

2019

中華郵政大數據競賽



郵政業務量預測與地區分群

下一隊

指導老師：謝邦昌 老師

# 目錄

## CONTENT

1 摘要

2 提案動機

3 處理環境

4 選用模型

5 分析流程說明

6 總結

## » 摘要

做出集群分類，發現未來郵局做決策時，應以六都和非六都作為決策區分。

K-Means  
Clustering

建立業務量預測，可以提早安排人力，使人均處理量相同。

ARIMA

利用現有數據規劃資源分配，例如i郵箱據點等等。

規劃資源  
分配

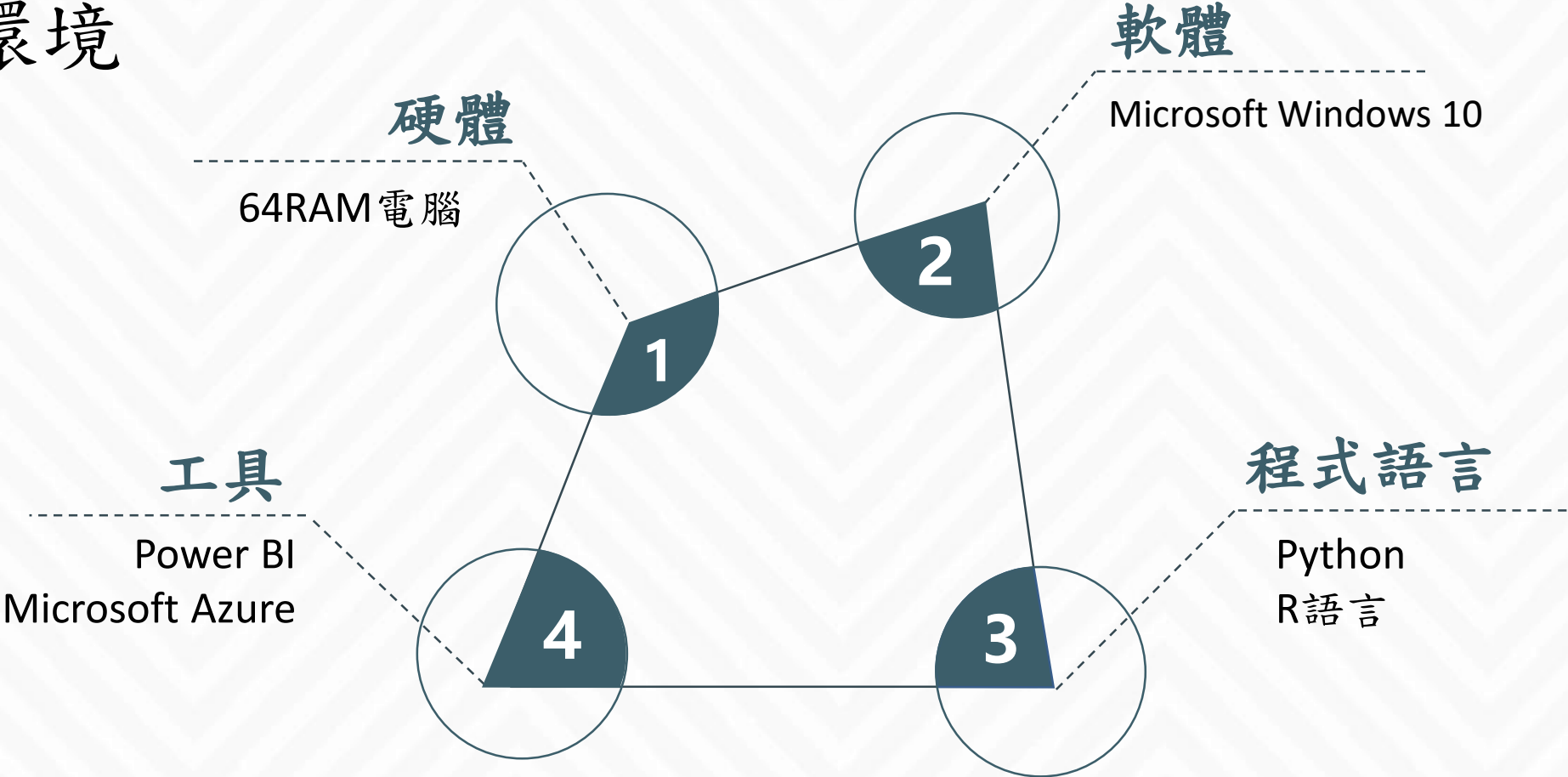
人口資料

透過全國人口資料，篩選出人口數較多村里，將郵務車路徑特別經過此區域的全家，也可以增設i郵箱來分擔郵局處理量。

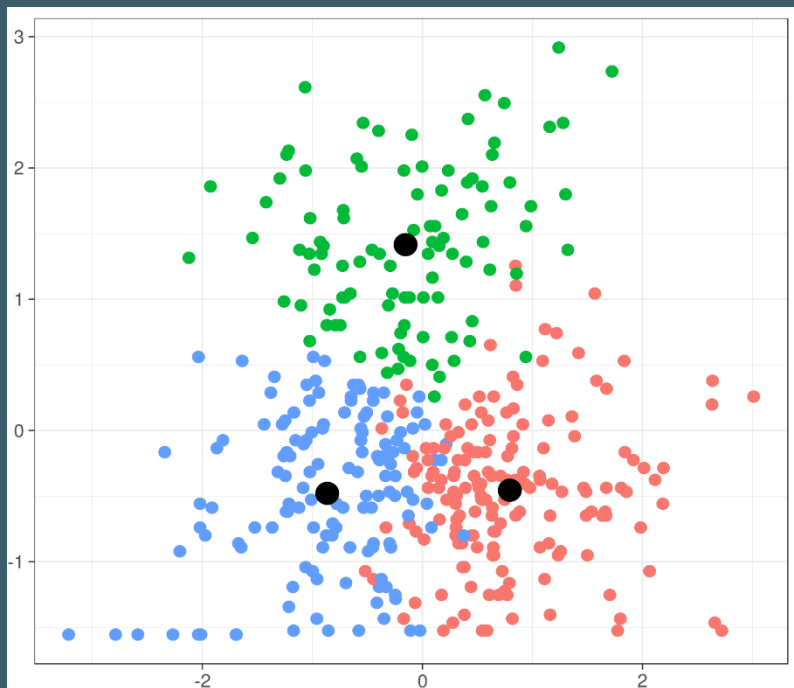
» 提案動機



