

# 公司員工離職預測

統計碩— 107354012  
陳冠廷

# 目錄

1

## 資料介紹

- 研究目的
- 變數解釋
- 資料來源
- 資料探索

2

## 資料整理

- 資料不平衡處理
- 變數處理
- 資料集拆分

3

## 模型配置

- 模型處理步驟
- 使用Lasso
- 使用逐步迴歸
- 確定最終模型

4

## 模型評估

- 準確度
- ROC Curve
- 混淆矩陣
- 模型解釋



# PART 1

## 資料介紹

- 研究目的
- 資料來源
- 變數解釋
- 資料探索

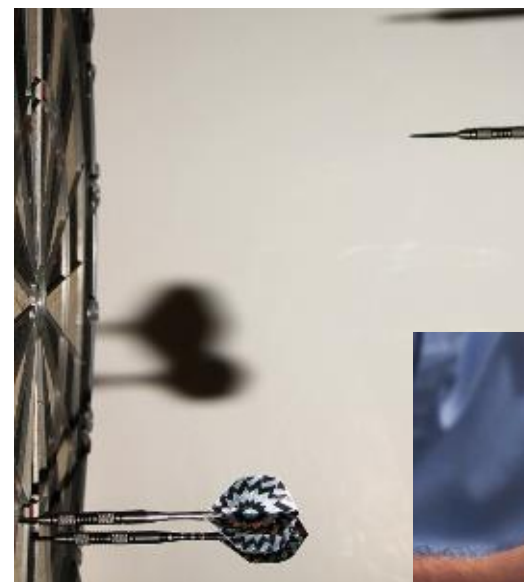


# 研究目的

## 資料介紹

員工一直是影響公司營運的重要因素之一，因此如何解決優秀員工流失的問題，是世界許多公司正面臨的挑戰。

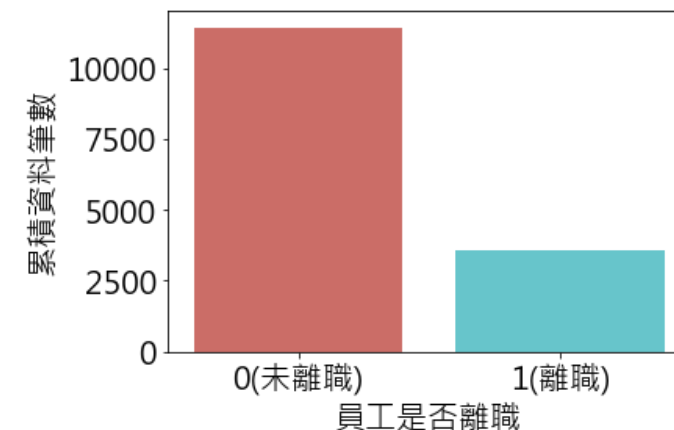
這次我透過對員工離職數據集的分析，嘗試用 Logistic Regression 來預測有很大機率將要離職的員工及其背後影響員工離職的關鍵（員工滿意度、薪資水平、平均每月工時等），在了解對離職員工產生重大影響的因素後，便可以建議公司採取適當的措施來改善這些因素以留住公司人才。



# 資料來源

## 資料介紹

- 原始出處: **kaggle**
- 資料網址: <https://www.kaggle.com/ludobenistant/hr-analytics>
- 資料提供者: Ludovic Bénistant
- 資料筆數: 14, 999筆員工資料(離職: 未離職=23. 8:76. 2)，10個變數(9個解釋變數，一個反應變數)



	left	satisfaction_level	last_evaluation	number_project	average_monthly_hours	time_spend_company	Work_accident	promotion_last_5years	position	salary
0	1	0.38	0.53	2	157	3	0	0	sales	low
1	1	0.80	0.86	5	262	6	0	0	sales	medium
2	1	0.11	0.88	7	272	4	0	0	sales	medium
3	1	0.72	0.87	5	223	5	0	0	sales	low
4	1	0.37	0.52	2	159	3	0	0	sales	low
5	1	0.41	0.50	2	153	3	0	0	sales	low
6	1	0.10	0.77	6	247	4	0	0	sales	low
7	1	0.92	0.85	5	259	5	0	0	sales	low

