

逢 甲 大 學

資訊工程學系

專題研究報告

簡易皮膚癌檢測系統

指導教授:林明言 教授

學生：戴劭心 D0593174

戴家齊 D0773186

中 華 民 國 一 百 零 九 年 十 一 月

誌謝

首先誠摯的感謝指導教授林明言教授，在專題上提供我們意見。也提供我們機器來訓練模型以及架設伺服器，讓我們能夠順利完成專題。

摘要

不論是誰，身上或多或少都會有一些痣，而幾乎99%以上的痣都是正常的，對於那1%如果是惡性的黑色素瘤，算是致死率相當高，也是相當容易轉移的癌症。在台灣邁入高齡社會的如今，罹患皮膚癌的人數也日益漸長，大部分的皮膚癌發生在50歲以後[1]，罹病的部位以臉部、頸部、耳朵、上手臂、四肢等易暴露在陽光下的地方為多。

而皮膚相較於內臟，是可以直接觀察的，也就是說對於皮膚癌，我們可以使用肉眼對其做事前的判斷。也就代表著，我們可以利用圖片辨識來對皮膚癌進行初步的檢查。所以我們利用 opencv 對圖片整理後，使用 CNN 卷積神經網路進行深度學習，最後將之架到網站上。來輔助民眾來做出是否進行更進一步檢查的判斷。

關鍵詞: 影像處理，卷積神經網路，羅吉斯回歸，k 鄰近演算法，支援向量機，決策樹

目錄

誌謝.....	i
摘要.....	ii
目錄.....	iii
圖目錄.....	v
表目錄.....	vi
第一章 緒論.....	1
1.1 研究背景.....	1
1.2 研究動機.....	2
1.3 研究目的.....	2
第二章 文獻回顧.....	3
2.1 CNN 卷積神經網絡.....	3
2.2 羅吉斯回歸.....	4
2.3 k 鄰近演算法.....	5
2.4 支援向量機.....	5
2.5 決策樹.....	6
第三章 研究方法.....	7
3.1 資料介紹.....	7
3.2 判斷圖片黑色素.....	8
3.3 圖片黑色素裁切.....	9
3.4 卷積神經網路.....	12
3.5 四種機器學習分類算法.....	13
3.5.1 羅吉斯回歸.....	13
3.5.2 k 鄰近演算法.....	14
3.5.3 支援向量機.....	14
3.5.4 決策樹.....	14
3.6 數據分析.....	15
第四章 系統介面.....	17
4.1 使用案例圖.....	17
4.2 網頁預測流程.....	18
4.3 網頁介面.....	19
4.3.1 簡介.....	19
4.3.2 皮膚癌簡介.....	19
4.3.3 皮膚癌檢測.....	20
第五章 結論.....	21

5.1 未來展望.....	21
5.1.1 準確率.....	21
5.1.2 APP 化	21
參考文獻.....	22

圖目錄

圖 1-1 A (上下或左右是否不對稱)[3]	1
圖 1-2 B (邊緣是否平滑)[3]	1
圖 1-3 C (顏色是否不均勻)[3]	2
圖 1-4 D (直徑大小是否超過 0.6 公分)[3]	2
圖 1-5 E (腫瘤是否快速隆起)[3]	2
圖 2-1 CNN 概念圖[4]	3
圖 2-2 CNN 概念圖 2[4]	3
圖 2-3 羅吉斯回歸[5]	4
圖 2-4 Sigmoid 函數[5]	4
圖 2-5 k 鄰近演算法[6]	5
圖 2-6 SVM 的運作原理[7]	5
圖 2-7 決策樹[8]	6
圖 3-1 metadata 內資料	7
圖 3-2 10000 張圖片的 RGB 的最大與最小標準差	8
圖 3-3 原始皮膚圖片	9
圖 3-4 將圖 3-3 轉為灰階並使用 Sobel 計算梯度後	9
圖 3-5 將圖 3-4 模糊後	10
圖 3-6 將圖 3-5 填滿空隙後	10
圖 3-7 將圖 3-6 侵蝕並擴增	11
圖 3-8 標記原圖	11
圖 3-9 裁切後	12
圖 3-10 CNN 的 model	12
圖 3-11 準確率	13
圖 4-1 使用案例圖	17
圖 4-2 網頁使用流程圖	18
圖 4-3 網頁介面-簡介	19
圖 4-4 網頁介面-皮膚癌簡介	19
圖 4-5 網頁介面-皮膚癌檢測	20
圖 4-6 網頁介面-皮膚癌檢測結果	20

表目錄

表 3-1 羅吉斯預測結果	13
表 3-2 k 鄰近演算法預測結果	14
表 3-3 支援向量機預測結果	14
表 3-4 決策樹預測結果	14
表 3-5 metadata 中以 50 歲為分界各患者統計	15
表 3-6 患病人數百分比	16

第一章 緒論

1.1 研究背景

根據衛生福利部於 2020/06/02 發布的資料，2018 年十大癌症發生人數(男女合計)依序為(1)大腸癌(2)肺癌(3)女性乳癌(4)肝癌(5)口腔癌(含口咽、下咽)(6)攝護腺癌(7)甲狀腺癌(8)皮膚癌(9)胃癌(10)食道癌[2]，其中皮膚癌為第八位。

而對於皮膚癌也產生了一套簡易判別方式為 ABCDE 原則如下。

A (asymmetry): 上下或左右是否不對稱(圖 1-1)[3]

B (border): 邊緣是否平滑(圖 1-2) [3]

C (color): 顏色是否不均勻(圖 1-3) [3]

D (diameter): 直徑大小是否超過 0.6 公分(圖 1-4) [3]

E (elevation): 腫瘤是否快速隆起(圖 1-5) [3]



圖 1-1 A(上下或左右是否不對稱)[3]



圖 1-2 B(邊緣是否平滑)[3]



圖 1-3 C (顏色是否不均勻)[3]

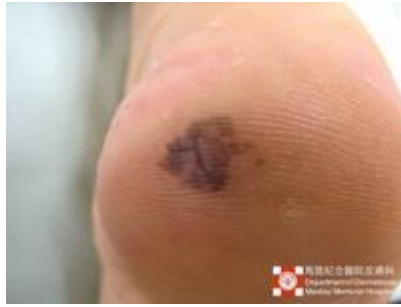


圖 1-4 D (直徑大小是否超過 0.6 公分)[3]



圖 1-5 E (腫瘤是否快速隆起)[3]