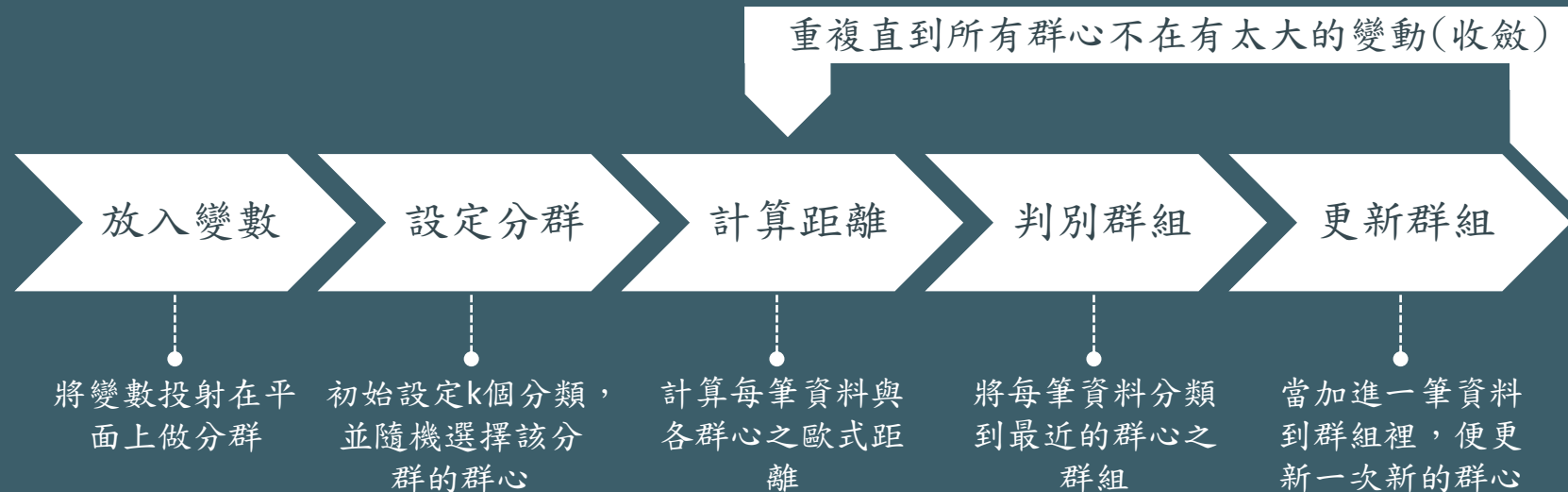
# 處理環境



# » K-Means Clustering

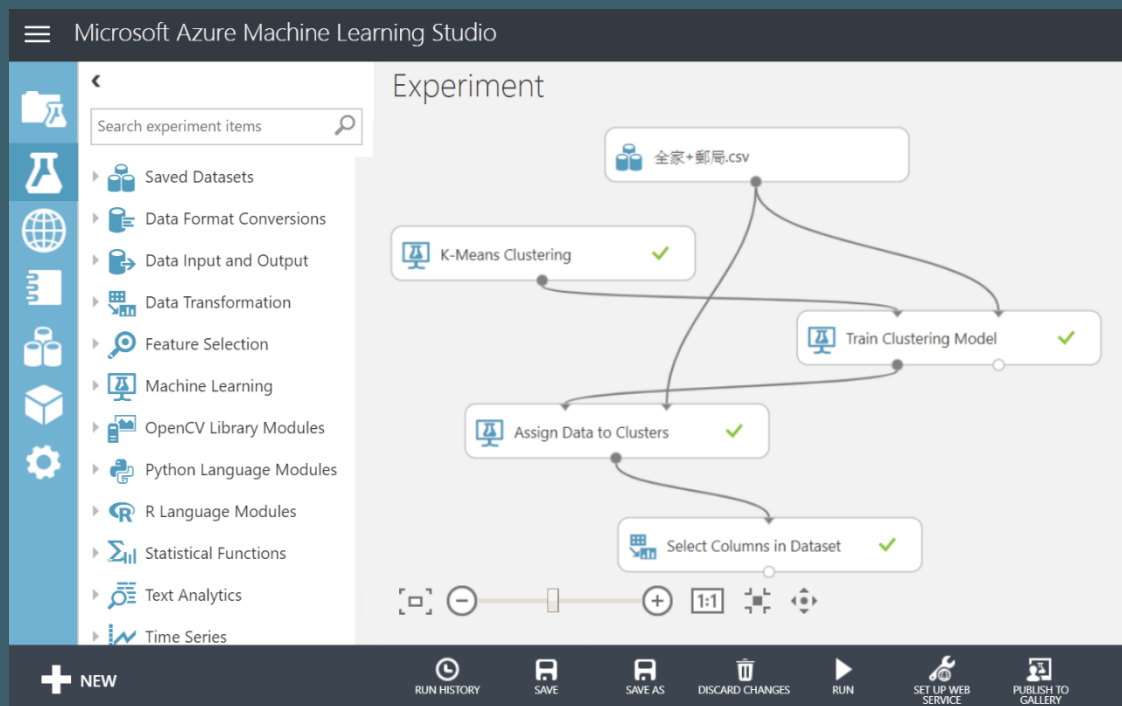


• K-Means概念圖



K-Means 集群分析的概念就是「物以類聚」，是非監督式學習 (Unsupervised learning)，意為得到的資料沒有任何Ground truth，必須透過資料本身去調整分群。

# » K-Means Clustering



- Microsoft Azure Machine Learning 介面

K-Means工具

## Microsoft Azure Machine Learning Studio

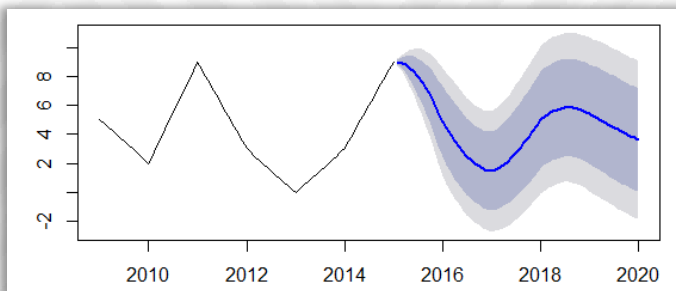
Microsoft Azure Machine Learning Studio採用容易使用的瀏覽器為基底，完全無須撰寫任何程式碼，按幾下就能將構想化成部署。可用來在雲端中快速建置、部署及共用預測性分析解決方案。

我們只要將資料匯入Microsoft Azure Machine Learning，提交指令給機器執行，便能計算出最佳分群結果。



## » ARIMA

ARIMA (p,d,q) 稱為差分自回歸滑動平均模型，  
為時間序列預測分析方法之一



AR為「自回歸」， $p$ 為自回歸階數  
MA為「滑動平均」， $q$ 為滑動平均階數  
 $d$ 為時間序列成為平穩時所做的差分次數

ARIMA (p,d,q) 模型先對階數非平穩的資料歷史數據 $Y_t$ 進行 $d$ 次差分處理，得到新的平穩的數據序列 $X_t$ ，將 $X_t$ 擬合ARMA (p,q) 模型，然後再將原 $d$ 次差分還原，便可以得到 $Y_t$ 的預測數據。