### 反應變數



Left

是否離職  
值為0或1  
0代表未離職  
1代表已離職

### 解釋變數

satisfaction\_level



員工滿意度  
Numeric:0~1

average\_monthly\_hours



平均每月工時  
Integer:96~310

number\_project



參與過的專案數  
Integer:2~7

last\_evaluation



績效評估  
Numeric:0~1

time\_spend\_company



進入公司的年數  
Integer:2~10

position



在公司所屬部門  
類別型變數  
共10種職位

Work\_accident



是否有過工作意外  
值為0或1  
0代表未發生  
1代表曾發生

promotion\_last\_5years



五年內是否升職  
值為0或1  
0代表未升職  
1代表已升職

salary



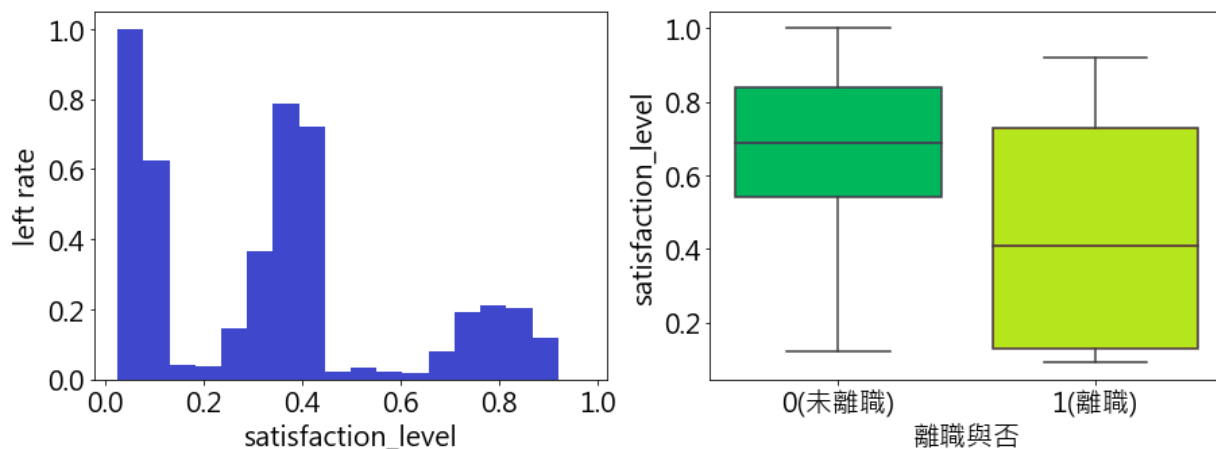
薪資水平  
類別型變數  
共分為  
High,Medium,Low  
三個水準