1.2 研究動機

基於以上，我們可以很清楚的看出，皮膚癌是可以透過肉眼來對其做出判斷。但即使如此，民眾對自己的判斷勢必不會太過相信，那如果此時有一個準確率更高，可以使人們信服的網站，想必會讓人安心許多。

現如今，台灣進入高齡社會，為了減少醫療資源的浪費，我們是否可以透過機器學習來分擔醫師的工作，亦或者是提高醫師的工作效率？所以對於那些皮膚出現異常的症狀卻不在乎懶的去醫院檢查亦或是過度緊張，我們希望透過這樣的簡單分析軟體讓使用者能初步了解自己的皮膚症狀從而進行進一步的處理。

1.3 研究目的

本專題將實做一個網站，將使用者上傳的圖片透過 CNN 來進行初步的判斷。若是使用者輸入更詳細的資料，我們也可以進行更近一步的運算，使得準確率能夠提升。

我們希望製做出來的網站能夠達成以下幾點。

1. 準確：不論如何，準確率一定要高，否則若是判斷錯誤，將造成無法挽回的後果。
2. 方便：不論是電腦或手機，只要連的上網路都能使用。
3. 快速：能馬上將結果回傳給使用者。

第二章 文獻回顧

2.1 CNN 卷積神經網絡

深度學習的一種方式，CNN 在影像識別方面的威力非常強大。而 CNN 也是少數參考人的大腦視覺組織來建立的深度學習模型。CNN 概念圖如圖 2-1、2-2[4]。

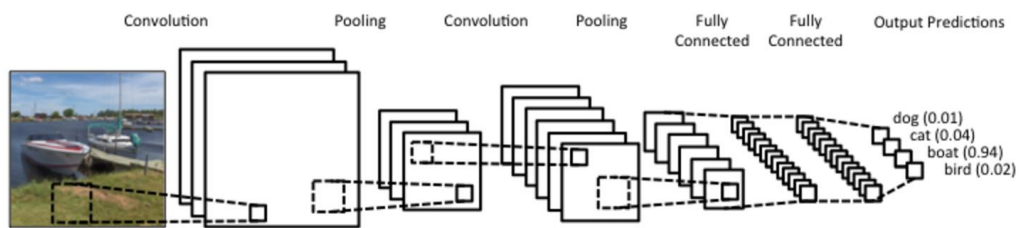


圖 2-1 CNN 概念圖[4]

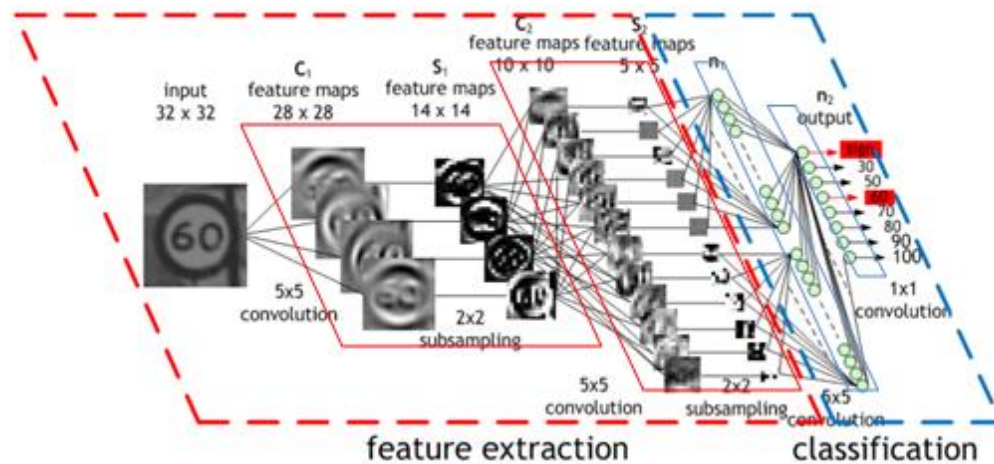


圖 2-2 CNN 概念圖 2[4]

主要分為 Convolution Layer 卷積層以及 Pooling Layer 池化層。卷積層就是將原始圖片的與特定的 Feature Detector(filter)做卷積運算，池化層主要是採用 Max Pooling，挑出矩陣當中的最大值。

2.2 羅吉斯回歸

羅吉斯回歸概念如圖 2-3 [5]，是尋找一條線，使得兩邊的數據能夠區隔開來。

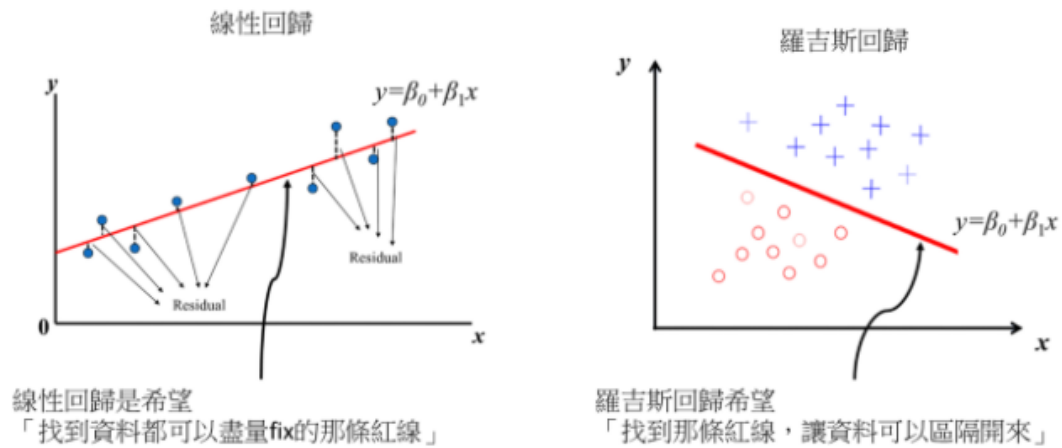


圖 2-3 羅吉斯回歸[5]

羅吉斯回歸用到的對數函數是 Sigmoid 函數(圖 2-4)[5]。

$$\sigma(f(x)) = \frac{1}{1 + e^{-f(x)}} = \frac{e^{f(x)}}{1 + e^{f(x)}}$$

圖 2-4 Sigmoid 函數[5]

