

電信用戶流失 流失風險預測與因素分析

資料集來源: [Telco Customer Churn](#)



[Linkedin](#)



[104人力銀行](#)

陳勁瑋 製作

AGENDA

1. 專案目標

4. 建造模型

**2. 數據洞
察**

5. 總結

3. 數據處理

6. 附錄

AGENDA

1.

專案目標

4.

建造模型

2.

數據洞
察

5.

總結

3.

策略建議

6.

附錄

專案目標 (情境假設)

假設我任職於一 **電信公司**，今日主管針對 **顧客流失率上升** 的問題，期望我提出改善方案

情境

任職於一電信公司 內，近期 **顧客流失率** 顯著上升

主管好奇顧客流失背後的原因，並期望我提出 **解方**

使用顧客數據進行分析，識別流失 / 未流失顧客的差異

專案脈絡

先理解顧客流失背後可能的原因，並建立 **假說**

針對假說進行驗證，並依據不同客群制定 **專屬** 的解方

專案目標 (現有痛點)

建立一個 **預測模型**，量化各因素對於客 戶**流失的影響程度**，並預測 **客戶的流失率**

現有痛點

電信客 戶逐漸流失，傳統分析方法無法揭示數據背後的深層意義

改善方式

使用 **機器學習** 分析客 戶行為數據，識別高風險流失族群

針對不同流失族群採取專屬的商業策略

專案目標

建立有效的預測模型，分析新顧客的 **未來流失風險**

找出顯著影響顧客是否 **流失的關鍵特徵**

AGENDA

1.

專案目標

4.

建造模型

2.

數據洞
察

5.

總結

3.

數據處理

6.