# 資料探索

## 資料介紹

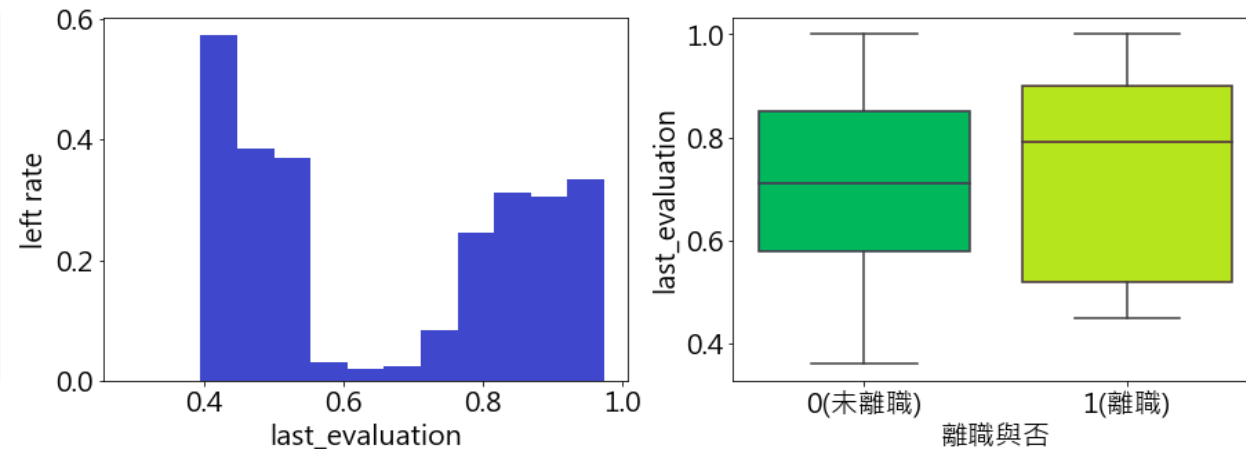
### satisfaction\_level

整體來看，員工滿意度越高，離職率越低



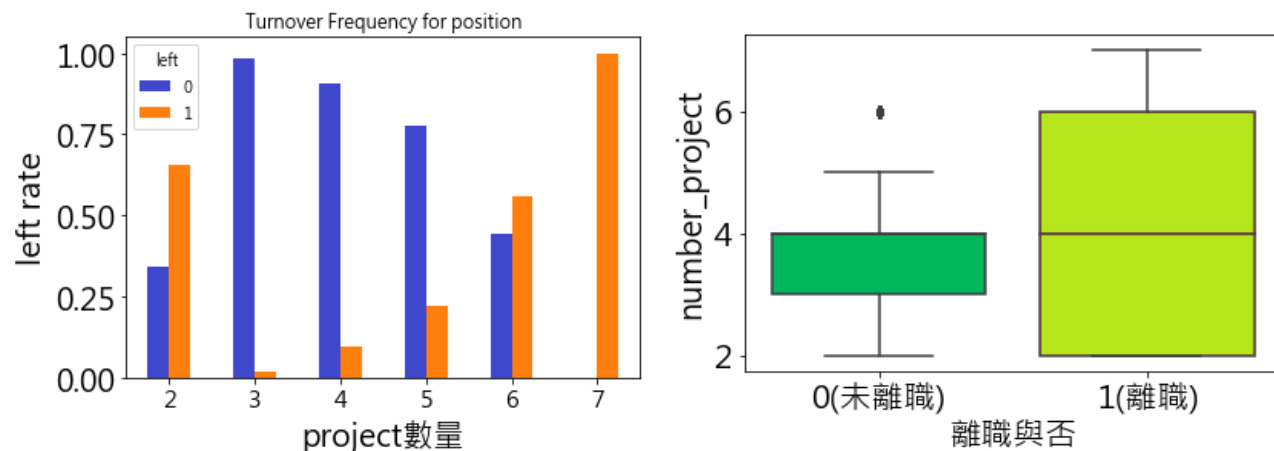
### last\_evaluation

離職率與績效評估看似沒有關聯



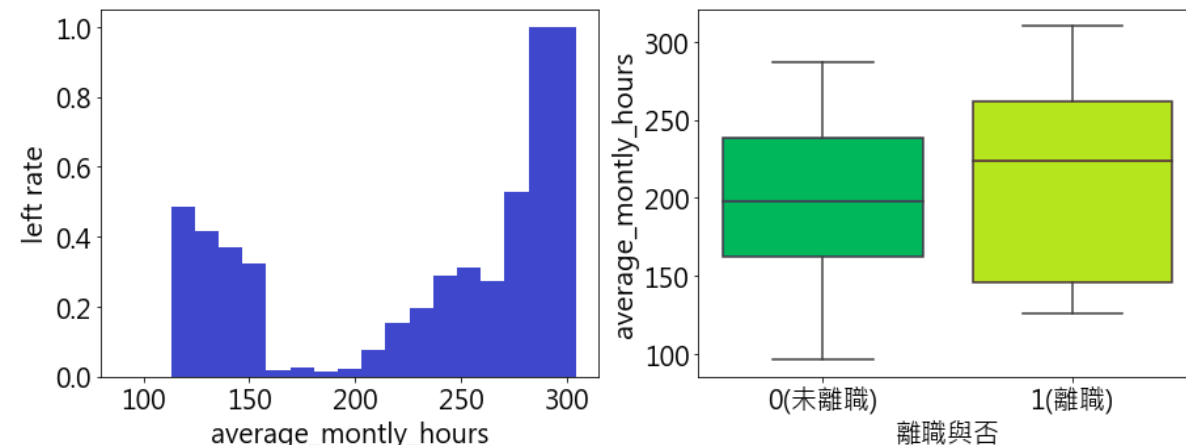
### number\_project

專案數越大，離職率越大



### average\_monthly\_hours

每月平均工時越多，離職率越大

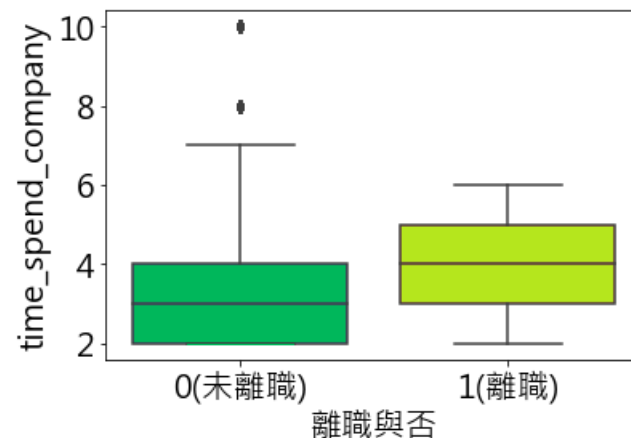
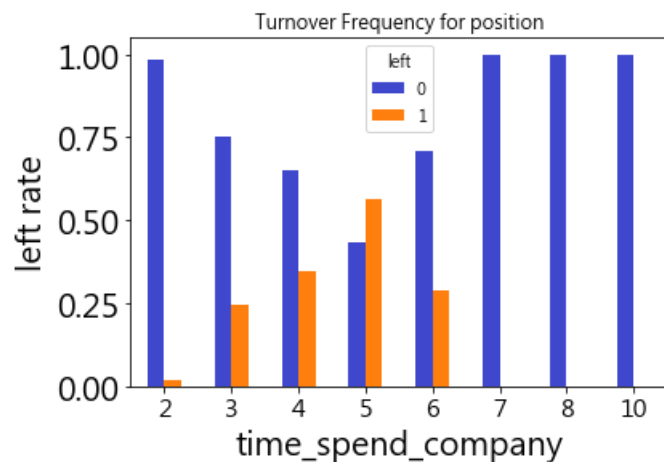