# » 分析流程



## 設定目標

預測業務量以及地區分群。

## 資料採集

郵政內部數據結合外部資料，包含全家、人口、工廠、商店、公司分布，並運用地理資訊圖資雲服務平台轉換經緯度。

## 資料清洗

運用R語言對資料預處理，以便建模。

## 數據建模分析

透過ARIMA時間序列和K-Means建模分析。

## 分析結果

郵政業務量預測與地區分群。

# 》地區分群 – 資料採集與清洗

## 全家地址

使用R語言將全家官網上的地址爬取下來，並將其地址轉為經緯度

## 郵局地址

政府開放資料平台

## 人口分布

中華資料採礦協加上全國達康利用行政院主計處普查資料推估

## 郵務車分布

GPS檔結合郵局地址，得知郵務車所在縣市

## 投遞失敗率

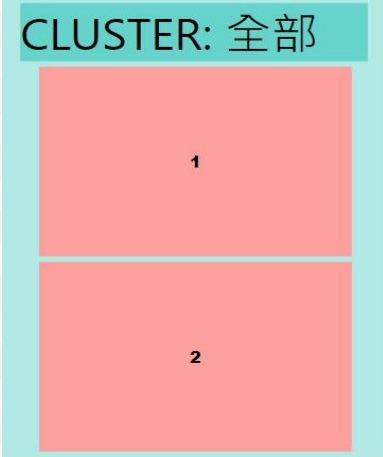
TTS檔中擷取投遞成功和失敗資料，並結合郵局地址將其分出縣市，加以運算

## 工廠、商店、公司

中華資料採礦協會節結合全國達康推估

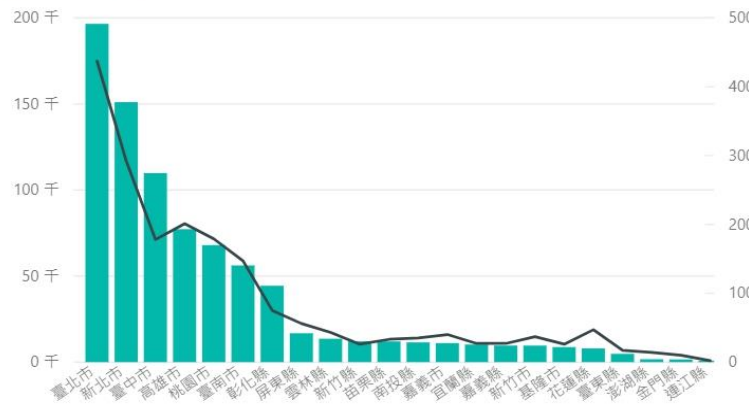
# 地區分群 - 資料發現

臺中市和彰化縣應增加郵務車  
以免造成每車平均郵務量過高

