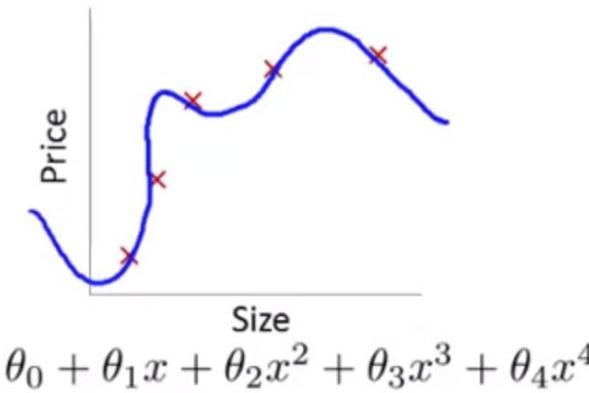
單純以遺傳性乳癌相關基因  
建立乳癌預測模型

完美



“Just right”

過擬合



High variance  
(overfit)

樣本數太少

特徵數/參數太多

# 從人的睡姿來預測床的形狀

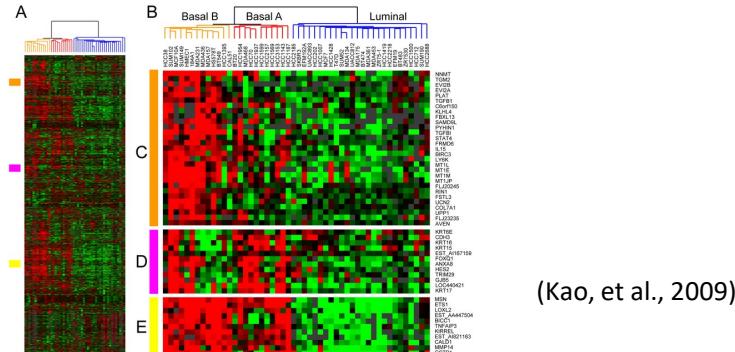


# Unsupervised Learning

無監督學習

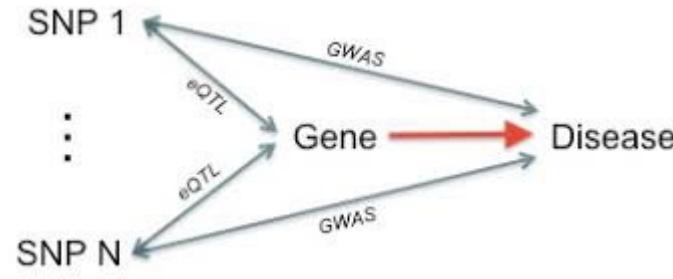
## Clustering (群集)

Microarray 基因表現分群



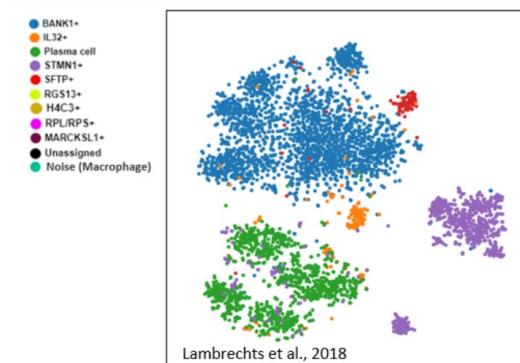
## Association (關聯)

GWAS (全基因組關聯分析)



## Dimensionality Reduce (降維)

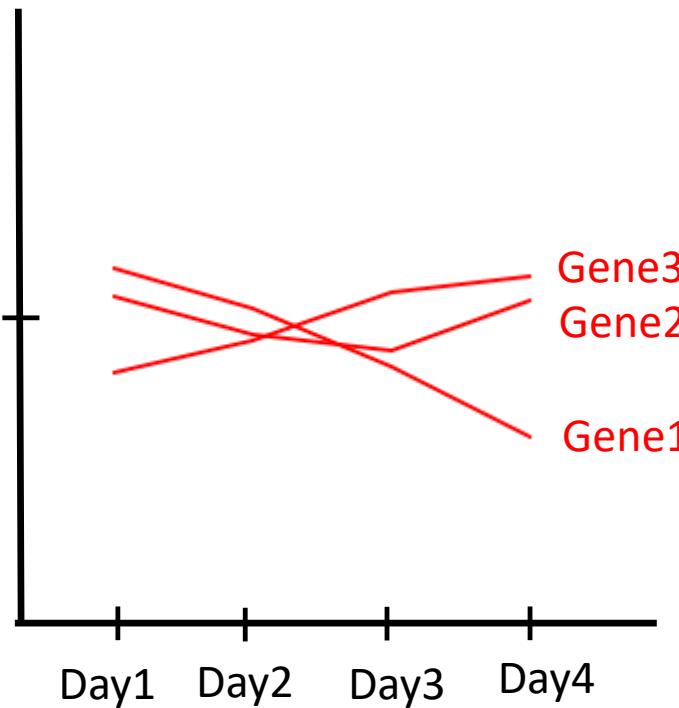
Single-cell analysis



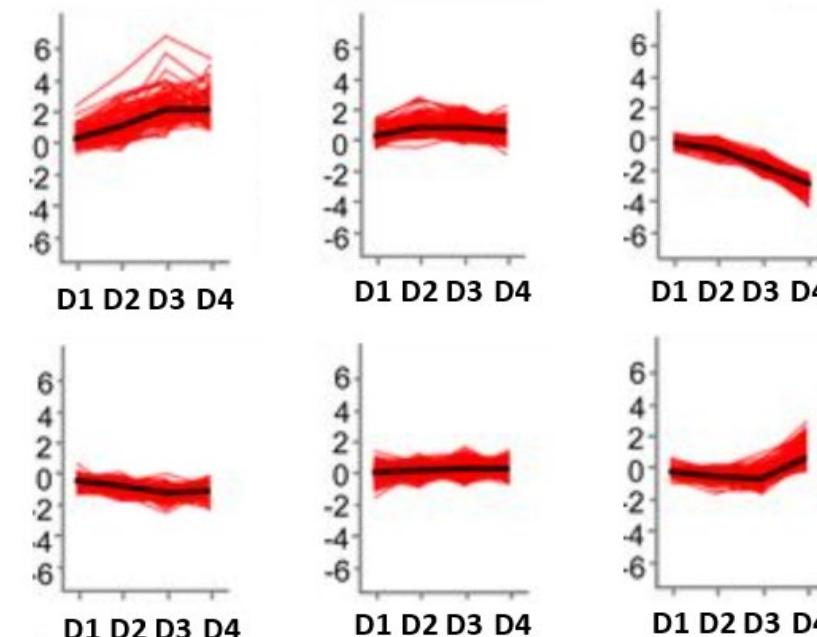
施打疫苗後基因產生的變化

## Clustering (群集)

Microarray/RNA-seq 基因表現分群



k-Means  
Clustering  
k=6



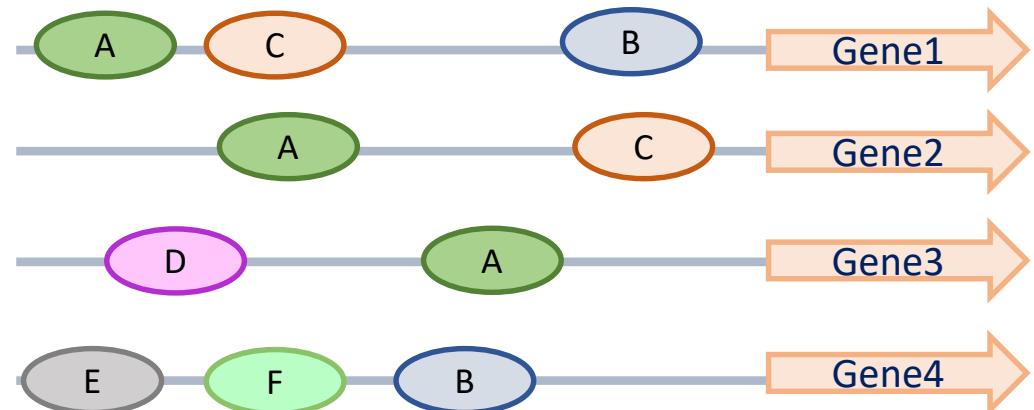
Day0: Control

Gene Functional  
Analysis

# 找出轉錄因子(Transcription factors)的組合

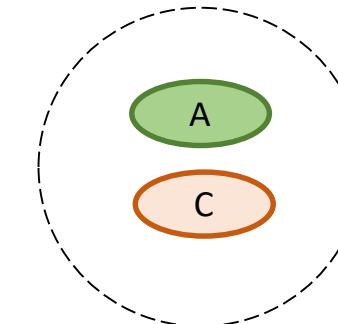
# Association (關聯)