線性回歸輸出是一個連續的實數，但羅吉斯回歸就是用線性回歸的輸出來判斷這個資料屬不屬於 target(二分類問題: target 和 non-target)。

2.3 k 鄰近演算法

k 鄰近演算法就是看離你最近的 K 個點，然後看哪個類別的點最多就把自己也當成那個類別，如圖 2-5[6]。

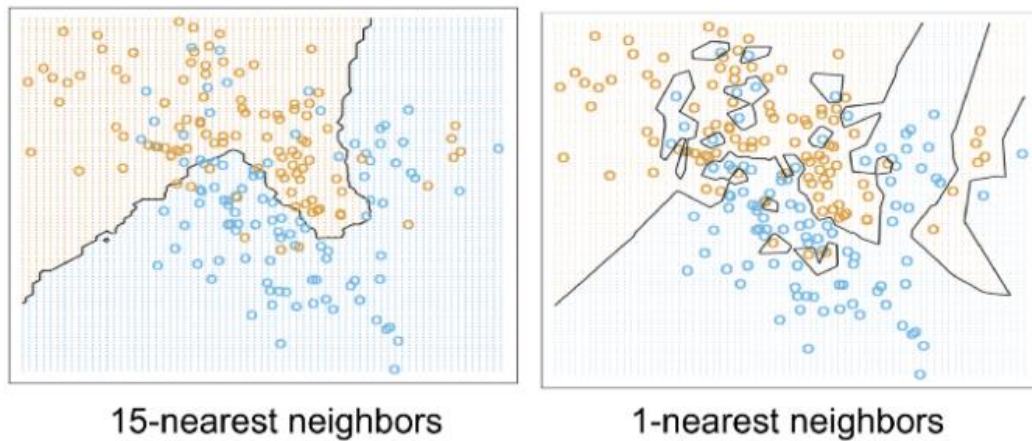


圖 2- 5 k 鄰近演算法[6]

2.4 支援向量機

支援向量機(Support Vector Machine, SVM)是一種基於統計學習理論基礎的機器學習模型。

SVM 就是將在低微度空間線性不可分的樣本映射到高維度空間去，找到一個超平面將這些樣本做有效的切割，而且，這個超平面兩邊的樣本要盡可能地遠離這個超平面。如圖 2-6[7]。

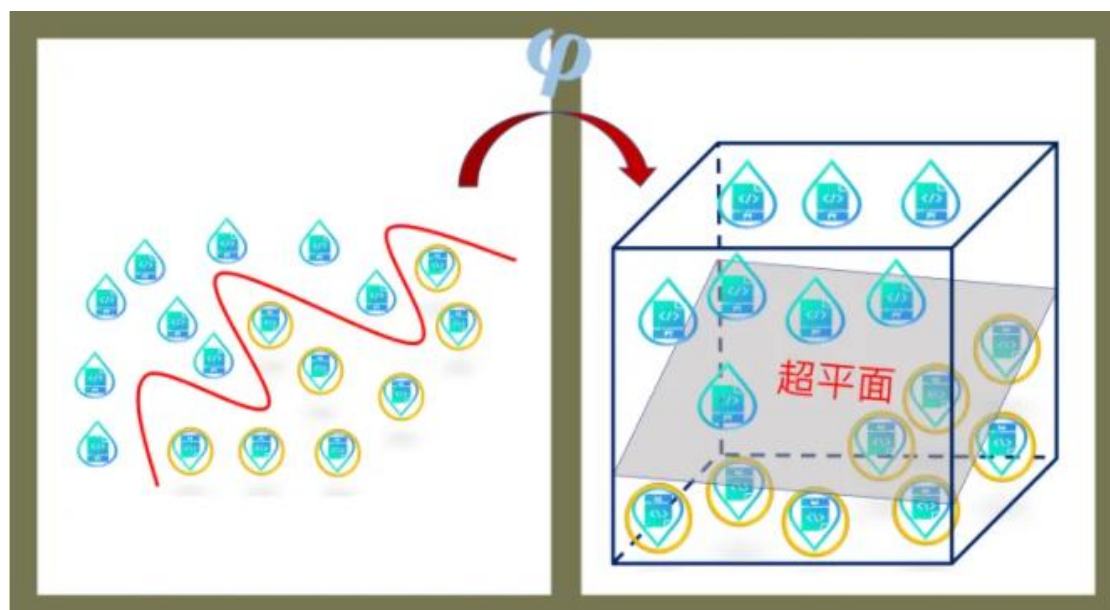


圖 2- 6 SVM 的運作原理[7]

2.5 決策樹

樣本通過「內部節點」的各個特徵的判斷，最終會到達屬於它的「葉子節點」，最後此樣本就被分在此類別，如圖 2-7[8]。

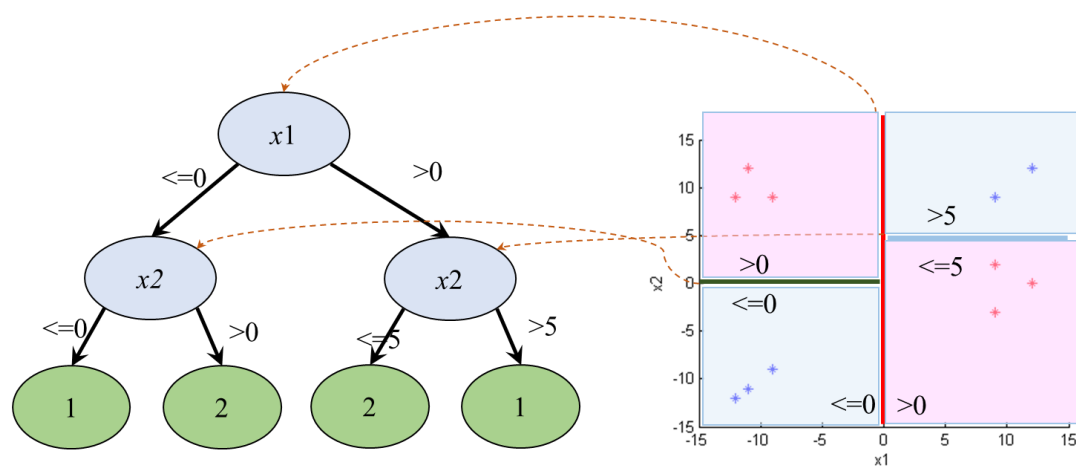


圖 2-7 決策樹[8]

第三章 研究方法

3.1 資料介紹

資料取自 kaggle[9]上的 Skin Cancer MNIST：HAM10000，裡面包含了 10015 張圖片，以及五個 csv 檔，分別為 metadata，8*8 的灰階像素 mnist，8*8 的彩色像素 mnist，28*28 的灰階像素 mnist，28*28 的彩色像素 mnist。