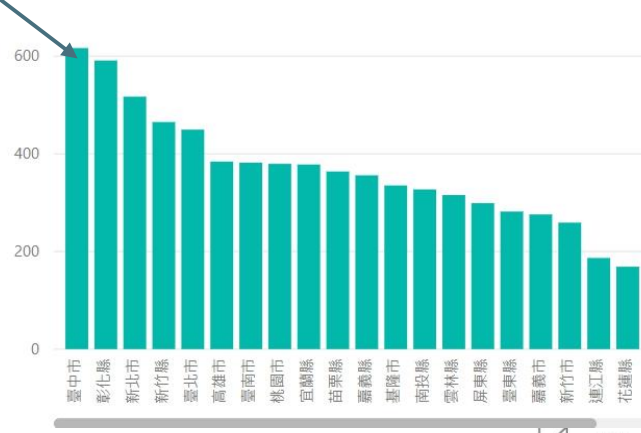


每日平均數量 與 車子數量 依據 縣市

● 每日平均數量 ● 車子數量

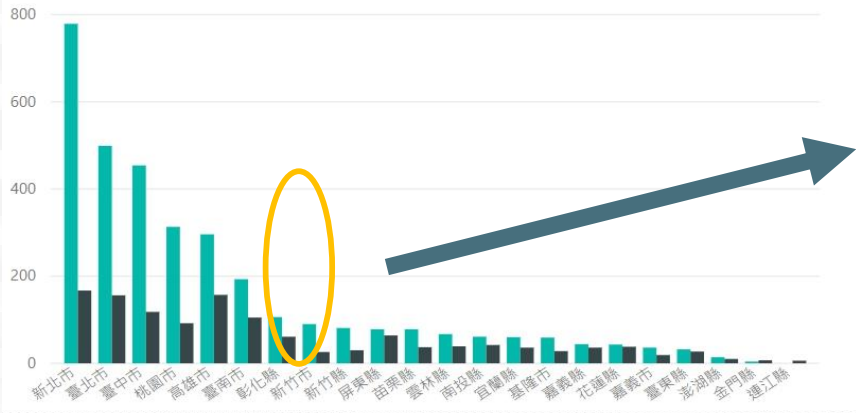


量/車 依據 縣市



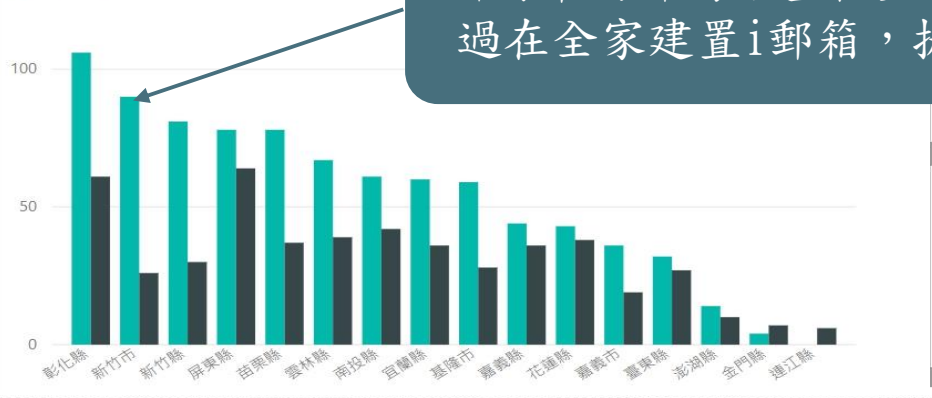
全家數量 與 郵局數量 依據 縣市

● 全家數量 ● 郵局數量



全家數量 與 郵局數量 依據 縣市

● 全家數量 ● 郵局數量



新竹市的郵局數量較少，但是可透過在全家建置i郵箱，提高普及率

# 地區分群 - 資料發現

全國人口數前五大里	人口數
高雄市左營區福山里	42927
高雄市左營區菜公里	34645
高雄市左營區新上里	30362
高雄市鼓山區龍水里	25208
高雄市楠梓區清豐里	25150

圖中為福山里的全家和郵局分佈，可發現這麼大的里，只有2家郵局，可藉由全家來補足；郵務車路徑規劃及i郵箱也可設置在此。

福山里的全家位置  
高雄市左營區福山里58鄰重愛路208號  
高雄市左營區福山里9鄰華夏路1761號  
高雄市左營區福山里60鄰華夏路1168號  
高雄市左營區福山里38鄰榮佑路67號  
高雄市左營區福山里51鄰榮總路225號