附錄

資料集介紹

類別變數

Gender

SeniorCitizen

Partner

Dependents

Contract

PhoneService

StreamingTV

MutipleLines

InternetService

OnlineSecurity

OnlineBackup

TechSupport

PaymentMethod

PaperlessBilling

DeviceProtection

StreamingMovies

連續變數

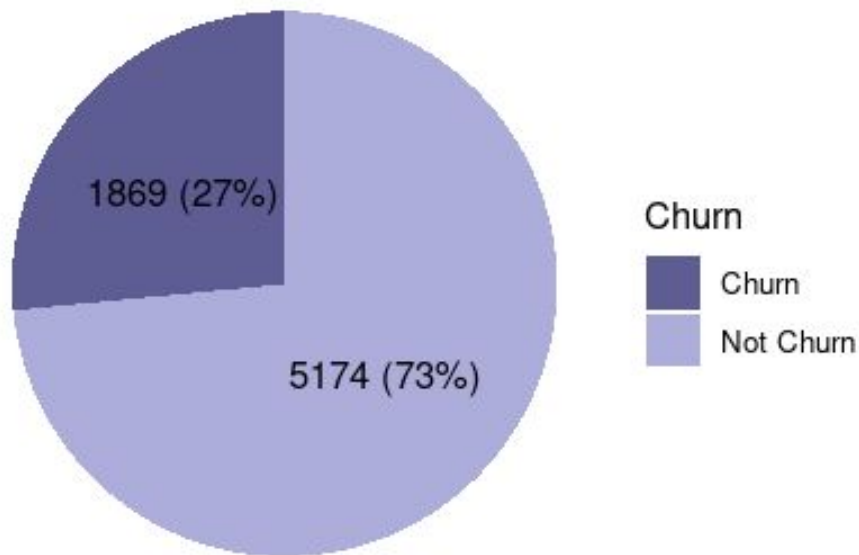
Tenure

MonthlyCharges

TotalCharges

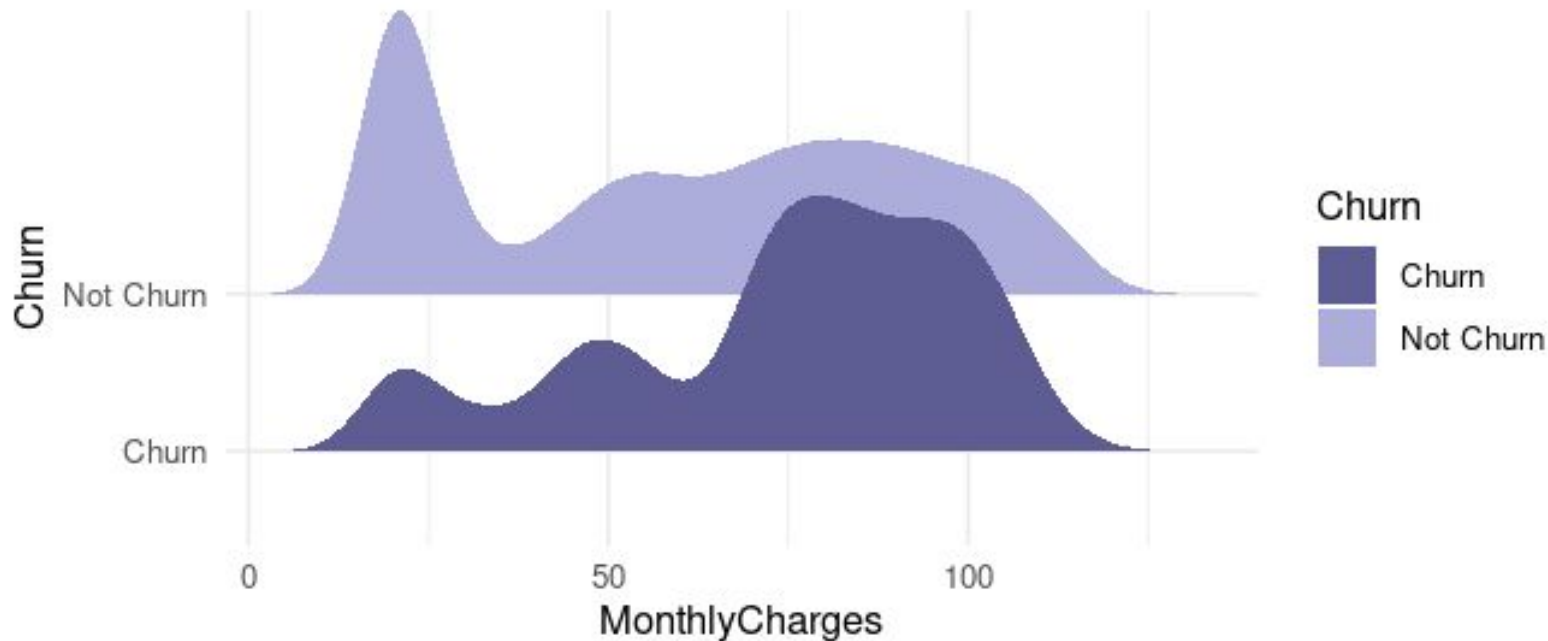
Churn

由下圖發現流失:未流失顧客約為 **1:3**

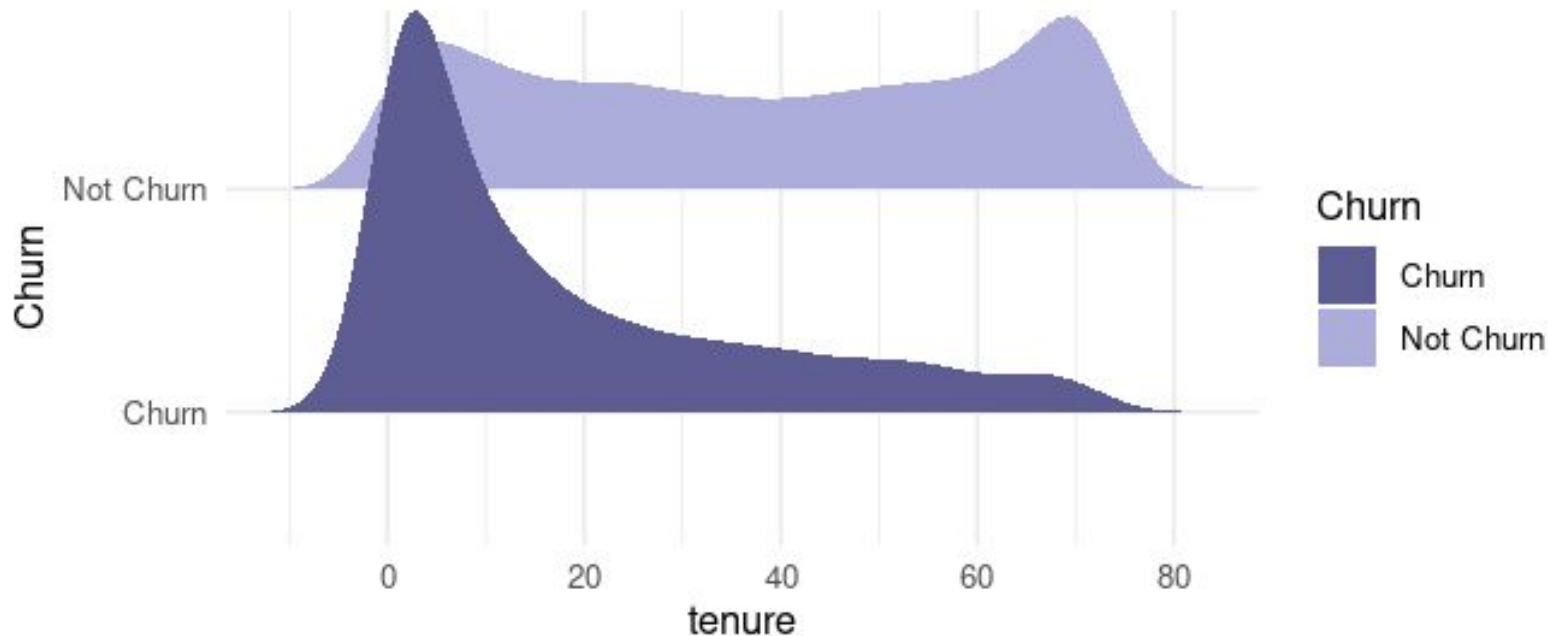


EDA分析

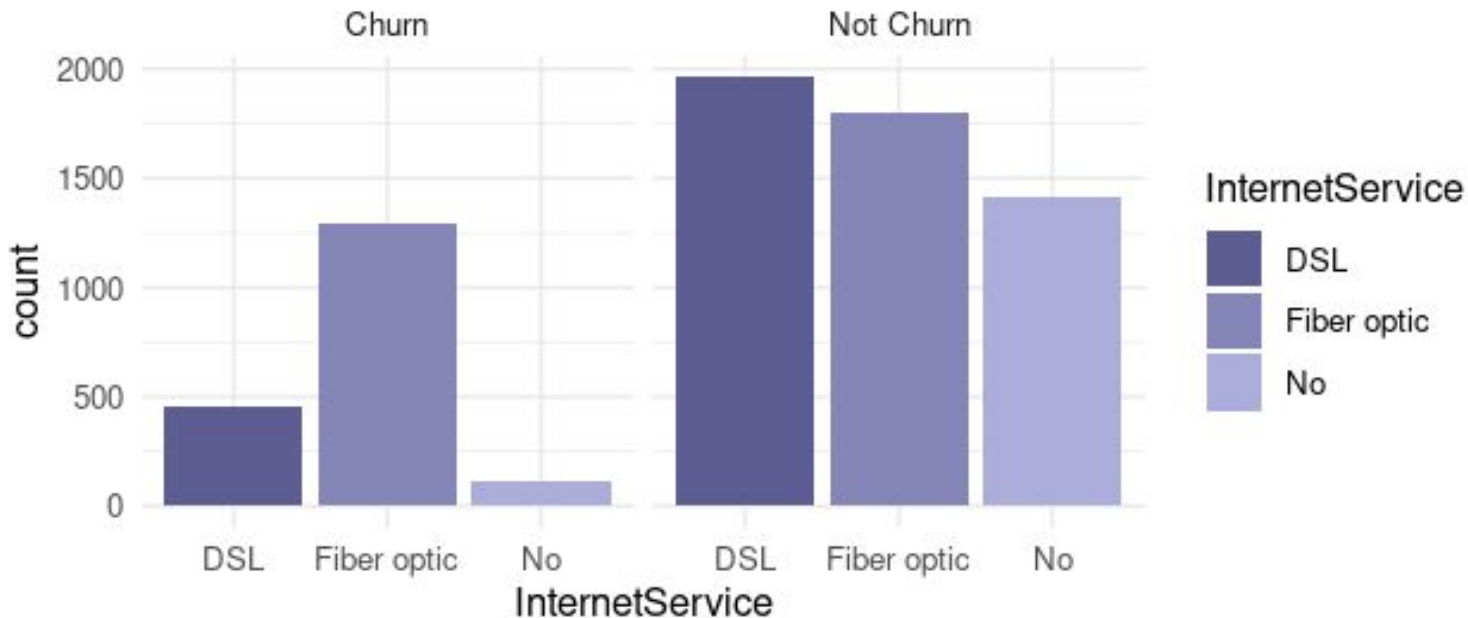
流失、未流失顧客皆呈現 **左偏**，但未流失顧客在較低月支付有 **較高的峰值**