Metadata 裡面資料有每一張圖片的病名，患者的年齡，色班的部位，如圖 3-1。

lesion_id	image_id	dx	dx_type	age	sex	localization
HAM_000	ISIC_0027	bkl	histo	80	male	scalp
HAM_000	ISIC_0025	bkl	histo	80	male	scalp
HAM_000	ISIC_0026	bkl	histo	80	male	scalp
HAM_000	ISIC_0025	bkl	histo	80	male	scalp
HAM_000	ISIC_0031	bkl	histo	75	male	ear
HAM_000	ISIC_0027	bkl	histo	75	male	ear
HAM_000	ISIC_0029	bkl	histo	60	male	face
HAM_000	ISIC_0029	bkl	histo	60	male	face
HAM_000	ISIC_0025	bkl	histo	70	female	back
HAM_000	ISIC_0025	bkl	histo	70	female	back
HAM_000	ISIC_0025	bkl	histo	55	female	trunk
HAM_000	ISIC_0029	bkl	histo	85	female	chest
HAM_000	ISIC_0025	bkl	histo	85	female	chest
HAM_000	ISIC_0025	bkl	histo	70	male	trunk
HAM_000	ISIC_0032	bkl	histo	70	male	trunk
HAM_000	ISIC_0031	bkl	histo	65	male	back
HAM_000	ISIC_0025	bkl	histo	75	male	upper extremity
HAM_000	ISIC_0031	bkl	histo	75	male	upper extremity
HAM_000	ISIC_0029	bkl	histo	70	male	chest
HAM_000	ISIC_0032	bkl	histo	70	male	chest

圖 3- 1 metadata 內資料

3.2 判斷圖片黑色素

對於一張的圖片，首先先判斷其是否帶有黑色素，而圖片是否足夠清晰。我們利用的是標準差。對一張圖我們將R G B三個值分別進行標準差計算，透過10000張照片，我們得到了六個值，分別為R、G、B的最大與最小值。也就是說，若一張圖片的標準差，小於這些值，我們可以判斷出他是一張接近完全平穩，不帶有任何黑色素的圖。而若大於這些值，很可能這張圖片不是單純的皮膚照片最後計算結果如圖 3-2。

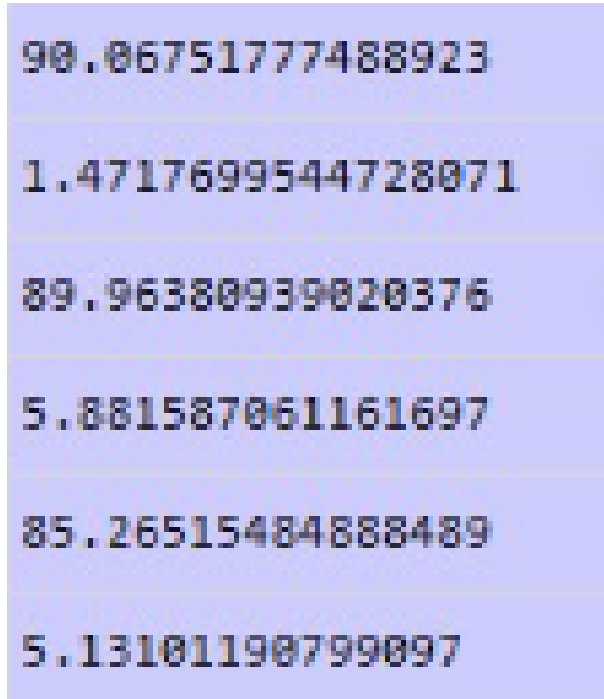


圖 3-2 10000 張圖片的 RGB 的最大與最小標準差

3.3 圖片黑色素裁切

機器學習若是有清楚的圖片將會有更高的準確率，對此我們決定將圖片不必要的部分裁切掉，僅留下黑色素的部分，對此我們先將圖 3-3 轉為灰階，並利用 opencv 的 Sobel 計算梯度進行影像處理如圖 3-4，將圖片留下高水平梯度和低垂直梯度的部分。