## promoter



Apriori algorithm

Transcription factors (TF)



Gene	TF
Gene1	A,B,C
Gene2	A,C
Gene3	A,D
Gene4	B,E,F

計算頻率

Frequent Items	Support
{A}	0.75
{B}	0.5
{C}	0.5
{A,C}	0.5
{E}	0.25
{F}	0.25
{A,B,C}	0.25
{A,D}	0.25
{B,E,F}	0.25

Minimum support 0.5  
Minimum confidence 0.5

For site combination  $A \Rightarrow C$ :

$$\text{Support} = \text{support}(\{A, C\}) = 0.5$$

$$\text{Confidence} = \text{support}(\{A, C\})/\text{support}(\{A\}) = 0.67$$

$$\text{Confidence} = \text{support}(\{A, C\})/\text{support}(\{C\}) = 1$$

**RESEARCH****Open Access**

# Functional analysis of transcription factor binding sites in human promoters

Troy W Whiteld<sup>1</sup>, Jie Wang<sup>1</sup>, Patrick J Collins<sup>2</sup>, E Christopher Partridge<sup>3</sup>, Shelley Force Aldred<sup>2</sup>, Nathan D Trinklein<sup>2</sup>, Richard M Myers<sup>3</sup> and Zhiping Weng<sup>1\*</sup>

*Genome Biology* 2012

Transcription factor binding resulted in transcriptional repression in more than a third of functional sites.

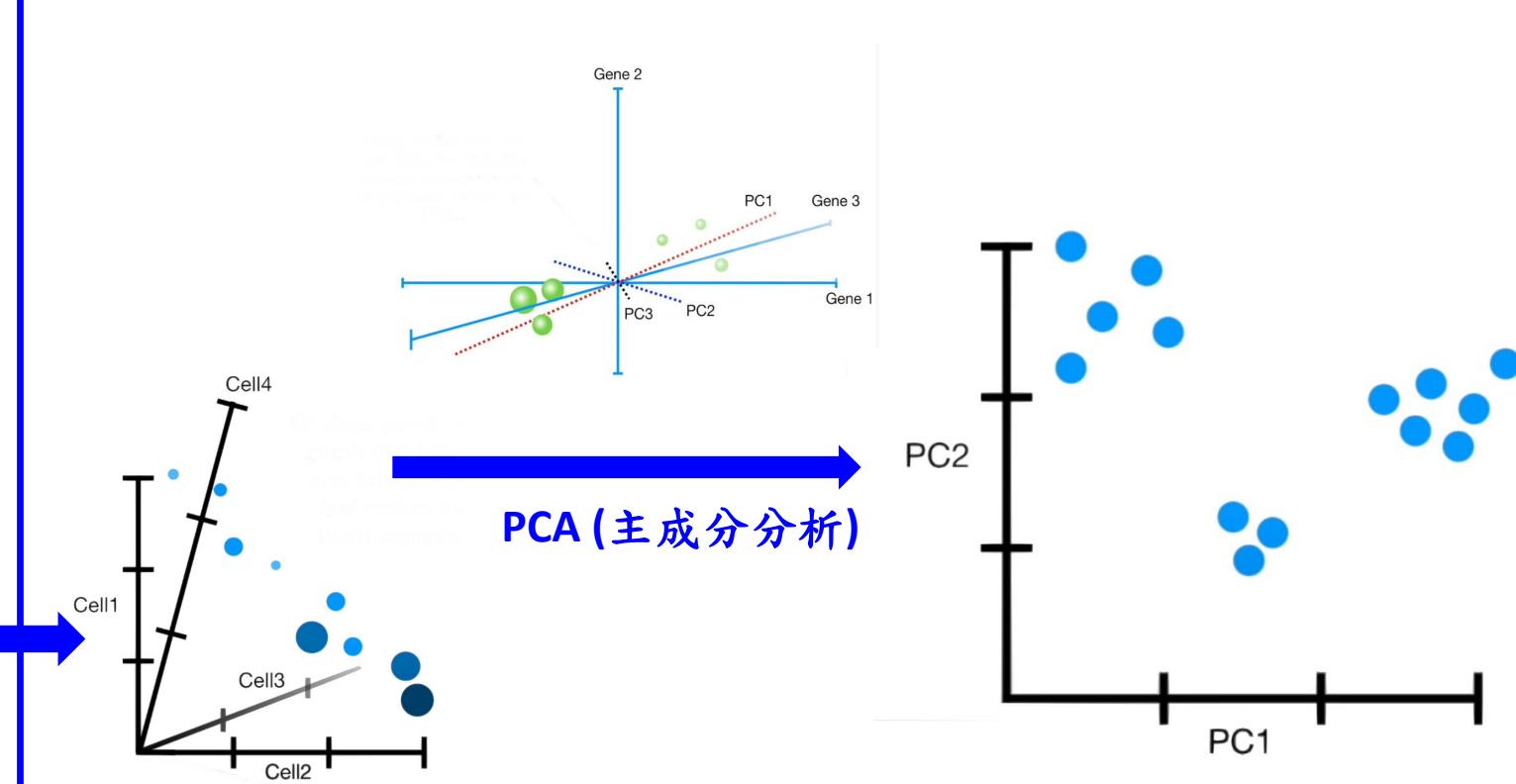
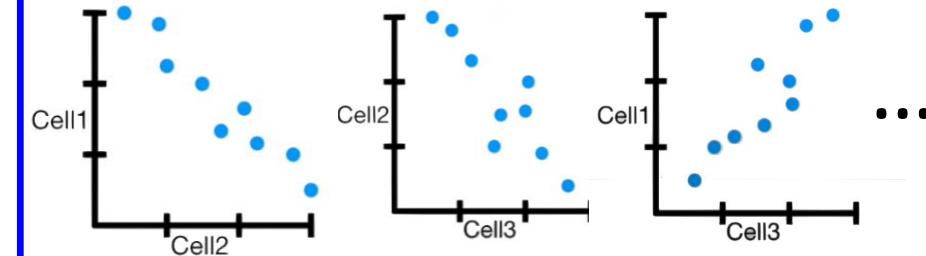
The activating TF binding sites are significantly closer than repressing TF binding sites to the TSS

Positive elements on human promoters between 40 and 350 bp away from the TSS, as well as the presence of negative elements from 350 to 1,000 bp upstream of the TSS (Cooper et al., *Genome Res*, 2006).

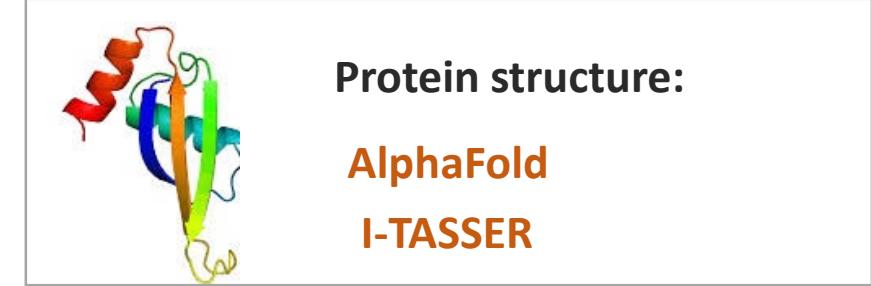
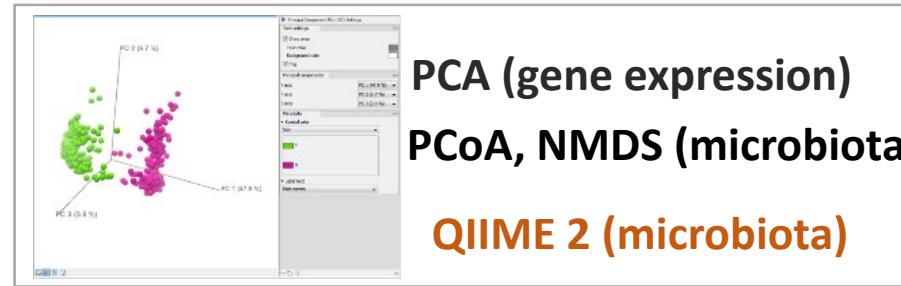
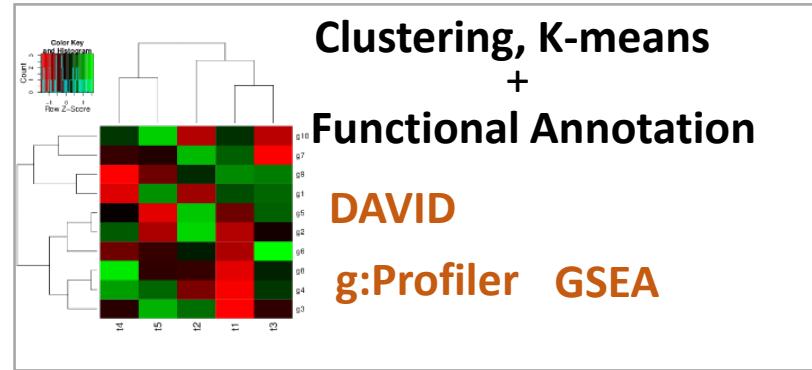
# 利用基因表現值將細胞分類

# Dimensionality Reduce (降維)

	Cell1	Cell2	Cell3	Cell4	...
Gene1	3	0.25	2.8	0.1	...
Gene2	2.9	0.8	2.2	1.8	...
Gene3	2.2	1	1.5	3.2	...
Gene4	2	1.4	2	0.3	...
Gene5	1.3	1.6	1.6	0	...
Gene6	1.5	2	2.1	3	...
Gene7	1.1	2.2	1.2	2.8	...
Gene8	1	2.7	0.9	0.3	...
Gene9	0.4	3	0.6	0.1	...