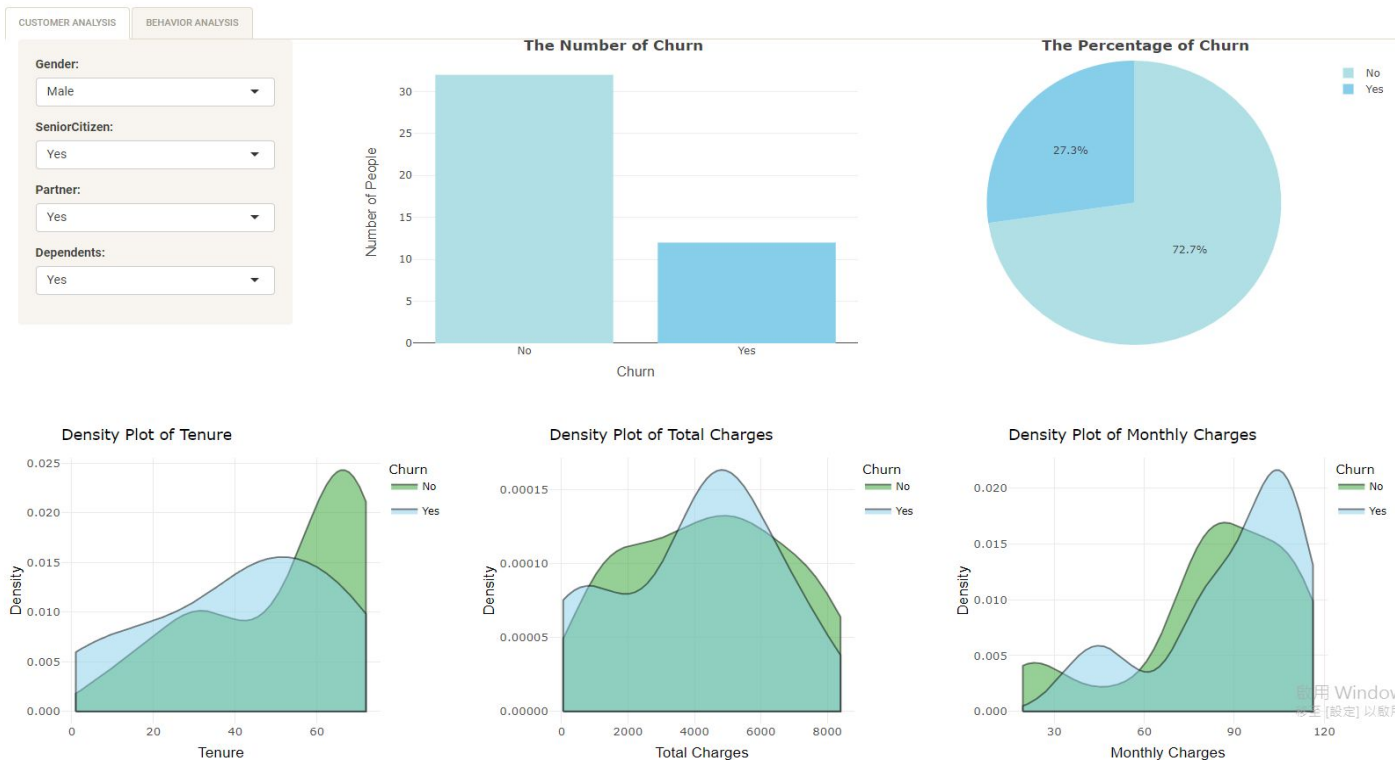


Tenure中的未流失顧客為 **雙峰分布**，流失顧客則呈現 **右偏**



在InternetService的流失顧客中，以 **Fiber optic**最顯著





AGENDA

1.