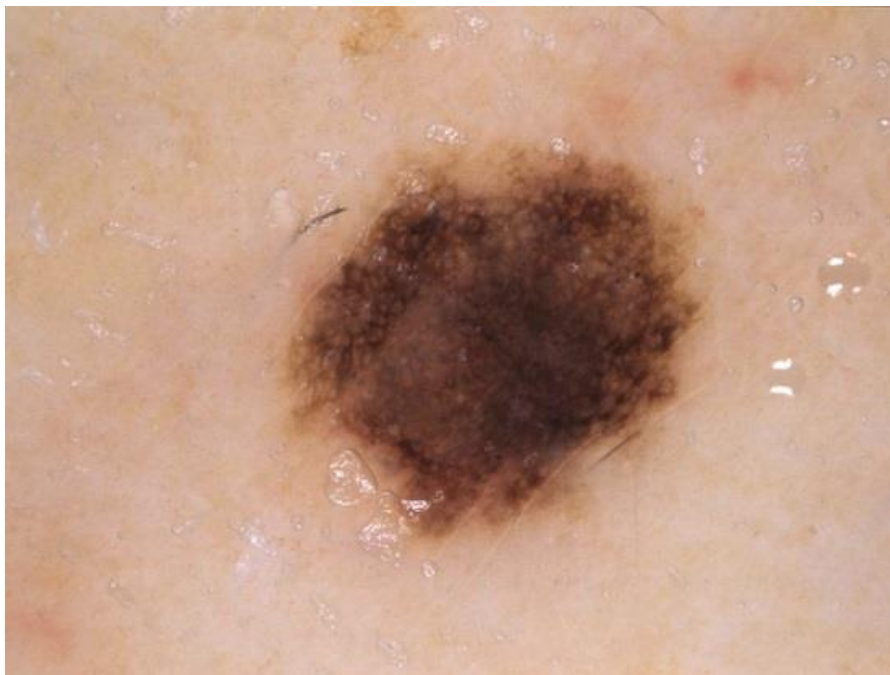


圖 3-3 原始皮膚圖片

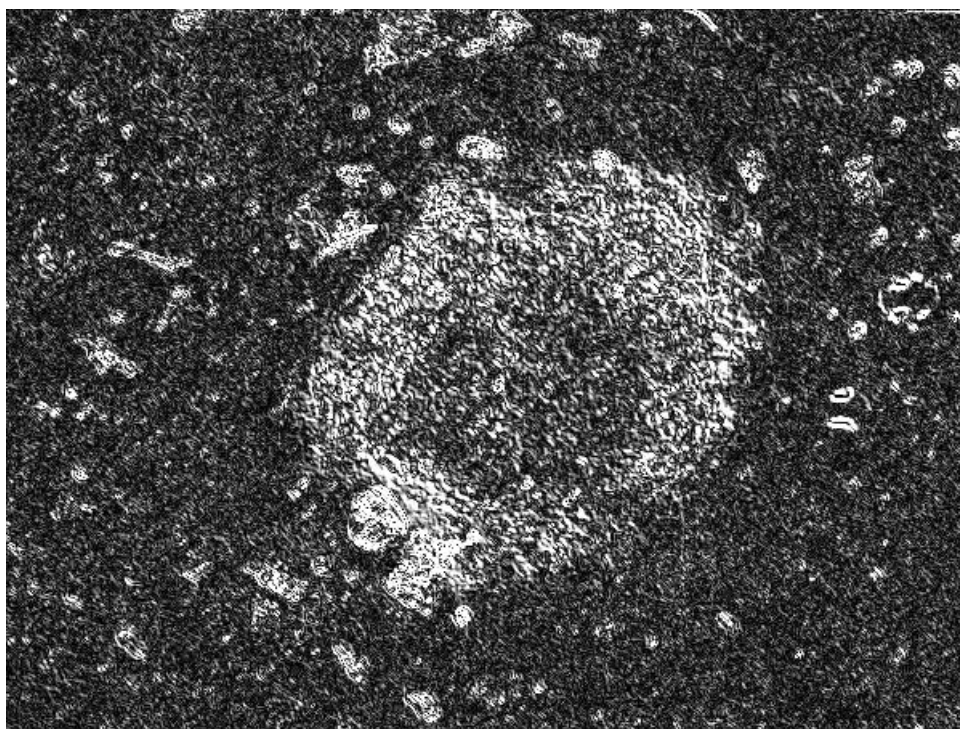


圖 3-4 將圖 3-3 轉為灰階並使用 Sobel 計算梯度後

接著將圖片模糊處理，降低周遭小區塊影響(圖 3-5)