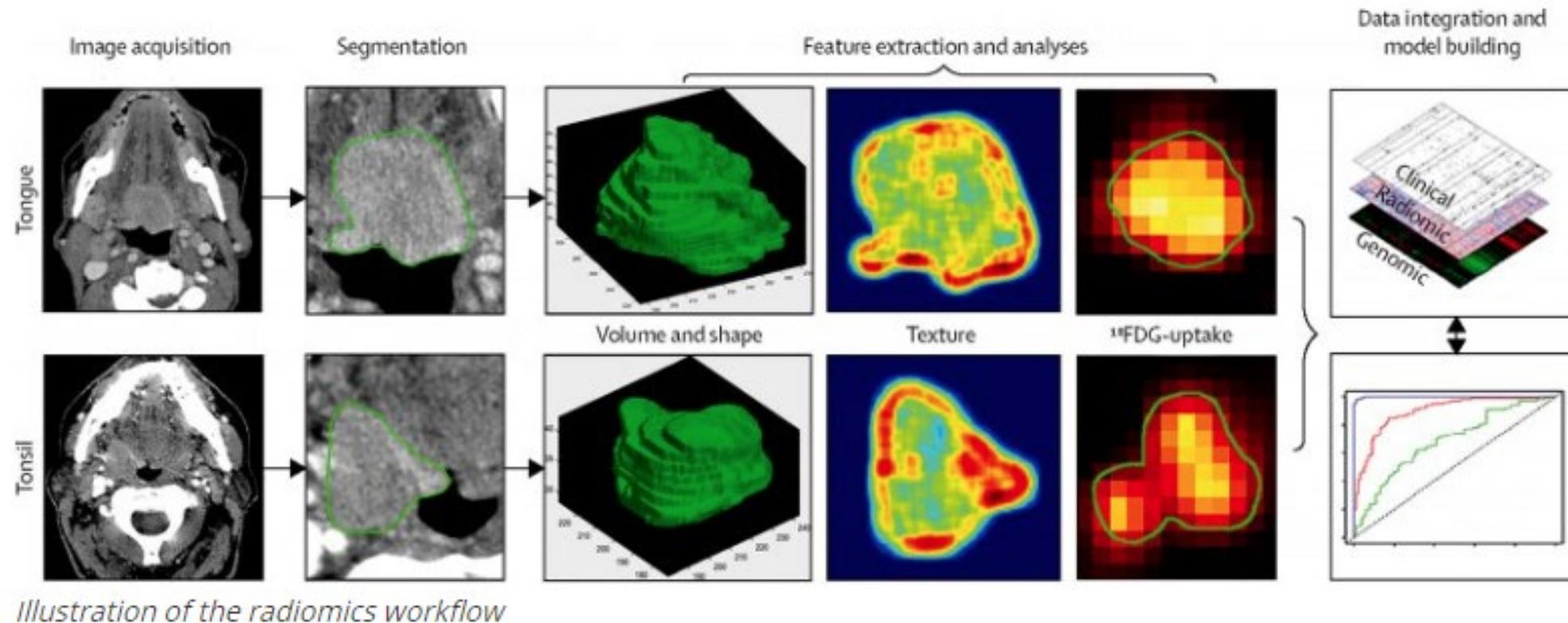
機器學習早  
已融入到生物資訊  
各個分析領域



# 深度學習 Deep Learning

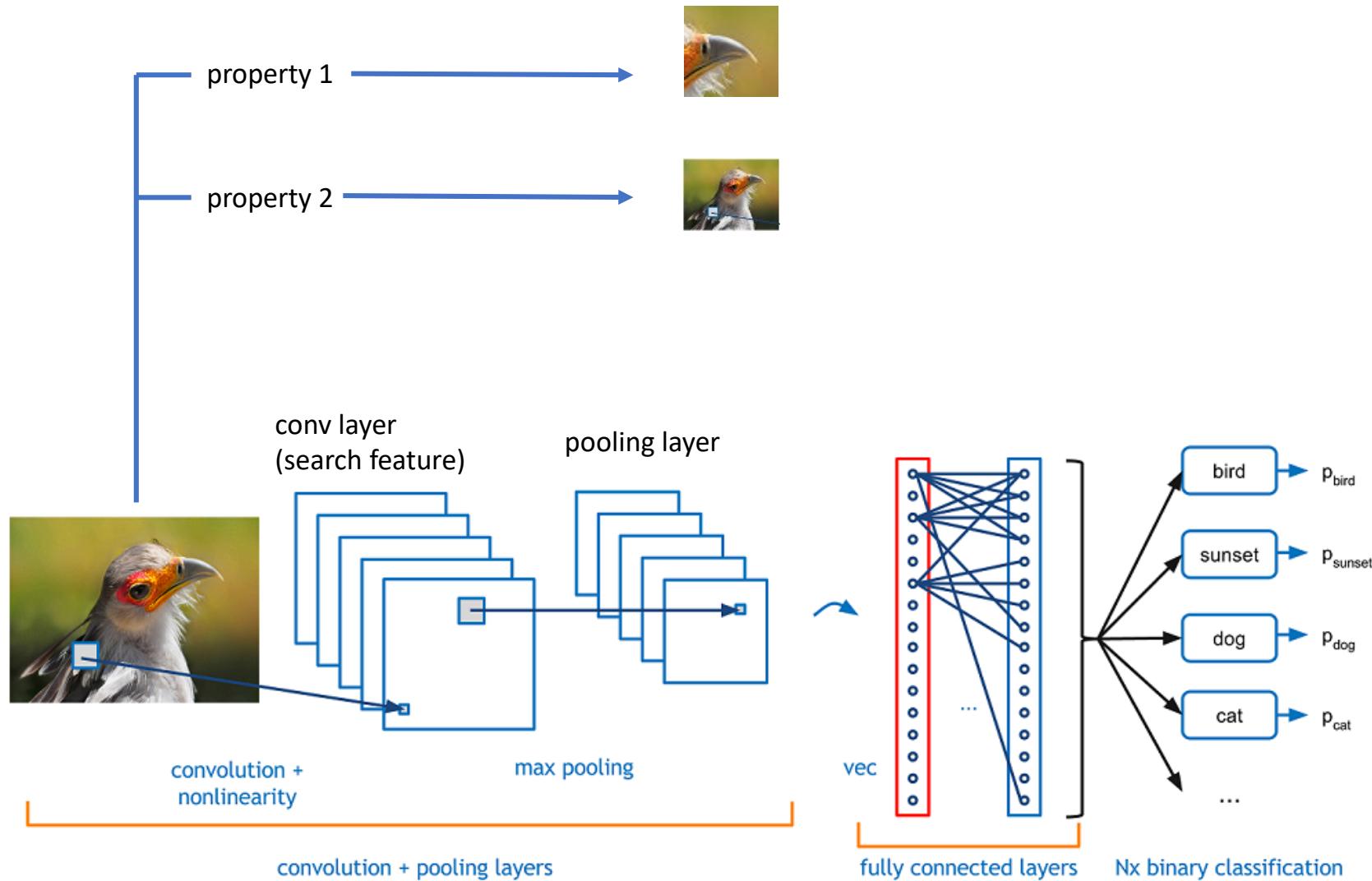
## Radiomics 放射組學

### 生物醫學影像分析流程



Philippe Lambin 2019

# Convolutional Neural Networks



## BMJ Open Diagnostic accuracy of X-ray versus CT in COVID-19: a propensity-matched database study

Aditya Borakati  <sup>1,2</sup>, Adrian Perera, <sup>2</sup> James Johnson, <sup>2</sup> Tara Sood <sup>2</sup>

To cite: Borakati A, Perera A, Johnson J, et al. Diagnostic accuracy of X-ray versus CT in COVID-19: a propensity-matched database study. *BMJ Open* 2020;10:e042946. doi:10.1136/bmjopen-2020-042946

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2020-042946>).