# 資料探索

# 資料介紹

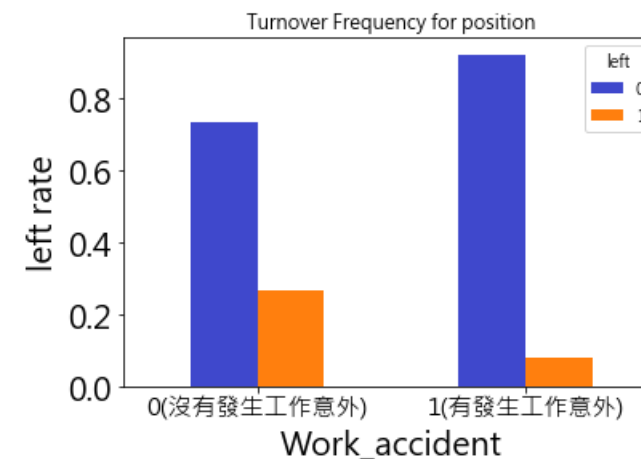
time\_spend\_company

工作五年內，離職率隨工作年數增加



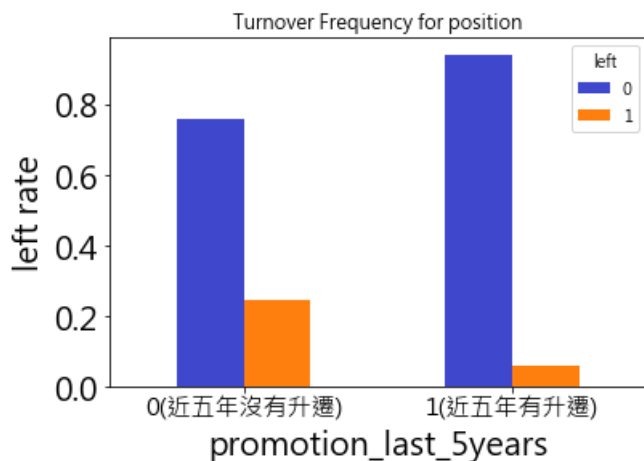
Work\_accident

有無發生意外與離職率沒有太大差異



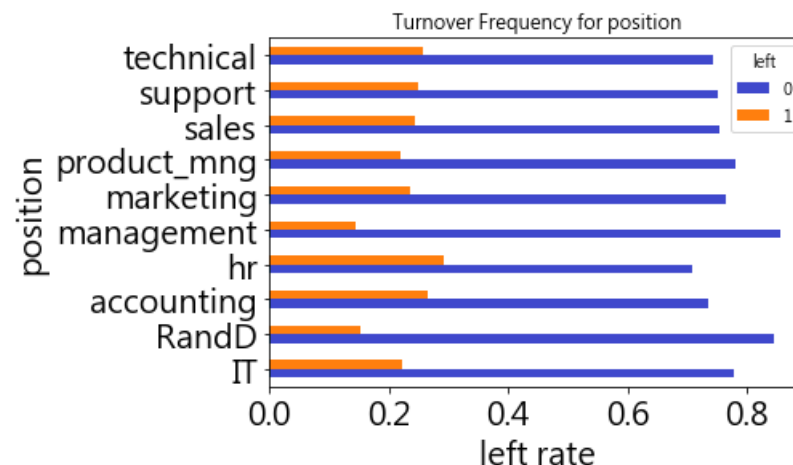
promotion\_last\_5years

有無升遷與離職率沒有太大差異



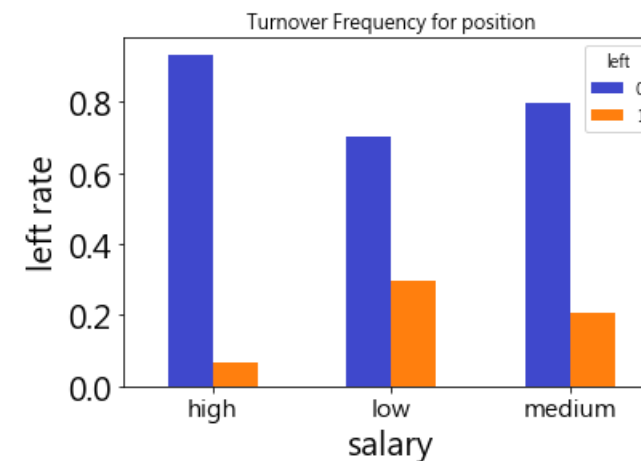
Position

所在部門與離職率沒有太大差異



Salary

薪水水平為高者較不易離職







## PART 2

## 資料整理

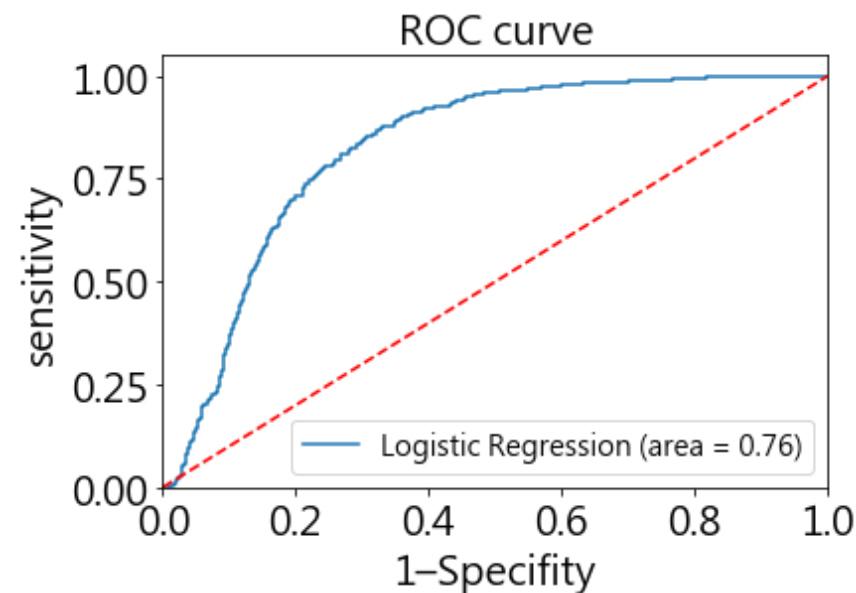
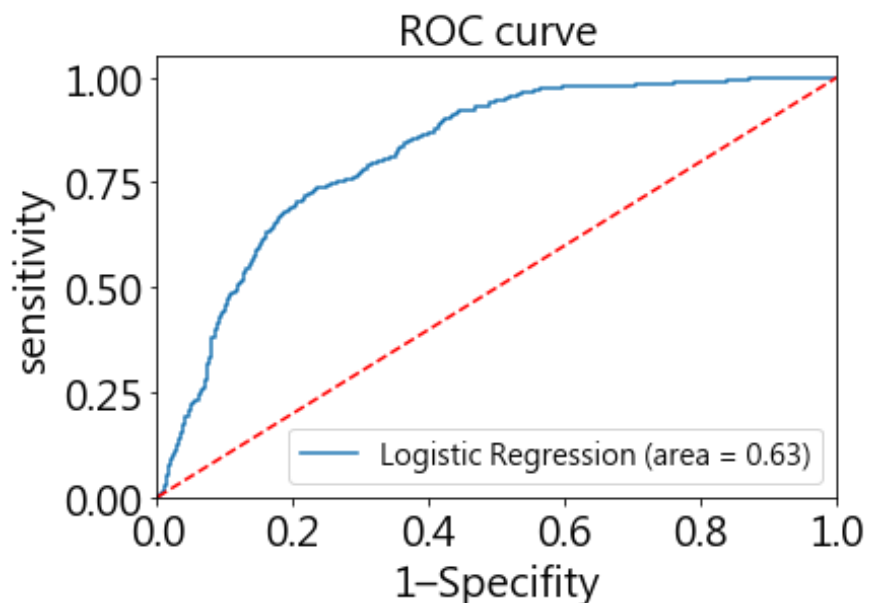
- 資料不平衡處理
- 資料集拆分
- 變數處理

# 資料不平衡處理

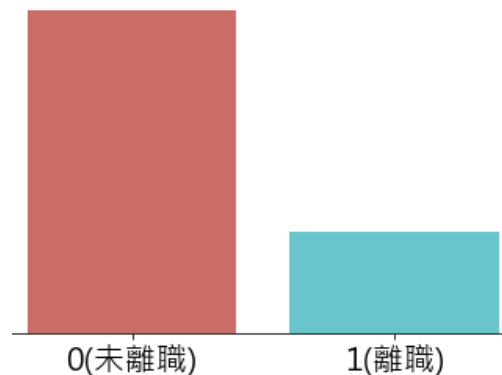
## 資料整理

此份資料為不平衡資料(離職:未離職=23.8:76.2)