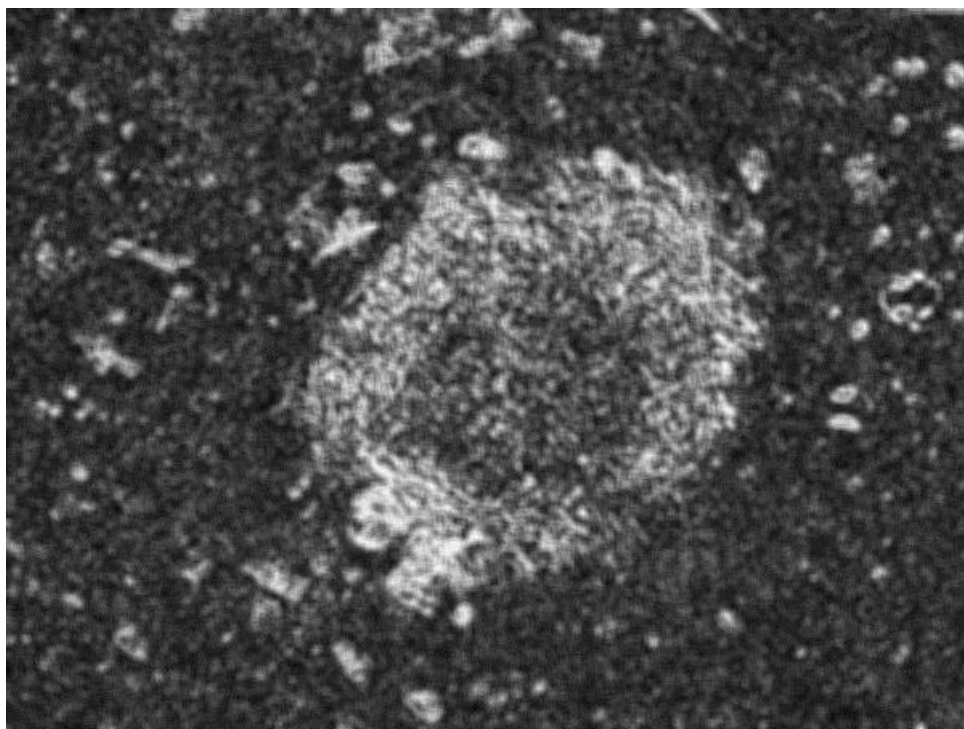


圖 3-5 將圖 3-4 模糊後

由於是取梯度的關係，圖像中間有部分空餘，將之填滿(圖 3-6)。



圖 3-6 將圖 3-5 填滿空隙後

外邊依舊有一些區塊為白色，將之侵蝕再擴增處理掉(圖 3-7)。

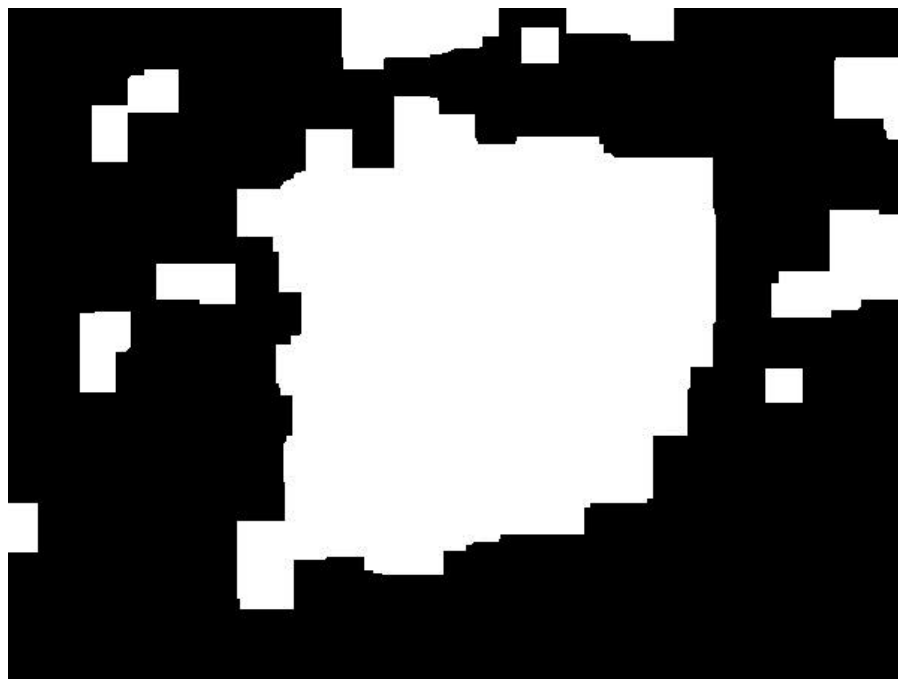


圖 3-7 將圖 3-6 侵蝕並擴增

最後找出輪廓取得四個角的座標繪製至原圖片(圖 3-8)。

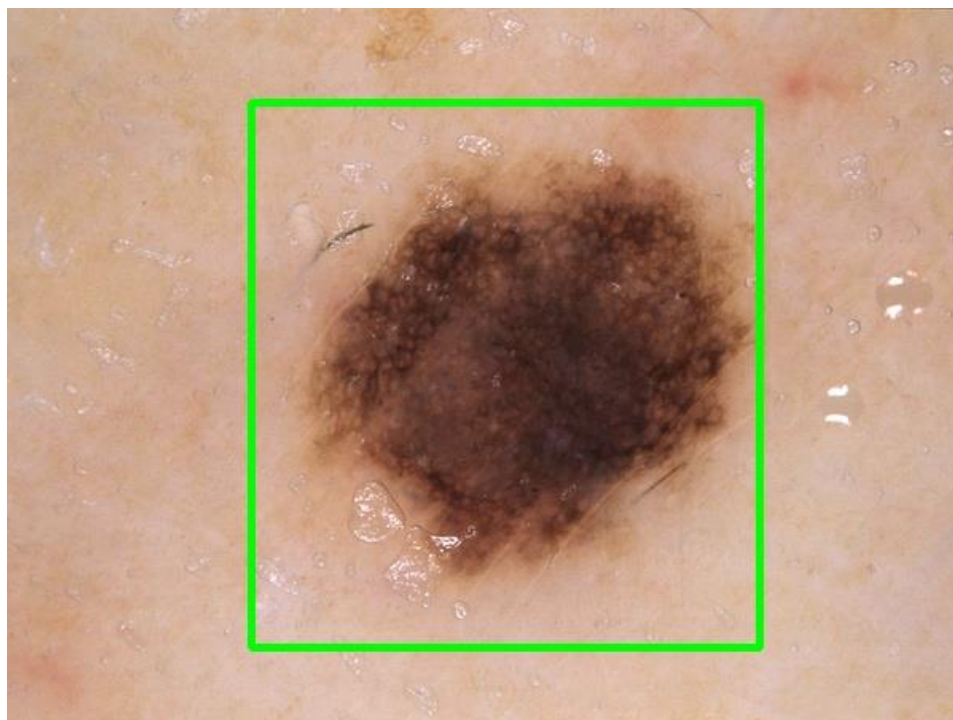


圖 3-8 標記原圖

最後裁剪圖片如圖 3-9

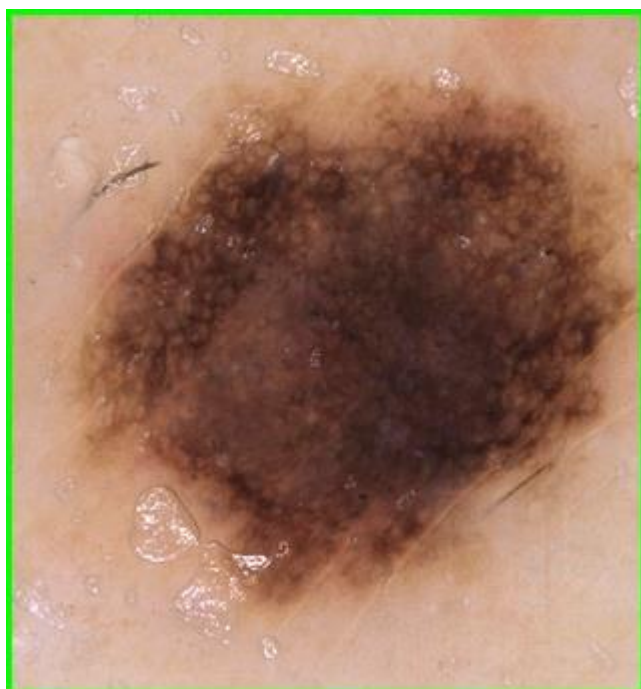


圖 3-9 裁切後

3.4 卷積神經網路

設定卷積層以及池化層(圖 3-10) [10]。

Model: "sequential_1"

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 28, 28, 32)	320
max_pooling2d_1 (MaxPooling2)	(None, 14, 14, 32)	0
conv2d_2 (Conv2D)	(None, 14, 14, 64)	18496
max_pooling2d_2 (MaxPooling2)	(None, 7, 7, 64)	0
flatten_1 (Flatten)	(None, 3136)	0
dense_1 (Dense)	(None, 1024)	3212288
dense_2 (Dense)	(None, 7)	7175
Total params: 3,238,279		
Trainable params: 3,238,279		
Non-trainable params: 0		
None		

圖 3-10 CNN 的 model

我將四個 csv 檔都嘗試過一次之後發現 28*28 的 RGB 是準確率最高的,所以就選用他來進行,將資料打散並將五個部分各自做為測試資料進行交叉驗證 (K-fold CV) 確認模型的好壞,得到五個準確率如圖 3-11。

```

2015/2015 [=====] - 1s 331us/step
accuracy= 0.7379652857780457 num= 8000 i= 0
2015/2015 [=====] - 1s 319us/step
accuracy= 0.762779176235199 num= 8000 i= 1
2015/2015 [=====] - 1s 271us/step
accuracy= 0.7598015069961548 num= 8000 i= 2
2015/2015 [=====] - 1s 282us/step
accuracy= 0.7583126425743103 num= 8000 i= 3
2015/2015 [=====] - 1s 273us/step
accuracy= 0.7617865800857544 num= 8000 i= 4

```

圖 3-11 準確率

準確率最低為 73.7% 而最高能到 76.2%。

3.5 四種機器學習分類算法

將預測出來資料與年齡放進四種模型中進行第二次預測,以提升準確率。
 以下在原先 accuracy= 0.755831241607666 訓練資料 8000 筆,測試資料 2015 筆的情況下。

