

# 資料前處理

## 一、連結資料表

經過kaggle Olist資料集的JOIN有關訂單的資料，獲得欄位共65欄，資料共99442筆。

## 二、無效資料處理

我們經過對資料集欄位做敘述統計分析，初步窺探資料分布情形並將極端值刪除後，與小組成員將所有變項逐一理解與討論，將資料分析較無關的欄位刪除，刪除後的剩餘的欄位數量共23個欄位。

## 三、運用舊有欄位產生新欄位

1.[seller\_level]將賣家產品銷售量以Q3、Q1為分界點，將Q3以上的賣家重新編碼為高銷售量賣家[1]；Q1以下的賣家重新編碼為低銷售量賣家[0]，並將其餘Q1-Q3之間的資料刪除，剩餘資料48040筆。

2.將olist平台產品利用單價、交易量進行潛在類別分析(LCA)，參照分析結果將相近的產品進行歸類，產品類型分類由71類別從新分為六大類，分別為時尚配件[1]、3C用品及小型家電[2]、藝文書籍音樂[3]、居家生活及辦公[4]、嬰幼童用品與休閒保健及食品[5]、五金工具及其他[6]。

3.利用RFM分析，R(最近的一次消費)，F(消費頻率)，M(消費金額)來評量客戶潛在價值：

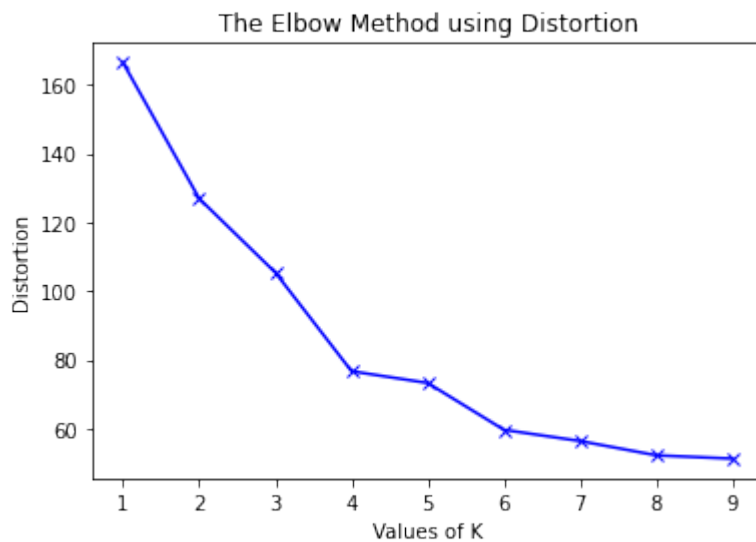
- (1) 最近消費日較近，消費頻率高、總消費金額皆高: 4級客戶
- (2) 最近消費日較遠，消費頻率低，總消費金額高: 3級客戶
- (3) 最近消費日較近，消費頻率高，總消費金額低: 2級客戶
- (4) 最近消費日較遠，消費頻率低，總消費金額低: 1級客戶

4. 利用交易量高低、產品單價高低生成出四大商品型態：

- (1) 交易量大、單價高：獲利商品:4級商品
- (2) 交易量大、單價低：薄利多銷商品:3級商品
- (3) 交易量小、單價高：奢侈商品: 2級商品
- (4) 交易量小、單價低：淘汰商品:1級商品

5. 將[total\_item\_price]、[total\_payment\_value]、[total\_freight\_value] 進行集群處理成新變項[Kmean\_value]:

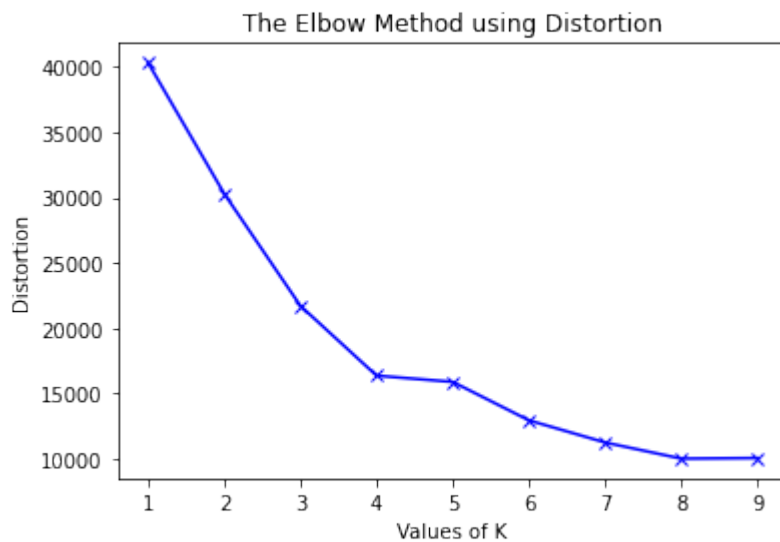
- (1)透過Elbow Method 匯出下圖，發現其值於第4類後的差異趨緩，故選擇分四群。



(2) 透過Kmean集群分析進行非監督式學習獲得新變項[Kmean\_value]。

6. 將[total\_freight\_value]、[total\_package\_volume]、[total\_package\_weight\_g]、[geo\_distance] 進行集群處理成新變項[Kmean\_package]:

(1)透過Elbow Method 匯出下圖，發現其值於第4類後的差異趨緩，故選擇分四群。



(2) 透過Kmean集群分析進行非監督式學習獲得新變項[Kmean\_package]。

## 四、空值資料處理

### 1. 確認資料是否有 NA 值

```
X = pd.read_csv('order_v17_XY.csv', header=0)
print(np.isnan(X).any())
```

```
item_count           False
total_item_price      False
total_payment_value   False
total_freight_value   False
review_score          True
review_answer_waiting_hours  True
```

### 2. 刪除空值

```
X.dropna(inplace=True)
```

```
item_count           False
total_item_price      False
total_payment_value   False
total_freight_value   False
review_score          False
review_answer_waiting_hours  False
```

## 五、資料變數說明

依照欄位內容將資料分為類別型及數值型資料，其變項名稱如下表所示：

類別型資料	數值型資料
item_count	total_item_price
order_purchase_time_day	total_payment_value
is_shipping_delayed	total_freight_value
is_delivered_delayed	review_score
customer_state_region_type	order_purchase_year
seller_state_region_type	order_purchase_month
review_type	order_purchase_day
RFM_type	order_purchase_dayofweek
seller_level	order_purchase_hour
product_type4	until_shipped_waiting_hours
product_category6	until_delivered_waiting_hours
Kmean_price	total_package_volume
Kmean_package	total_package_weight_g
	delivery_efficiency
	total_delivered_waiting_day

	geo_distance
	total_logistics_using_hours
	estimated_logistics_using_hours
	logistics_delay_hours

## 機器學習BOOSTING

	X:26 原始model	X:9	X:10	X:12
模型	測試準確率	測試準確率	測試準確率	測試準確率
C4.5 - CART	0.688	0.689	0.689	0.682
ExtraTrees	0.702	0.694	0.699	0.709
XGBOOST	0.799	0.774	0.774	0.792
CatBOOST	0.726	0.707	0.720	0.728