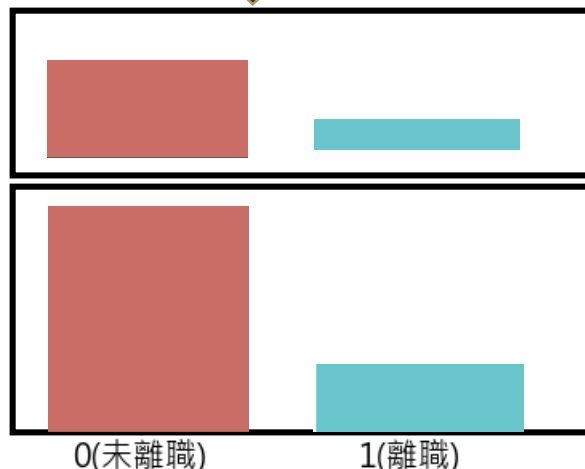
- 採取SMOTE方法(Synthetic Minority Over-Sampling Technique)模擬生成樣本，  
從離職樣本中生成樣本使得離職與未離職樣本數相同
- SMOTE運作概述: 1.對於每一個少數樣本 $s$ ，以歐式距離計算該樣本與其他少數樣本的距離  
2.以此距離找出與該樣本 $s$ 最近的 $K$ 個樣本  
3.從 $K$ 個樣本中選取 $M$ 個樣本  
4.將該樣本 $s$ 在 $M$ 個樣本中的每一個樣本 $r_i$ 透過 $S_{\text{new}} = \lambda s + (1 - \lambda) r_i$ 生成新樣本 $S_{\text{new}}$ ，  
其中 $\lambda$ 為0~1的任意數
- 採用SMOTE方法後，AUC獲得改善(0.63→0.76)



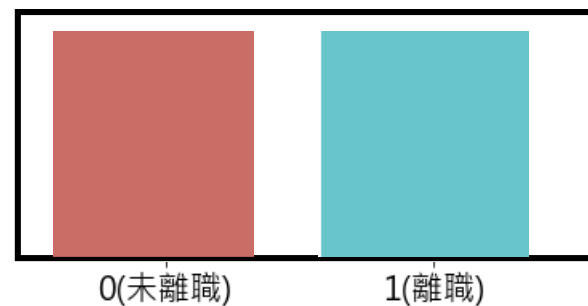
# 資料集拆分



分別抓出離職以及未離職之員工資料

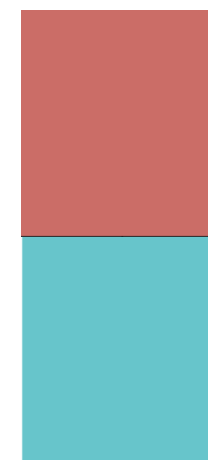
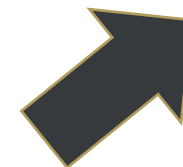
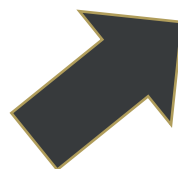