# 地區分群 - 資料發現

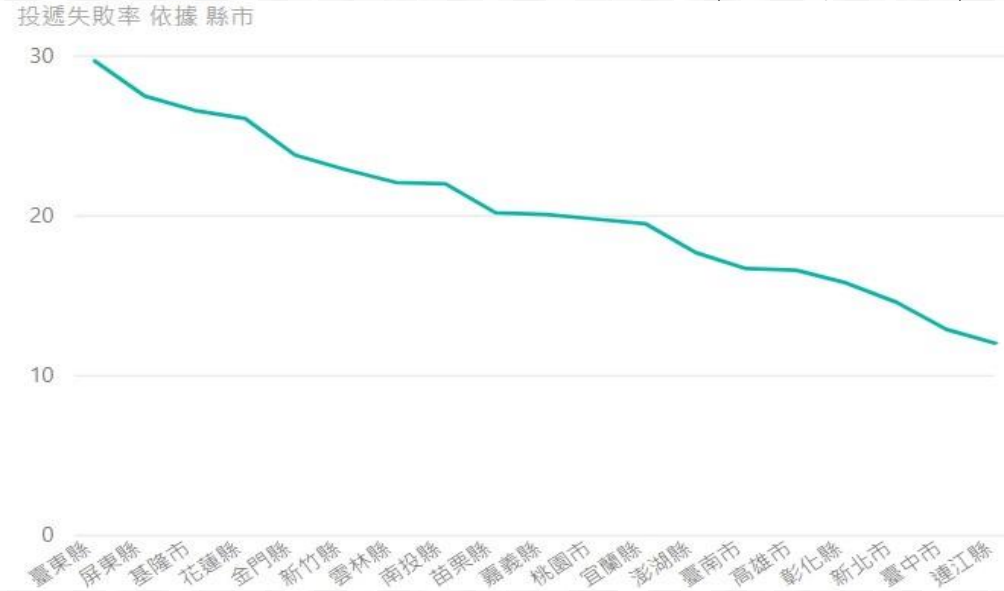
全家、郵局數量跟投遞失敗率



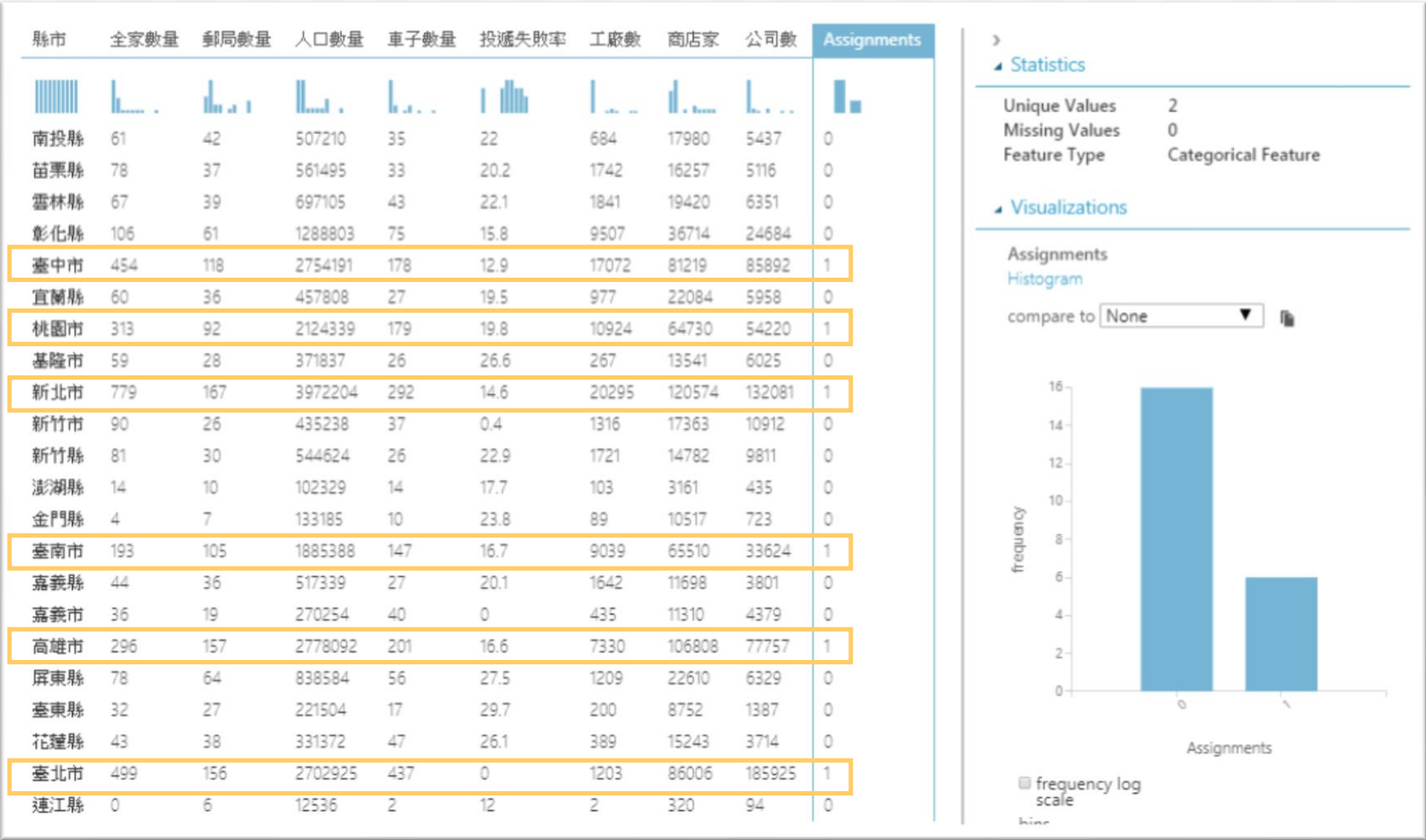
投遞失敗率較高的地方可以藉由全家和郵局共同設置i郵箱，居民領取時間不受限制

投遞失敗率高的地方主要在非六都，可能原因為從事行業的類型及居住型態

各縣市投遞失敗率



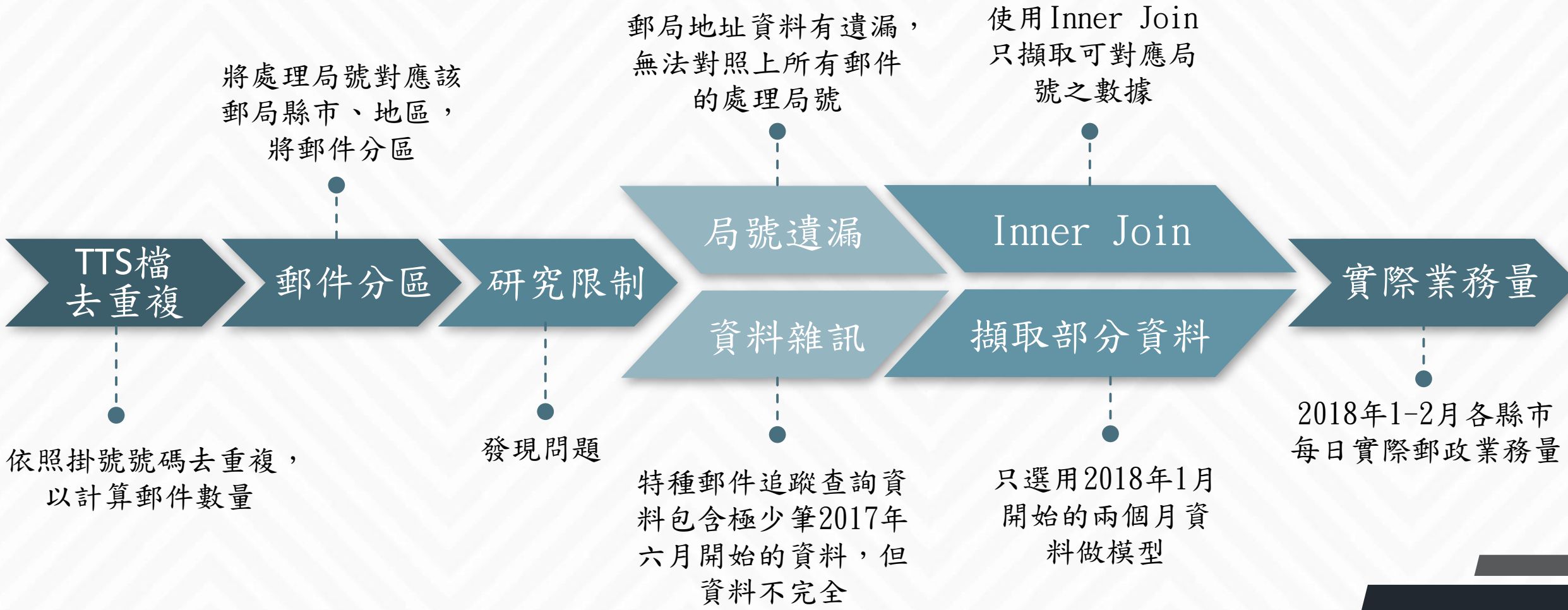