Received 20 July 2020  
Revised 07 October 2020  
Accepted 08 October 2020

### ABSTRACT

**Objectives** To identify the diagnostic accuracy of common imaging modalities, chest X-ray (CXR) and CT, for diagnosis of COVID-19 in the general emergency population in the UK and to find the association between imaging features and outcomes in these patients.

**Design** Retrospective analysis of electronic patient records.

**Setting** Tertiary academic health science centre and designated centre for high consequence infectious diseases in London, UK.

**Participants** 1198 patients who attended the emergency department with paired reverse transcriptase PCR (RT-PCR) swabs for SARS-CoV-2 and CXR between 16 March and 16 April 2020.

**Main outcome measures** Sensitivity and specificity of CXR and CT for diagnosis of COVID-19 using the British Society of Thoracic Imaging reporting templates.

Reference standard was any RT-PCR positive nasopharyngeal swab within 30 days of attendance. ORs of CXR in association with vital signs, laboratory values and 30-day outcomes were calculated.

**Results** Sensitivity and specificity of CXR for COVID-19

### Strengths and limitations of this study

- Large, appropriately powered, study population consisting of all patients attending the emergency department rather than those solely with confirmed COVID-19; this allowed assessment of specificity for the imaging modalities and applicability to the general population who may attend medical personnel with other complaints, but have underlying SARS-CoV-2 infection.
- Comprehensive statistical analyses were conducted to address confounding in reporting of X-rays including propensity score matching and logistic regression to give a 'doubly robust' model.
- Low amount of missing data and for secondary covariates only; multiple imputation was performed with a good fit, however, observed data would be preferable to imputed data.
- Single centre, retrospective study; potential for inter-reporter and intercentre variability in reporting.
- Large proportion of patients excluded due to not having a reverse transcriptase PCR swab, predominantly those with previous reported no negative test.

**Results** Sensitivity and specificity of CXR for COVID-19 diagnosis were 0.56 (95% CI 0.51 to 0.60) and 0.60 (95% CI 0.54 to 0.65), respectively. For CT scans, these were 0.85 (95% CI 0.79 to 0.90) and 0.50 (95% CI 0.41 to 0.60),

# 靠人工判讀X光片無法正確診斷是否確診新冠肺炎

只從X光片，確診者被診斷出確診的比率

Sensitivity = 0.56

只從X光片，非確診者被診斷非確診的比率

Specificity = 0.60

©  2020 The Author(s). *BMJ Open* published by BMJ Publishing Group Ltd on behalf of BMJournals Ltd. This is an open access article distributed under the terms of the [Creative Commons Attribution Non-Commercial-ShareAlike license](https://creativecommons.org/licenses/by-nc-sa/4.0/), which permits unrestricted use, distribution, and reproduction in other forms, provided the original author(s) and publisher are credited and a link is made to the published article on *bmjopen* and the full terms of the license are followed.

<sup>1</sup>Division of Surgery and Interventional Science, University College London, London, UK

<sup>2</sup>Emergency Department, Royal Free Hospital, London, UK

Correspondence to  
Dr Aditya Borakati;  
a.borakati@doctors.org.uk

# 影像分析: 胸腔X光判讀感染肺炎

正常(n = 1,349)



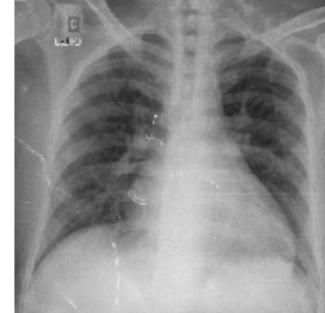
非COVID-19病毒性肺炎(n=1,345)



細菌性肺炎(n=2,538)



COVID-19肺炎(n=3616)



資料來源: Kang Zhang et. al., 2018, *Cell*  
Chowdhury et. al., 2020. *IEEE*

開發框架:

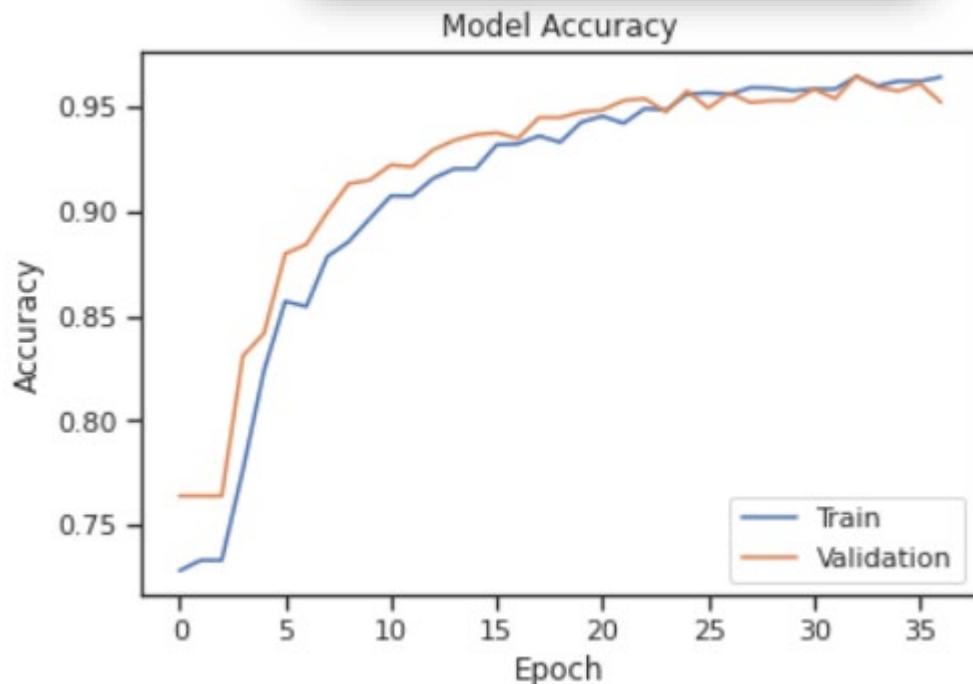
- TensorFlow-21.02-tf2-py3

硬體設定:

- a100.2g.20gb

# Convolutional neural networks(CNN)

```
Model: "sequential_1"
Layer (type)      Output Shape       Param #
=====
conv2d_55 (Conv2D)    (None, 68, 68, 128)   3584
max_pooling2d_2 (MaxPooling2D) (None, 34, 34, 128) 0
dropout_3 (Dropout)    (None, 34, 34, 128) 0
conv2d_56 (Conv2D)    (None, 32, 32, 64)    73792
max_pooling2d_3 (MaxPooling2D) (None, 16, 16, 64) 0
dropout_4 (Dropout)    (None, 16, 16, 64) 0
conv2d_57 (Conv2D)    (None, 14, 14, 64)    36928
flatten_1 (Flatten)   (None, 12544)        0
dense_2 (Dense)      (None, 16)           200720
dropout_5 (Dropout)   (None, 16)           0
dense_3 (Dense)      (None, 2)            34
=====
Total params: 315,058
Trainable params: 315,058
Non-trainable params: 0
=====
CPU times: user 39.5 ms, sys: 4.12 ms, total: 43.6 ms
Wall time: 42.6 ms
```



13808 samples

全部運行完畢時間: 約4min 30s

開發框架:

- TensorFlow-21.02-tf2-py3

硬體設定:

- a100.2g.20gb

只從X光片，確診者被診斷出確診的比率

**Sensitivity = 0.93**

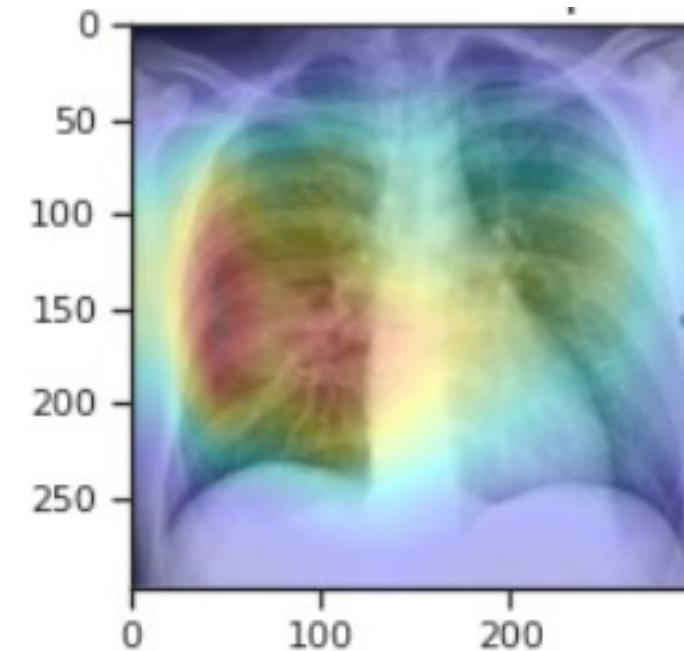
只從X光片，非確診者被診斷非確診的比率

**Specificity = 0.96**

# 找出AI判讀X光片為確診的依據



X-ray



Grad - CAM Algorithm

偏紅色為AI重要的參考區域

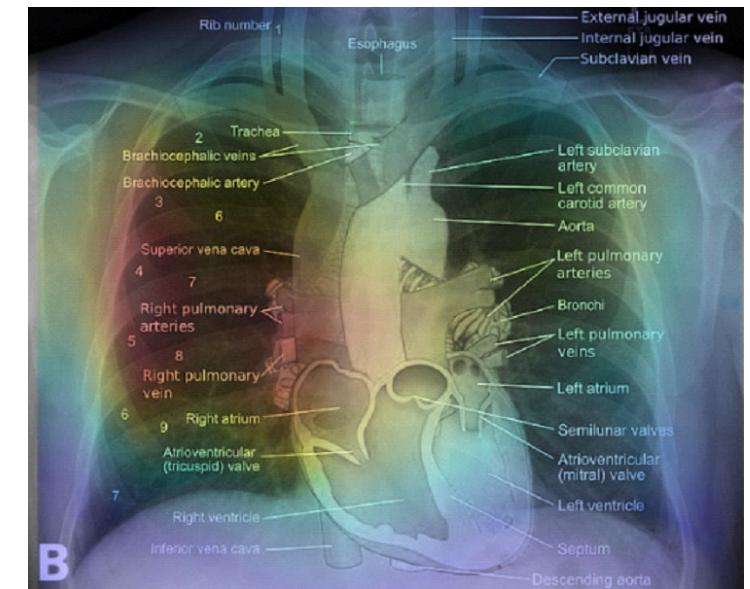
```
root@unlime-service-py: ~ binary_ADO920:pynb detection-of-covid-positi GB_SNP_0308:pynb Untitled3:pynb
+-----+
| / Jupyter / | +-----+
| Name | Last Modified |
| detec... | 41 minutes ago |
| digital_... | 22 days ago |
| digital_... | 22 days ago |
| digital_... | 22 days ago |
| E4_f... | 3 days ago |
| E4_g... | 3 days ago |
| E4.txt | 3 days ago |
| E4corr_... | 3 days ago |
| elephant... | 10 days ago |
| GB_SNP... | 18 hours ago |
| GB_SNP... | 22 days ago |
| manhatt... | 14 days ago |
| Manhatta... | 21 days ago |
| Manhatta... | 21 days ago |
| Manhatta... | 21 days ago |
| nonE4_rf... | 3 days ago |
| nonE4_x... | 3 days ago |
+-----+
[119]: imag.append(cv2.imread(img_path))
imag.append(cv2.imread("./cam.jpg"))

for i in range(len(img_path)):
    save_and_display_graddam(img_path[i], covid_noncovid_heatmap[i])

titles_list = ["Cov1 original", "Cov1 heatmap", "Cov2 original", "Cov2 heatmap", "Cov3 original", "Cov3 heatmap", "Cov4 original", "Cov4 heatmap"]
plot_multiple_img(img, titles_list, ncols = 4, main_title = "GRAD-CAM COVID-19 Heatmap")
GRAD-CAM COVID-19 Heatmap
```

Heatmap

加上註釋



# 神經網路與深度學習

Neural networks and deep learning

1

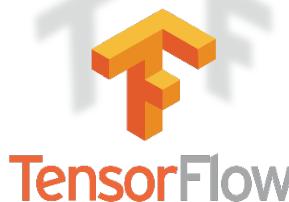


**Python or Julia**

Program language



2



TensorFlow



PyTorch

**TensorFlow or PyTorch**

Software library for machine learning

3



Keras

**Keras**

Deep learning API

# 神經網路在基因體學的應用



## Predicting Splicing from Primary Sequence with Deep Learning

2019

# SpliceAI

nature  
REVIEWS GENETICS

Research Highlight | Published: 25 January 2019

RNA

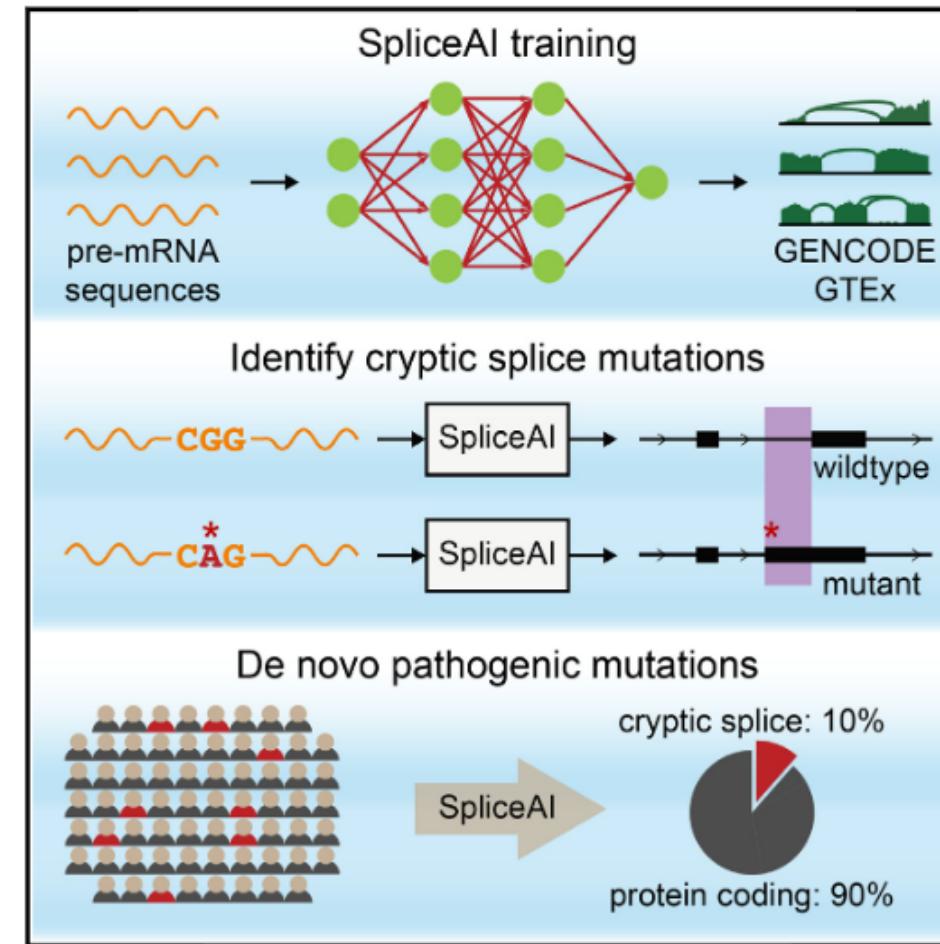
### Learning the language of splicing

Dorothy Clyde

Nature Reviews Genetics 20, 132–133 (2019) | Download Citation

Precise splicing of pre-mRNAs is essential for proper gene function, and defective splicing can lead to disease. However, our understanding of the sequence features that underlie the accuracy of this process remains incomplete, which makes it difficult to identify genetic variation that could disrupt it. Now, a paper in *Cell* describes SpliceAI, a deep neural network that accurately predicts not only splice sites in previously unseen pre-mRNA sequence but also the effects of sequence variants on splicing patterns.

### Graphical Abstract



### Authors

Kishore Jaganathan,  
Sofia Kyriazopoulou Panagiotopoulou,  
Jeremy F. McRae, ..., Serafim Batzoglou,  
Stephan J. Sanders, Kyle Kai-How Farh

### Correspondence

kfarh@illumina.com

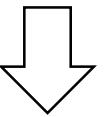
### In Brief

A deep neural network precisely models mRNA splicing from a genomic sequence and accurately predicts noncoding cryptic splice mutations in patients with rare genetic diseases.