分別將兩種資料按3:7  
拆分測試集及訓練集



訓練集

將離職以及未離職之員工訓練集進行SMOTE使兩種資料樣本數相同



模型訓練集

合併離職以及未離職之員工SMOTE後之訓練集作為模型之訓練集



模型測試集

合併離職以及未離職之員工之測試集作為模型訓練之測試集  
(維持原數據集中離職與未離職比)

- 將數值型變數歸一化，使變數範圍縮小到0~1之間，避免因數值大小使模型係數差異過大
- 類別型變數轉換：
  - 1.position欄位中分為IT,RandD,accounting,hr,management,marketing,produc\_r\_mng, sales,support,technical共10種，因此轉換為9個dummy variables
  - 2.salary欄位中分為High,Medium,Low共三種，因此轉換為2個dummy variables
  - 3.轉換後變數總數為18個
- 增加交互作用項：
  - 1.將類別型變數轉換後的變數兩兩相乘做一個交互作用項，增加 $C_2^{18} = 153$ 個變數
  2. 轉換後變數總數為171個
- 增加高次項：
  1. 將原數據集中的5個連續變數satisfaction\_level, last\_evaluation, average\_monthly\_hours, number\_project,time\_spend\_company增加二次項及三次項，共增加10個變數
  - 2.轉換後變數總數為181個



# PART 3

## 模型配置

- 模型處理步驟
- 使用Lasso
- 使用逐步迴歸
- 確定最終模型



拆分資料集

將Lasso篩選後的變數放入  
Logistic Regression 模型裡  
並去掉不顯著的變數

評估最終模型



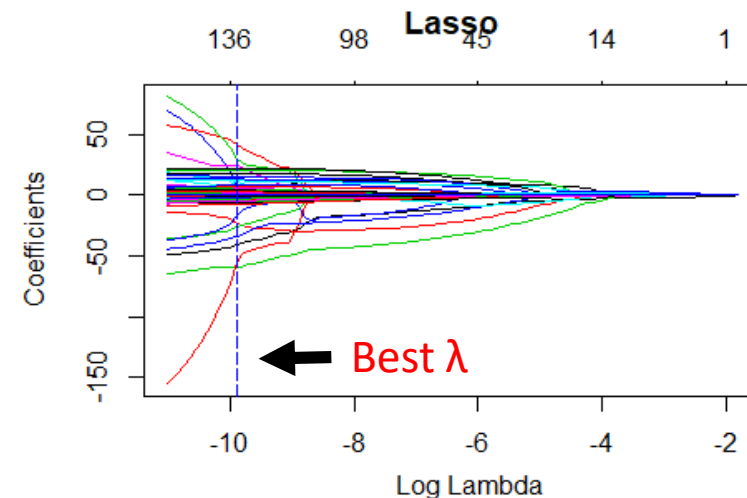
使用Lasso篩選  
變數

將剩餘變數分別帶入向前、  
向後、雙向逐步迴歸中，  
選擇AIC最低的作為最終  
模型

# 使用Lasso

## 模型配置

- 最佳 $\lambda$ 為 $5.05202 \times 10^{-5}$
- 使用Lasso後變數由181個變為130個
- 將剩餘130個變數帶入Logistic Regression 模型裡  
並刪除P-value>0.01的變數
- 刪除後剩餘32個變數



Lasso變數收斂情形(列出部分)

	Estimate
(Intercept)	27.20006492
satisfaction_level	-40.95892796
last_evaluation	-23.65009906
number_project	-59.27808348
average_monthly_hours	8.66861291
promotion_last_5years	.
Work_accident	-1.83021398

刪除P-value>0.01的變數(列出部分)

	Estimate	Pr(> z )
<del>(Intercept)</del>	<del>7.118e-01</del>	<del>8.298828e-01</del>
satisfaction_level	-55.20267092	3.022989e-53
<del>last_evaluation</del>	<del>12.19882159</del>	<del>3.928042e-02</del>
number_project	-74.95249596	2.257499e-28
average_monthly_hours	148.20130101	8.639483e-25
time_spend_company	20.46260773	2.272918e-05
<del>Work_accident</del>	<del>-13.53175921</del>	<del>9.466631e-01</del>