3.5.1 羅吉斯回歸

表 3-1 羅吉斯預測結果

	CNN 預測正確	CNN 預測錯誤
羅吉斯預測正確	1512	41
羅吉斯預測錯誤	11	451

由表 3-1 可以看出來,羅吉斯回歸即使有將原先預測錯誤的改為正確,但也有將正確的改為錯誤,整體來說在 2015 筆測試資料中,僅僅將正確的資料增加了 30 筆,使正確率增加了 1.4% 左右。

但如果就修正率來說,在 1523 筆正確資料裡,他僅將 11 筆修正為錯誤的,機率約為 0.72%。但在 492 筆錯誤資料裡,他將 41 筆修正為正確的,機率約為 8.33%。可以看出改正的機率是明顯大於改錯的機率的。

3.5.2 k 鄰近演算法

表 3-2 k 鄰近演算法預測結果

	CNN 預測正確	CNN 預測錯誤
KNN 預測正確	1479	57
KNN 預測錯誤	44	435

由表 3-2 可以看到，經過 k 鄰近演算法修正的資料，預測正確的筆數僅僅增加了 13 筆，使得正確率增加了約 0.6%。

其修正率，在 1523 筆正確資料裡，他將 44 筆修正為錯誤的，機率約為 2.88%。但在 492 筆錯誤資料裡，他將 57 筆修正為正確的，機率約為 11.58%。

即使將錯誤的改正許多，但也因為將正確的改為錯誤的過多，使得整體的準確率並沒有提升太多。

3.5.3 支援向量機

表 3-3 支援向量機預測結果

	CNN 預測正確	CNN 預測錯誤
clf 預測正確	1465	39
clf 預測錯誤	58	453

由表 3-3 可以看出在支援向量機重新預測過後，整體的準確率是下降的。

3.5.4 決策樹

表 3-4 決策樹預測結果

	CNN 預測正確	CNN 預測錯誤
clftree 預測正確	1391	91
clftree 預測錯誤	132	401

由表 3-4 可以看出在支援向量機重新預測過後，整體的準確率是下降的。即使 clftree 預測正確的量在四個方法裡面是最高的，但同時預測錯誤的卻也最多。

3.6 數據分析

表 3-5 metadata 中以 50 歲為分界各患者統計

	50 歲以下(含)	50 歲以上
脂漏性角化症	228	861
黑色素瘤	312	799
鱗狀細胞癌	39	288
基底細胞癌	89	425
皮膚纖維瘤	59	56
黑素細胞痣	4554	2106
血管病變	66	76
總計	5347	4611

由表 3-5 可見,以 50 歲為分界,皮膚病的分類是有明顯差異的。

透過表 3-6 應該能更明確的看出對於病症,年齡是有一定的影響的。

表 3-6 患病人數百分比

	50 歲以下(含) (人數/總計)	50 歲以上比值 (人數/總計)
脂漏性角化症	4.26%	18.67%
黑色素瘤	5.83%	17.32%
鱗狀細胞癌	0.72%	6.24%
基底細胞癌	1.66%	9.2%
皮膚纖維瘤	1.10%	1.21%
黑素細胞痣	85.16%	45.67%
血管病變	1.23%	1.64%
總計	100%	100%

我們可以看出，在 50 歲以下的人群裡，絕大多數都是良性的黑素細胞痣。但在 50 歲以上的人群裡，雖說黑素細胞痣依舊占了接近百分之 50，但很明顯的下降了許多。而其他六種皮膚病，除了皮膚纖維瘤以及血管病變，都是有極大幅度的提升。