# Modeling

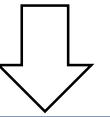
Input: pre-mRNA sequences

CATGCATCTCGA **G** CATGCGTACACT  
40 nt                  40 nt  
200 nt                200 nt  
1000 nt              1000 nt  
5000 nt              5000 nt



Convolutional Neural Networks (CNNs)

Residual neural network (ResNet)



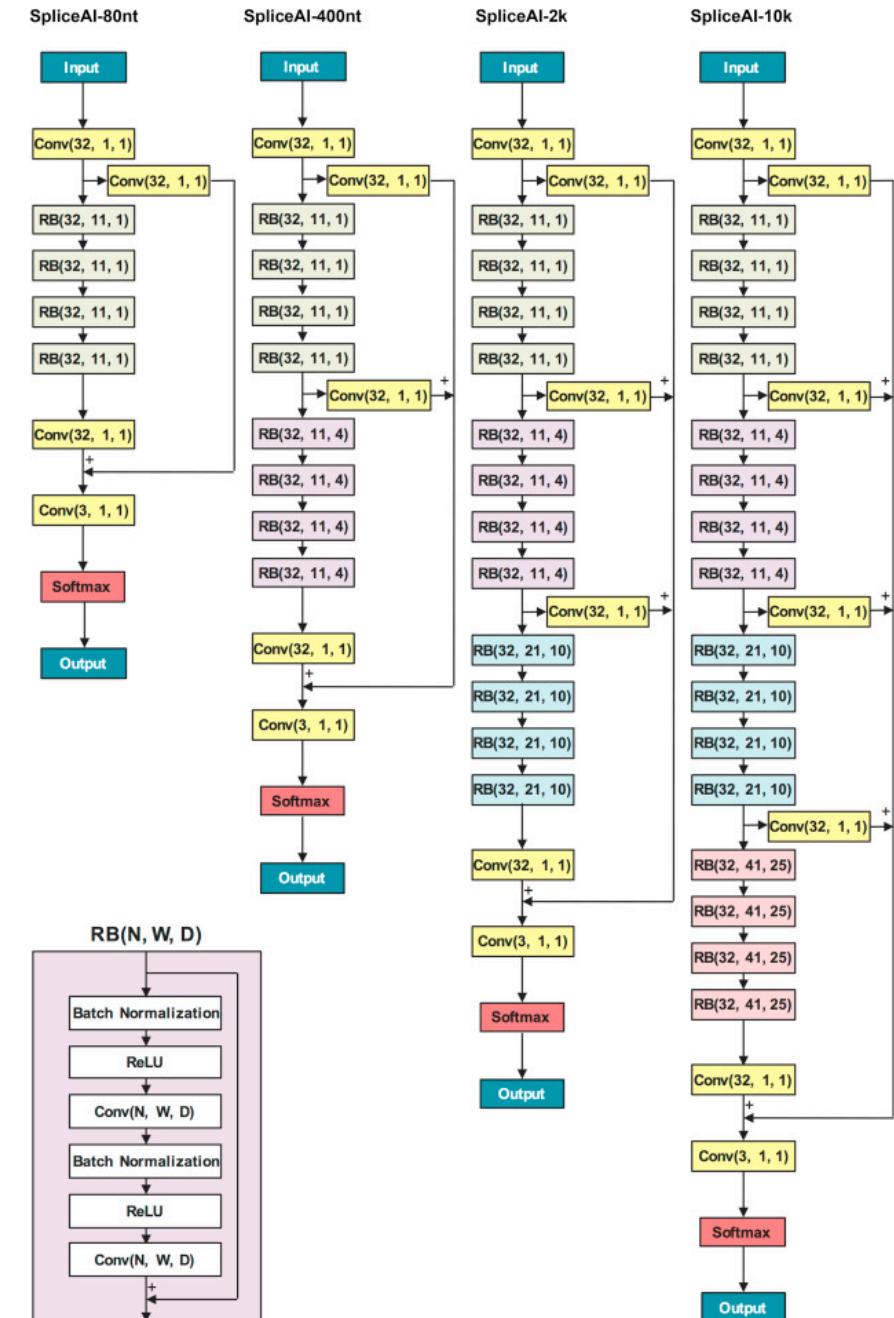
Output: probability

splice acceptor Loss/Gain  
splice donor Loss/Gain  
neither

variant: chr6-2317763-T-A

Acceptor Loss

pre-mRNA position: 106 bp



(Jaganathan et al., Cell, 2019)

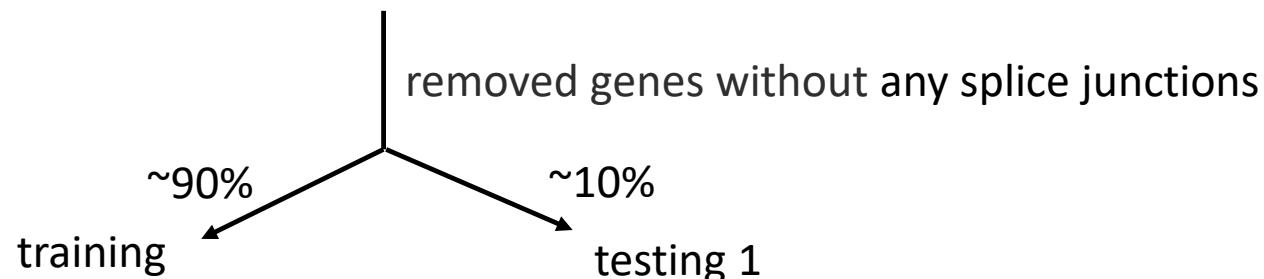
# Training and Testing

- hg19/GRCh37 Genome assembly
- GENCODE Version 24lift37



bin	name	chrom	strand	txStart	txEnd	cdsStart	cdsEnd	exonCount	exonStarts	exonEnds	score	name2	cdsStartStat	cdsEndStat	exonFrames
585	ENST00000473358.1	chr1	+	29553	31097	29553	29553	3	29553,30563,30975,	30039,30667,31097,	0	RP11-34P13.3	none	none	-1,-1,-1,
585	ENST00000469289.1	chr1	+	30266	31109	30266	30266	2	30266,30975,	30667,31109,	0	RP11-34P13.3	none	none	-1,-1,
585	ENST00000417324.1	chr1	-	34553	36081	34553	34553	3	34553,35276,35720,	35174,35481,36081,	0	FAM138A	none	none	-1,-1,-1,
585	ENST00000461467.1	chr1	-	35244	36073	35244	35244	2	35244,35720,	35481,36073,	0	FAM138A	none	none	-1,-1,
585	ENST00000335137.3	chr1	+	69090	70008	69090	70008	1	69090,	70008,	0	OR4F5	cmpl	cmpl	0,
585	ENST00000466430.5	chr1	-	89294	120932	89294	89294	4	89294,92090,112699,120774,	91629,92240,112804,120932,	0	RP11-34P13.7	none	none	-1,-1,-1,-1,
585	ENST00000495576.1	chr1	-	89550	91105	89550	89550	2	89550,90286,	90050,91105,	0	RP11-34P13.8	none	none	-1,-1,
585	ENST00000477740.5	chr1	-	92229	129217	92229	92229	4	92229,112699,120720,129054,	92240,112804,120932,129217,	0	RP11-34P13.7	none	none	-1,-1,-1,-1,
585	ENST00000471248.1	chr1	-	110952	129173	110952	110952	3	110952,112699,129054,	111357,112804,129173,	0	RP11-34P13.7	none	none	-1,-1,-1,
73	ENST00000610542.1	chr1	-	120724	133723	120724	120724	4	120724,120873,129054,133373,	120869,120932,129223,133723,	0	RP11-34P13.7	none	none	-1,-1,-1,-1,

20,287 protein-coding gene



chromosomes 2, 4, 6, 8, 10-22, X, and Y  
(13,384 genes, 130,796 donor-acceptor pairs)

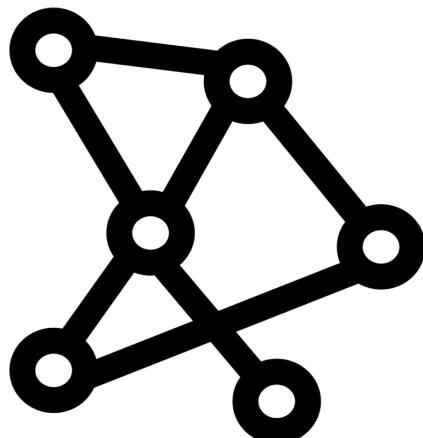
chromosomes 1, 3, 5, 7, and 9 without any paralogs  
(1,652 genes, 14,289 donor-acceptor pairs)

95% accuracy !