# 地區分群 - 分析結果



將全家數量、郵局數量、人口數量、車子數量、投遞失敗率、工廠數、商店數、公司數匯入 Microsoft Azure Machine Learning，計算出的分群（如左圖），恰好分出六都與非六都。



# 業務量預測 - 資料清洗

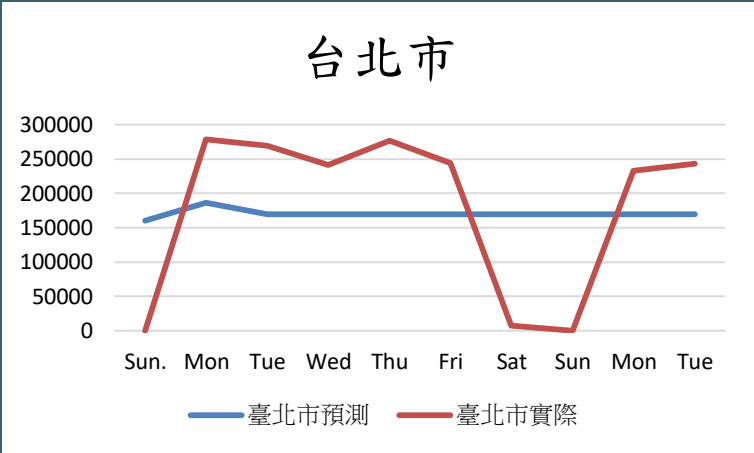
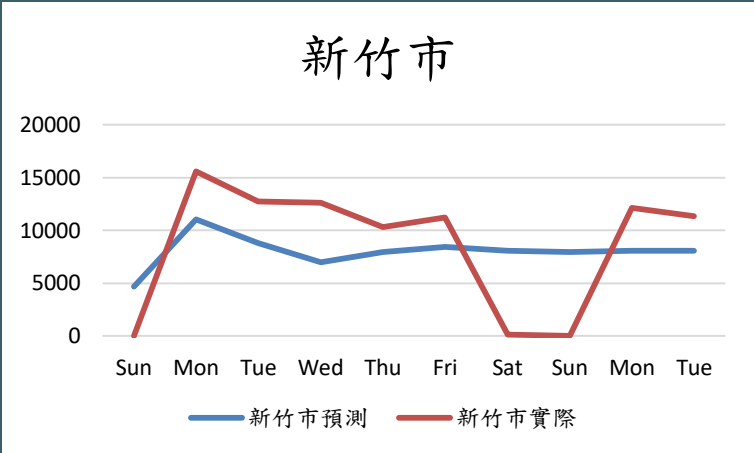