第四章 系統介面

4.1 使用案例圖

如圖 4-1，使用者能查看關於各類皮膚癌資料。

使用者能上傳照片，並根據是否有輸入年齡給予不同預測結果。

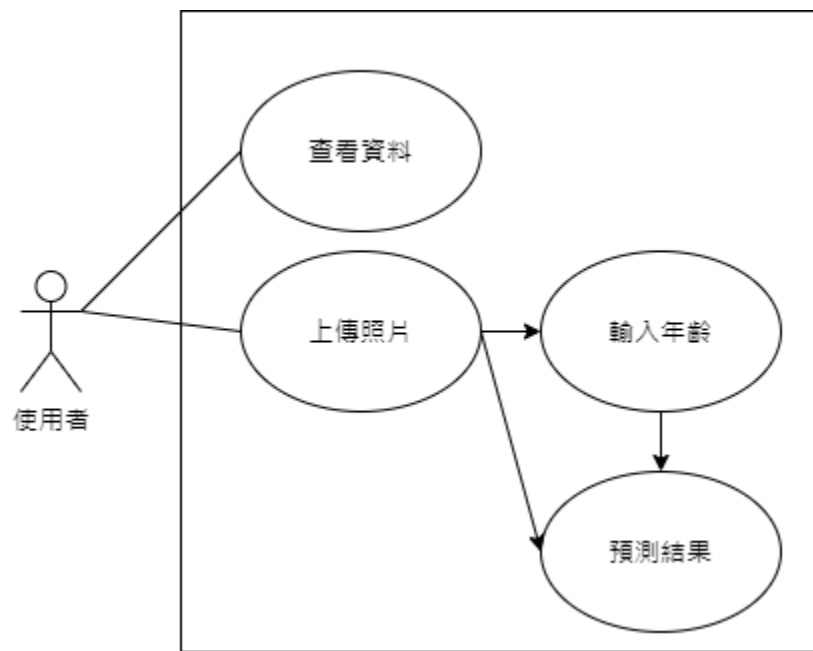


圖 4-1 使用案例圖

4.2 網頁預測流程

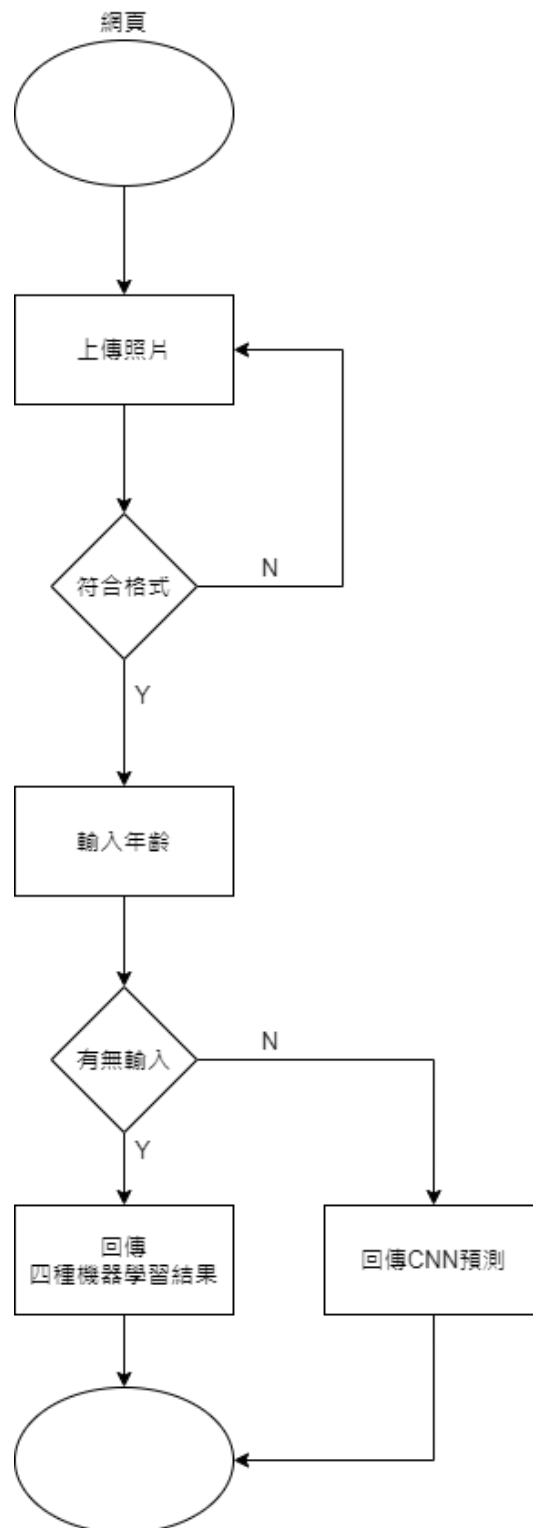


圖 4-2 網頁使用流程圖

圖 4-2 為網頁使用流程圖，網頁會先判斷使用者上傳的圖片是否符合格式。並依據年齡的輸入與否，回傳不同的預測結果。

4.3 網頁介面

4.3.1 簡介



圖 4-3 網頁介面-簡介

如圖 4-3 在這裡將簡單介紹此網頁的功能。

4.3.2 皮膚癌簡介



圖 4-4 網頁介面-皮膚癌簡介

如圖 4-4，讓使用者對於皮膚癌能有初步的認識。

4.3.3 皮膚癌檢測



圖 4-5 網頁介面-皮膚癌檢測

如圖 4-5，使用者可以上傳圖片以及輸入年齡，在一切資料準備完畢後將會開始預測最後結果如圖 4-6。



方法	病名
DecisionTreeClassifier	黑素細胞癌
LinearSVC	黑素細胞癌
KNeighbors	黑素細胞癌
LogisticRegression	黑素細胞癌

圖 4-6 網頁介面-皮膚癌檢測結果

第五章 結論

目前就圖像辨識來說，CNN依舊是最好的選擇。而在機器學習分類算法，那四種除了羅吉斯回歸，準確率都沒有太大的改變，甚至還有下降，如果有更完善的資料，不只是年齡，應當可以使之發揮更大的效用。

5.1 未來展望

5.1.1 準確率

不論何種機器學習，準確率都是必須優先提升的一點。就我們這個專題來說，若是因為誤判而造成延誤就醫，將是非常危險的。

5.1.2 APP 化

此專題僅建立在 web 平台上，若有可能，希望未來可以將之移植至 ios 以及 android 上，使得使用上更加的方便。

參考文獻

- [1] 長庚醫院。面對皮膚癌。2020年12月1日，
<https://www1.cgmh.org.tw/khcc/datas/%AD%B1%B9%EF%A5%D6%BD%A7%C0%F9.pdf>
- [2] 衛生福利部。衛生福利部公布癌症發生資料。2020年12月1日，
<https://www.hpa.gov.tw/Pages/Detail.aspx?nodeid=4141&pid=12682>
- [3] 馬偕紀念醫院 吳育弘醫師。認識皮膚癌。2020年12月1日，
<http://www.mmh.org.tw/taitam/derma/academic/article/untitled-2.html>
- [4] Yeh James。卷積神經網絡介紹。2020年12月1日，
<https://medium.com/jameslearningnote/%E8%B3%87%E6%96%99%E5%88%86%E6%9E%90-%E6%A9%9F%E5%99%A8%E5%AD%B8%E7%BF%92-%E7%AC%A5-1%E8%AC%9B-%E5%8D%B7%E7%A9%8D%E7%A5%9E%E7%B6%93%E7%B6%B2%E7%B5%A1%E4%BB%8B%E7%B4%B9-convolutional-neural-network-4f8249d65d4f>
- [5] Tommy Huang。機器/統計學習：羅吉斯回歸(Logistic regression)。2020年12月1日，
<https://medium.com/@chih.sheng.huang821/%E6%A9%9F%E5%99%A8-%E7%B5%B1%E8%A8%88%E5%AD%B8%E7%BF%92-%E7%BE%85%E5%90%89%E6%96%AF%E5%9B%9E%E6%AD%B8-logistic-regression-aff7a830fb5d>
- [6] 林昱北。[Machine Learning] kNN 分類演算法。2020年12月1日，
<https://medium.com/@NorthBei/machine-learning-knn%E5%88%86%E9%A1%9E%E6%BC%94%E7%AE%97%E6%B3%95-b3e9b5aea8df>
- [7] PyInvest。[機器學習首部曲] 支援向量機 SVM。2020年12月1日，
<https://pyecontech.com/2020/03/24/%E6%A9%9F%E5%99%A8%E5%AD%B8%E7%BF%92%E9%A6%96%E9%83%A8%E6%9B%B2-%E6%94%AF%E6%8F%B4%E5%90%91%E9%87%8F%E6%A9%9F-svm/>
- [8] Tommy Huang。機器學習：決策樹 (Decision Tree)。2020年12月1日，
<https://chih-sheng-huang821.medium.com/%E6%A9%9F%E5%99%A8%E5%AD%B8%E7%BF%92-%E6%B1%BA%E7%AD%96%E6%A8%B9-decision-tree-ed102ee62dfa>
- [9] K Scott Mader。Kaggle-Skin Cancer MNIST: HAM10000。2020年12月1日，
<https://www.kaggle.com/kmader/skin-cancer-mnist-ham10000>
- [10] 張哲誠老師。機器學習基礎理論與實作課程。2020年12月1日，
<https://courses.openedu.tw/courses/course-v1:FCUx+QA+20002/course/>