正確觀念



正確方法



有效的精準醫療



不預設立場

多面向觀察

邏輯正確

結果符合生物學常識

# 肥胖是一種疾病，是心血管疾病與糖尿病元凶？

Dr. Jaime Guevara-Aguirre

Dr. Valter Longo



厄瓜多小人  
(西班牙後裔)

一般厄瓜多人

肥胖  
Obese

20%

12%

糖尿病  
Diabetes

~0%

5%

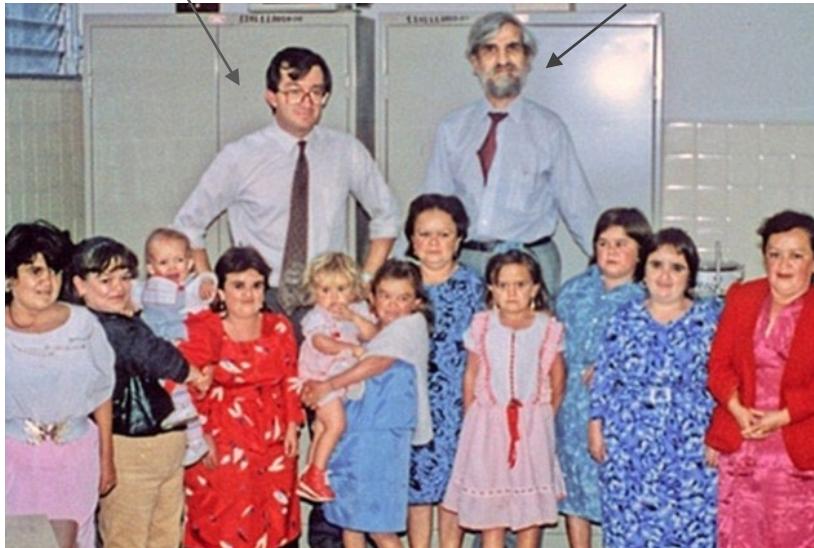
癌症  
Cancer

~0%

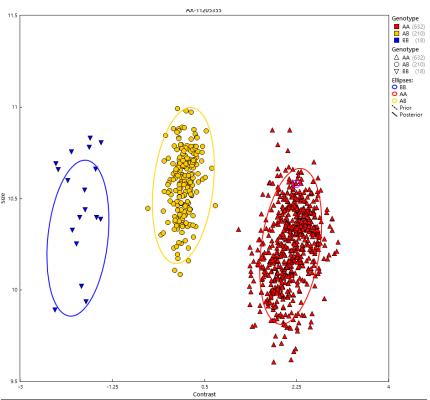
22%

Dr. Jaime Guevara-Aguirre

Dr. Arlan Rosenbloom



## Genotyping



基因體學

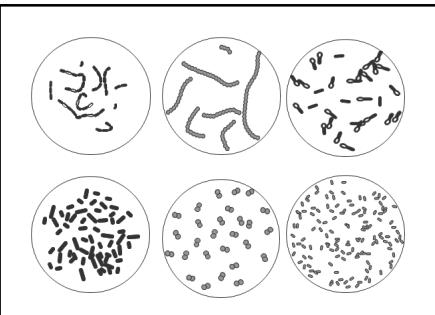
代謝體學

腸道菌體學

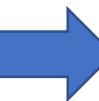
## Metabolites data

Fat No	C0	C2	C3	C4	C5	C5-L-DC	C7-DC	C8	C10	C12	C14	C16	C18	C18-L	C18-2	
C000400	-0.28348	0.17467	-0.12124	0.075949	0.24419	0.05205034	0.059697	0.269461	0.06987	0.96384	0.254183	0.747937	1.082462	-0.11301	0.18447	
C000402	0.297065	0.10545	0.118549	0.812814	0.296573	0.054008	0.01294	-0.21853	-0.35443	-0.05934	-0.04474	0.258312	0.009848	0.422455		
C000407	-0.23575	0.09711	-0.50371	-0.33498	0.54643	0.09746284	0.029967	0.18133	-0.05914	-0.39142	-0.639	-0.94986	1.28911	-0.84997	0.06568	
C000427	-0.10524	0.446385	-0.50371	-0.64448	-0.53305	0.05205034	0.057485	0.146841	0.18464	0.498521	0.789405	0.477322	0.244804	0.020266		
C000435	0.532347	0.31154	0.544424	0.30872	0.73282	0.05205034	0.040854	0.04599	0.05689	0.083256	0.17433	0.204992	-0.21299	-0.44668	0.000097	
C000444	-0.12487	0.342307	0.49344	-0.57525	-0.49881	0.107456284	0.040193	0.012824	-0.02026	0.096676	0.081563	0.304272	0.304855	0.17722	0.065971	
C000446	0.17075	0.39668	0.361338	-0.16506	0.317485	0.28010791	0.042005	0.29166	-0.14622	0.014236	0.047141	0.010971	-0.3505	-0.07030	-0.06075	
C000448	0.173947	0.78189	0.171136	-0.07805	0.152	0.1925034	0.07372	0.91254	0.86197	-0.06379	-0.74582	-0.35769	-0.38702	-0.68282	-0.09212	
C000455	0.17079	0.09293	0.20935	0.27703	0.253757	0.05205034	0.154723	0.31195	0.40195	0.083256	0.303948	0	0	0	-0.11301	
C000460	0.31044	0.65546	0.37761	0.153651	-0.20044	0.05205034	0.04995	0.05357	-0.05333	-0.02833	-0.2014	0.17931	-0.11784	-0.42164	0.144975	
C000474	0.23376	0.05538	1.047724	0.25763	0.04372	0.029967	-0.05962	-0.11451	-0.21527	0.047141	0.085516	0.304855	0.019629	0.029444		
C000479	-0.37455	0.05995	0.49344	-0.61451	0.736977	0.107456284	0.016474	-0.14202	-0.03519	-0.07325	0.254183	0.057995	-0.62930	0.019629	0.033382	
C000483	0.251935	0.147166	0.39214	-0.30256	-0.03368	0.387032132	0.731511	0.854423	0.791542	0.49143	-0.13883	0.375968	-0.21299	0.132104	0.245577	
C000484	-0.08586	0.265715	0.31107	-0.23217	-0.57942	0	0.042827	0.567685	0.010185	0.071072	0.049662	0.01918	0.055495	-0.13454	-0.22487	
C000487	0.34544	0.65546	0.69494	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	0.05205034	
C000497	-0.34544	0.643454	0.22699	-0.1173	-0.67647	0.01457	-0.61753784	0.03161	1.280108	1.432215	0.319068	0.026262	-0.26258	-0.14886	-0.11301	-0.066918
C000494	0.186743	-0.1594	0.138494	-0.05395	0.056989	0.074745284	0.042827	-0.05593	-0.14622	0	-0.09576	-0.05164	0.210218	-0.1554	-0.55532	
C000496	-0.04417	-0.0223	0.62177	0.171588	-0.20219	0.05205034	0.040599	1.09311	-0.99422	-0.75132	-0.044626	-0.33398	-0.31487	-0.43494	-0.87199	

## Gut microbiota data



## 機器學習



## 精準醫療

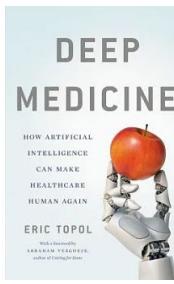


## Metabolites data

性別  
年齡  
人種  
族群（肥胖族群，抽菸族群）



# 參考書籍



Deep Medicine: How Artificial Intelligence Can Make Healthcare Human Again, Topol, Eric, 2019



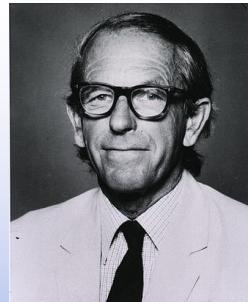
精準醫學：早期預防癌症，破解基因迷思對症下藥，曾欽元，2020

## 最後彩蛋



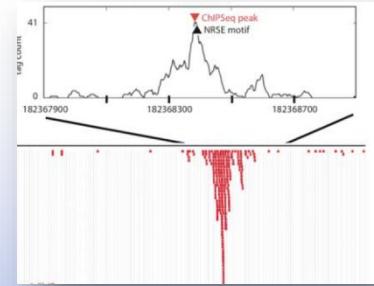
# 1. 人工智慧與定序都是老技術新發展，並開創熱潮

Sequencing



Frederick Sanger  
1975

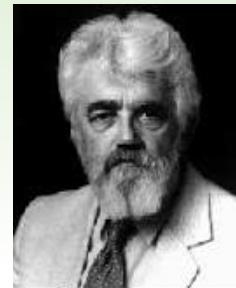
Dideoxynucleotide termination



NGS

(Johnson et al, Science 2007)  
2007 first real application of  
Next Generation Sequencing

AI



John McCarthy  
1956

coined the term "artificial intelligence" (AI)