# 業務量預測 - 調整模型

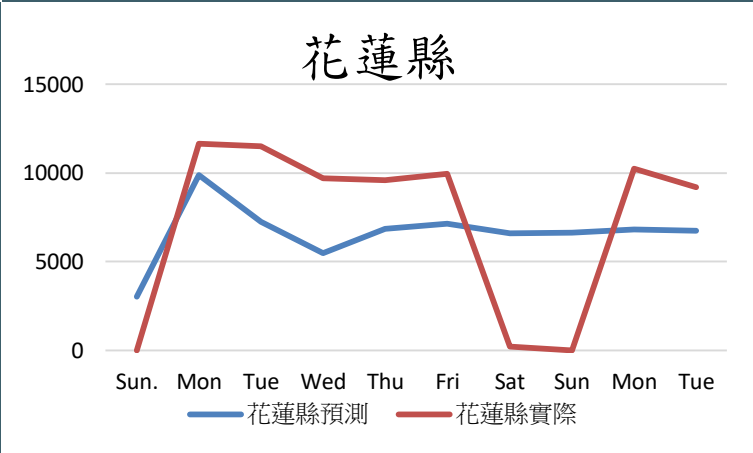
均方誤	高雄市	花蓮縣	基隆市	嘉義市	嘉義縣	金門縣	連江縣	苗栗縣	南投縣	澎湖縣	屏東縣
59天	60729.9	5477.504	6340.904	7745.272	7028.792	1071.549	275.3558	8470.449	8057.513	1062.837	11813.9
63天	44770.01	4070.233	4880.158	6833.327	5171.652	1044.682	232.0883	6517.305	6400.556	1047.199	8977.246
均方誤	臺北市	臺東縣	臺南市	臺中市	桃園市	新北市	新竹市	新竹縣	宜蘭縣	雲林縣	彰化縣
59天	142357.6	3136.975	43505.94	86661.42	51150.54	111141.5	6969.583	8352.504	7095.231	9672.656	35658.24
63天	114377.6	2934.852	30410.74	59928.82	35620.64	83895.38	5055.367	6332.854	5775.613	7261.724	23290.1

原先只選用2018年1月開始的2個月總共59天的資料做模型，但發現業務量以一星期為一個週期，所以我們將模型略做調整，改為以63天的資料做模型，其均方誤便降低許多。另外，我們也試過許多方式調整模型，例如將實際值取log或標準化再預測，最後還是以原始值預測的均方誤為最小。

# 業務量預測 - 預測結果



僅有一季資料，其中又遇二月年假，使非常態性資料佔據不小比例，導致預測偏差



資料不足，故預測值很快就趨向固定值

## 》》 總結

### 地區分群

01

透過K-Means Clustering將台灣分成六都與非六都，未來可運用分群結果決策。



### 業務量預測

02

透過ARIMA時間序列預測業務量，以方便進行郵政人員調動，提高效率。



### 分析建議

03

臺中市和彰化縣可增加郵務車，新竹市可在全家建置i郵箱。



### 未來期望

04

補足遺漏值，分析可以縮小到鄉鎮。資料量從季節變分年，提高準確性。且期望未來能做到預測送達時間。



## » 資料來源

行政院主計總處統計：<https://www.dgbas.gov.tw/np.asp?ctNode=2824>

中華資料採礦協會：<https://www.cdms.org.tw/>

全國達康：<http://www.trend-go.com/>

地理資訊圖資雲服務平台：[https://www.tgos.tw/TGOS/Web/TGOS\\_Home.aspx](https://www.tgos.tw/TGOS/Web/TGOS_Home.aspx)

全家官網：<https://www.family.com.tw/Marketing/index.aspx>

模型理論：

<https://docs.microsoft.com/zh-tw/azure/machine-learning/service/concept-azure-machine-learning-architecture>

<https://medium.com/@chih.sheng.huang821/機器學習-集群分析-k-means-clustering-e608a7fe1b43>

<https://wiki.mbalib.com/zh-tw/ARIMA模型>





**The End**