# 剩餘變數32個

## 模型配置

原始變數4個，高次項變數9個，交互作用項19個

- |   |   |
|---|---|
| [1] "satisfaction_level"                        | [17] "last_evaluation*number_project"           |
| [2] "number_project"                            | [18] "last_evaluation*average_monthly_hours"    |
| [3] "average_monthly_hours"                     | [19] "last_evaluation*time_spend_company"       |
| [4] "time_spend_company"                        | [20] "last_evaluation*position_sales"           |
| [5] "satisfaction_level^2"                      | [21] "last_evaluation*salary_low"               |
| [6] "satisfaction_level^3"                      | [22] "last_evaluation*salary_medium"            |
| [7] "last_evaluation^2"                         | [23] "number_project*average_monthly_hours"     |
| [8] "last_evaluation^3"                         | [24] "number_project*time_spend_company"        |
| [9] "number_project^2"                          | [25] "number_project*position_accounting"       |
| [10] "average_monthly_hours^2"                  | [26] "average_monthly_hours*time_spend_company" |
| [11] "average_monthly_hours^3"                  | [27] "time_spend_company*position_IT"           |
| [12] "time_spend_company^2"                     | [28] "time_spend_company*position_management"   |
| [13] "time_spend_company^3"                     | [29] "time_spend_company*position_sales"        |
| [14] "satisfaction_level*last_evaluation"       | [30] "time_spend_company*salary_medium"         |
| [15] "satisfaction_level*average_monthly_hours" | [31] "position_management*salary_low"           |
| [16] "satisfaction_level*time_spend_company"    | [32] "position_management*salary_medium"        |

- 使用向前逐步迴歸，分析後之變數剩29個，減少average\_monthly\_hours^2, average\_monthly\_hours^3, number\_project\*position\_accounting，AIC為6476.3
- 使用向後逐步迴歸，分析後之變數剩31個，減少number\_project\*position\_accounting，AIC為6252.2
- 使用雙向逐步迴歸，分析後之變數剩29個，減少average\_monthly\_hours^2, average\_monthly\_hours^3, number\_project\*position\_accounting，與向前逐步迴歸相同，AIC為6476.3
- 選擇AIC最低之向後逐步迴歸所得模型，作為最終模型

Logit( $\hat{left}$ ) =

-8.38-49.14\*satisfaction\_level-57.13 \*number\_project +159.28\*average\_monthly\_hours  
+31.86\*time\_spend\_company + 54.04\*satisfaction\_level^2 -35.72\* satisfaction\_level^3  
-58.6\*last\_evaluation^2 + 29.61\*last\_evaluation^3 +19.36\*number\_project^2  
-302.6\* average\_monthly\_hours^2 +154.98\* average\_monthly\_hours^3 -81.86\* time\_spend\_company^2  
+ 31.07\* time\_spend\_company^3 + 12.87\*satisfaction\_level\*last\_evaluation  
+ 11\*satisfaction\_level\*average\_monthly\_hours + 20.18\*satisfaction\_level\*time\_spend\_company  
+ 19.62\* last\_evaluation\*number\_project + 16.74\*last\_evaluation\*average\_monthly\_hours  
+ 13.81\*last\_evaluation\*time\_spend\_company + 1.79\*last\_evaluation\*position\_sales  
+ 2.34\*last\_evaluation\*salary\_low +2.78\*last\_evaluation\*salary\_medium  
+ 26.31\*number\_project\*average\_monthly\_hours + 10.55\*number\_project\*time\_spend\_company  
+ 7.94\*average\_monthly\_hours\*time\_spend\_company -0.53\*time\_spend\_company\*position\_IT  
-5.64\*time\_spend\_company\*position\_management -2.97\*time\_spend\_company\*position\_sales  
-2.05\*time\_spend\_company\*salary\_medium +2.47\*position\_management\*salary\_low  
+ 2.29\*position\_management\*salary\_medium