Geoffrey Hinton  
2006  
Deep learning

Deep  
Learning

## 2. Deep learning與NGS第一次結合是2016年，Google贏得SNP variant (from DNA-seq ) 預測大賽



precisionFDA



HIGHEST  
SNP Performance  
in the precisionFDA Truth Challenge



AWARDED TO

**Verily Life Sciences**

Ryan Poplin Mark DePristo  
Verily Life Sciences Team

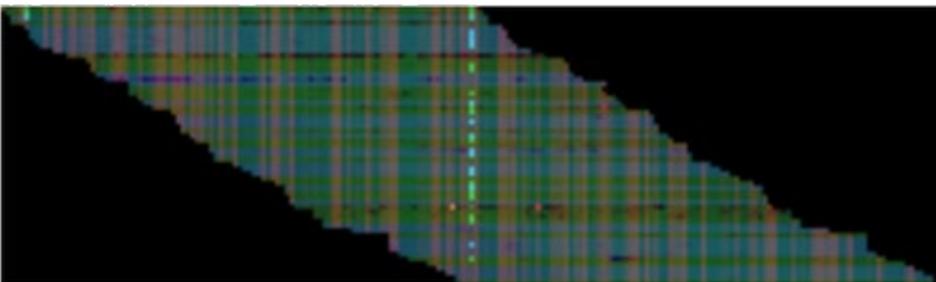


Google  
Google ❤️ Open Source

# DeepVariant

(Poplin et al., Nature Biotechnology, 2018)

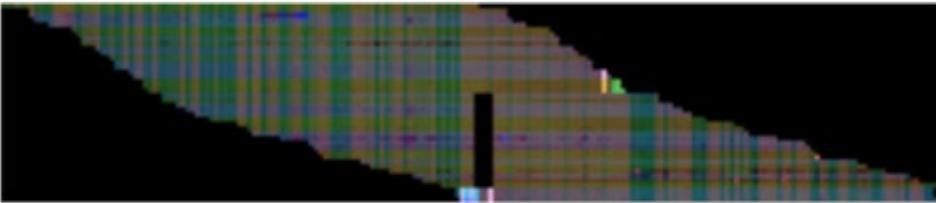
A



**SNP**

*true SNP on one chromosome pair*

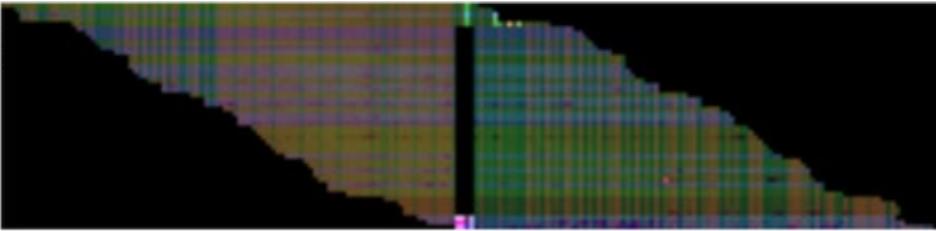
B



**Hetero Deletion**

*a deletion on one chromosome*

C



**Homo Deletion**

*a deletion on both chromosomes*

D

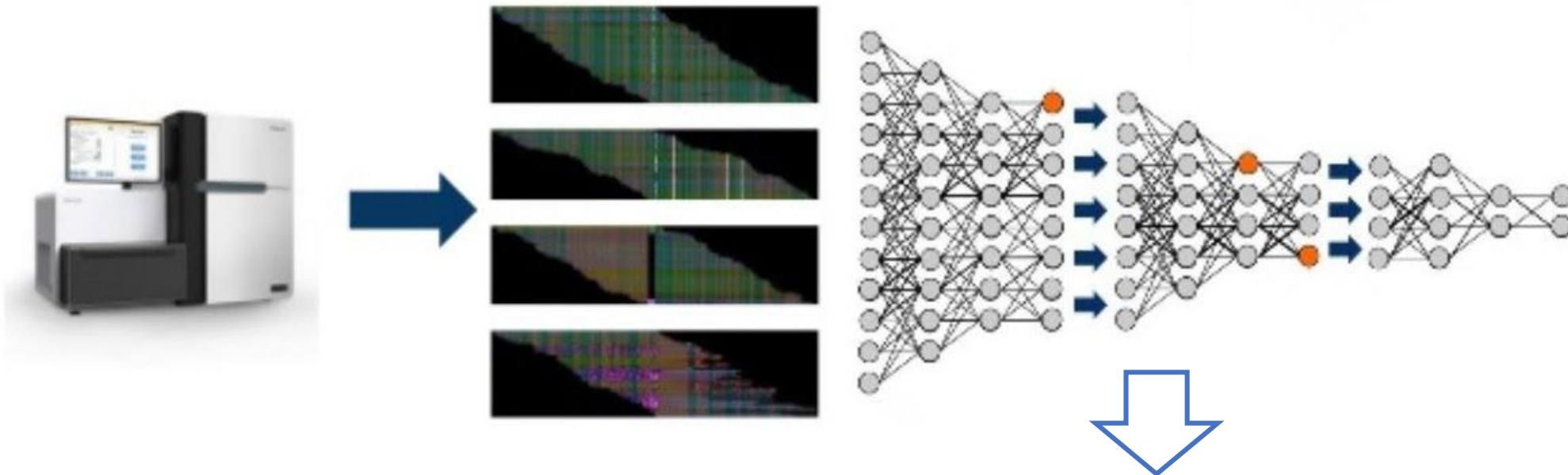


**False Positive**

*a false variant caused by errors*

Red = {ACGT}   Green={quality score}   Blue={read strand}

## Convolutional neural network (CNN)



Genotype

Homozygous reference: 0.01

Homozygous variant: **0.95**

Heterozygous: 0.04

99.9% accuracy

Homozygous variant call

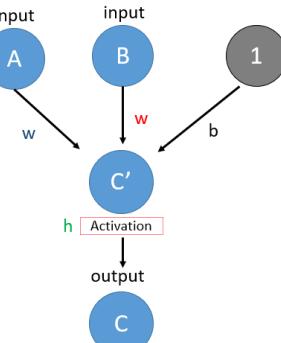
### 3. 1958年，Frank Rosenblatt當初為了解開大腦的奧秘，模仿神經元設計了感知器，是人工智慧的先驅，也是AI神經網絡最小單位



Frank Rosenblatt  
(American psychologist)



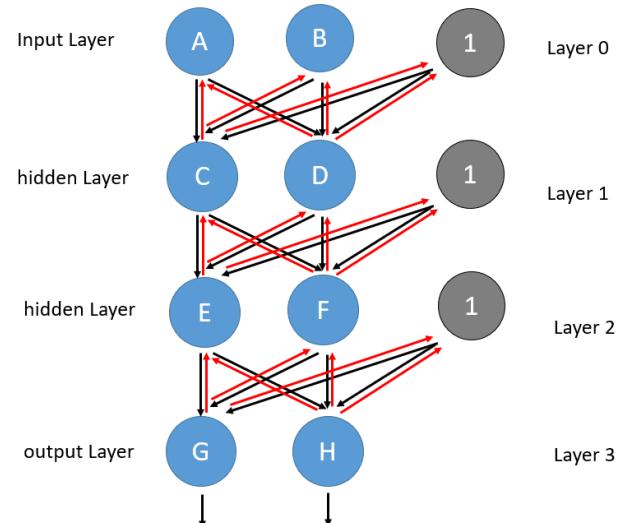
#### 感知器 Perceptron



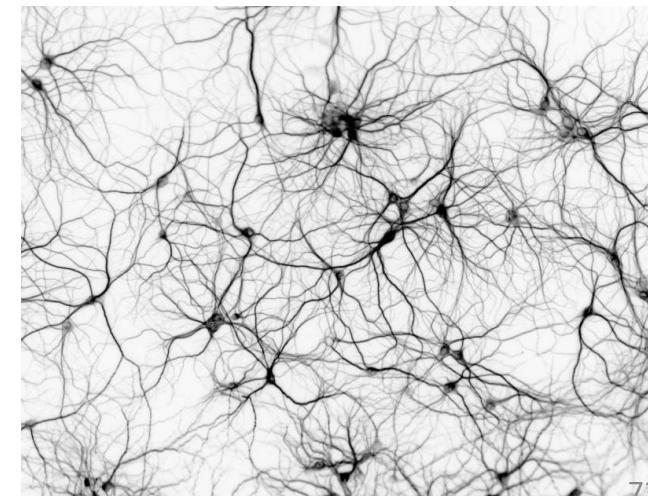
#### 神經元 Neuron



#### AI神經網絡 Neural Network



#### 神經網路 Nerve Net



[nature](#) > [articles](#) > [article](#)

[Published: 01 January 1966](#)

## Transfer of Conditioned Responses from Trained Rats to Untrained Rats by Means of a Brain Extract

[FRANK ROSENBLATT](#), [JOHN T. FARROW](#) & [WILLIAM F. HERBLIN](#)

[Nature](#) **209**, 46–48 (1966) | [Cite this article](#)

103 Accesses | 43 Citations | [Metrics](#)

**Thanks for your attention~**

My Email: paul@nhri.edu.tw