# PART 4

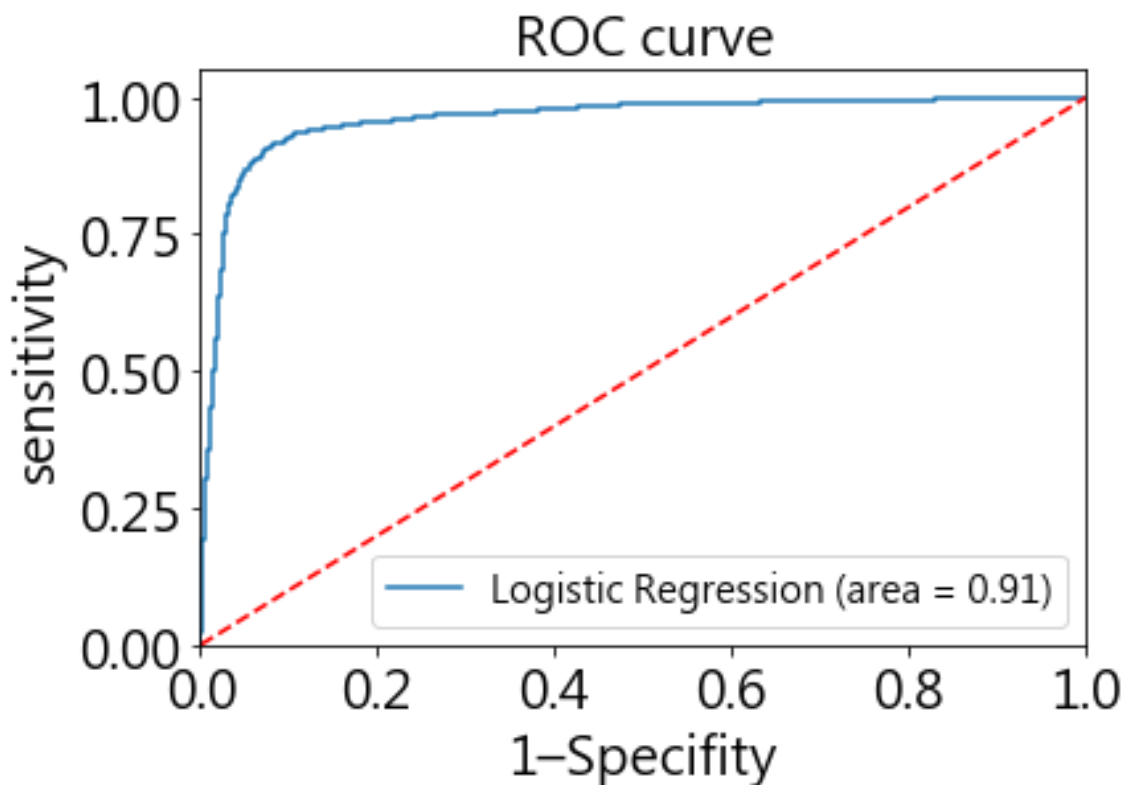
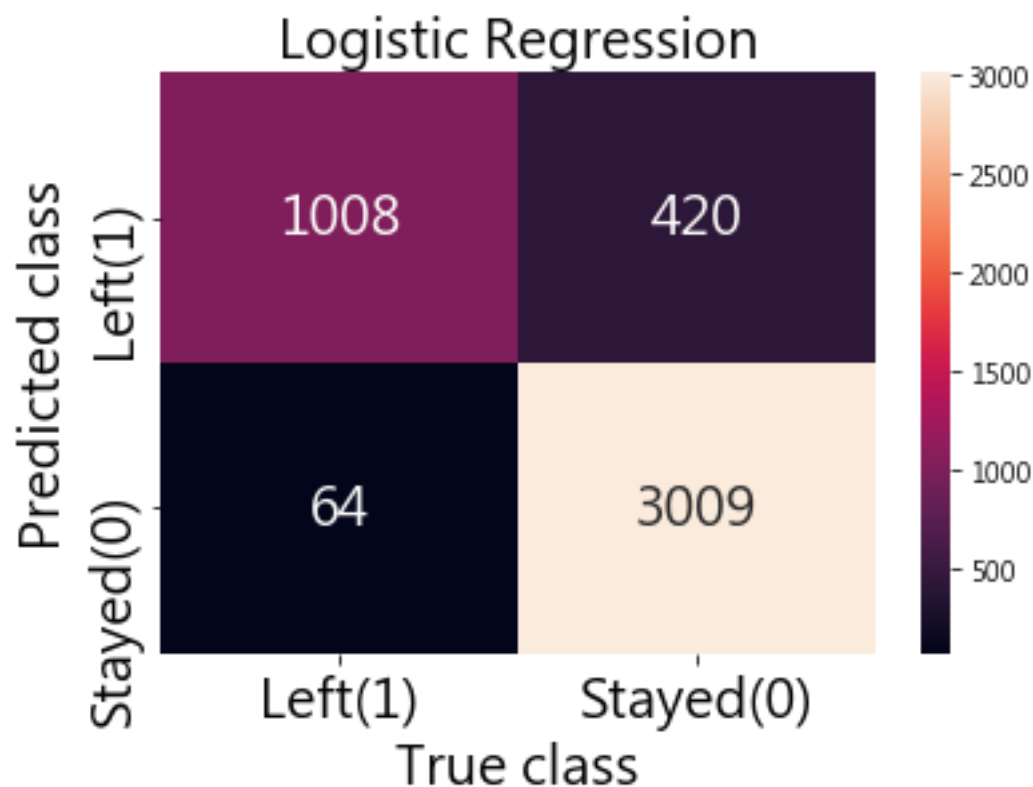
## 模型評估

- 準確度
- 混淆矩陣
- ROC Curve
- 模型解釋

# 準確度/混淆矩陣/ROC Curve

## 模型評估

- 訓練集配適準確度:91.03% ，測試集配適準確度:89.24%
- 敏感度= $1008/(1008+64)=0.94$  ， 特異度= $3009/(3009+420)=0.88$
- AUC:0.91



正影響(值越大，員工越容易離職)

負影響(值越大，員工越不易離職)

number\_project(參與過的專案數)  
average\_monthly\_hours(平均每月工時)  
time\_spend\_company(進入公司的年數)

satisfaction\_level(員工滿意度)

- satisfaction\_level (員工滿意度)又會受到last\_evaluation(績效評估), average\_monthly\_hours(平均每月工時)及time\_spend\_company(進入公司的年數)的影響
- last\_evaluation (績效評估)主要受到number\_project(參與過的專案數),average\_monthly\_hours(平均每月工時)及time\_spend\_company(進入公司的年數)的影響
- 若公司想要挽留人才，可以想辦法降低員工的number\_project(參與過的專案數)及average\_monthly\_hours(平均每月工時)，以提高他們的satisfaction\_level(員工滿意度)，進而降低離職的可能性

**THANK YOU FOR LISTENING**