專案目標

4.

建造模型

2.

數據洞
察

5.

總結

3.

數據處理

6.

附錄

步驟一：去除極端 值

Numeric

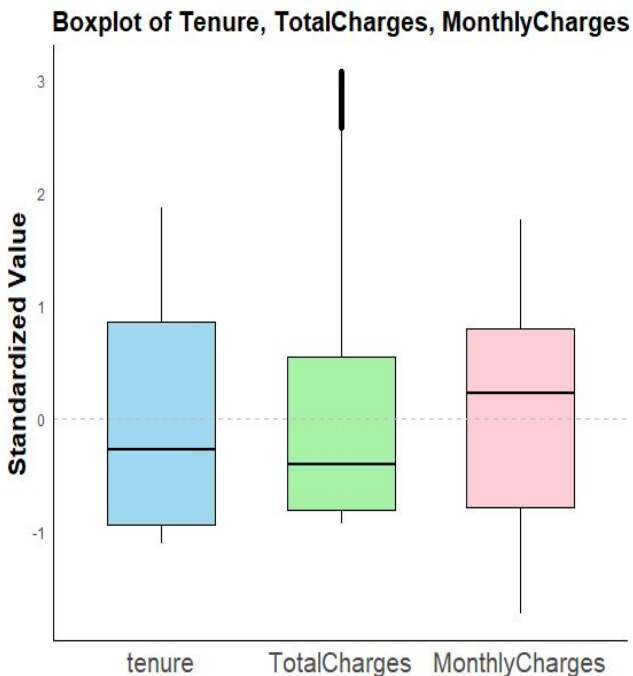
針對連續變數的分析，需考慮是否為極端 值

去除方法

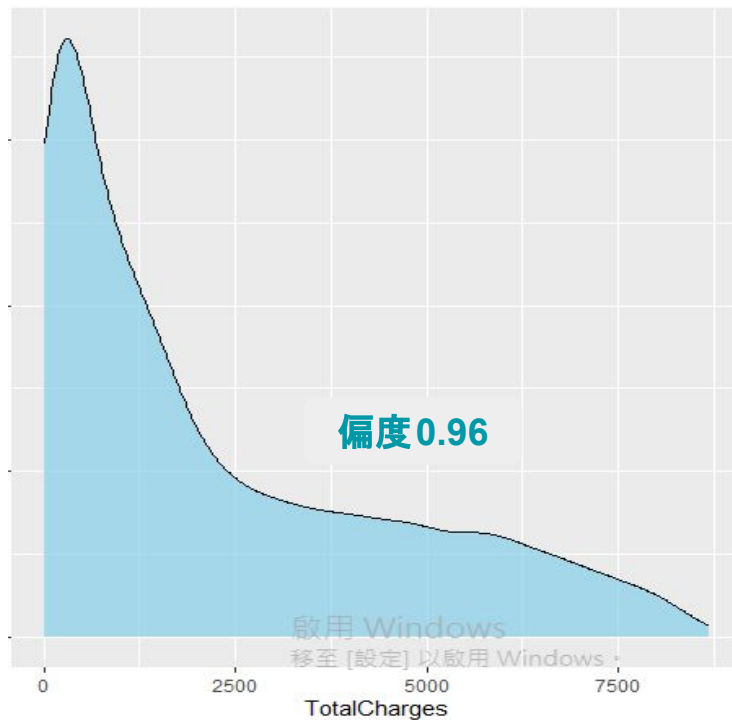
使用IQR法清理資料。當有數 值在第1、3四分位距外的
 $1.5 \times \text{IQR}$ 距離時，視為極端 值

結果

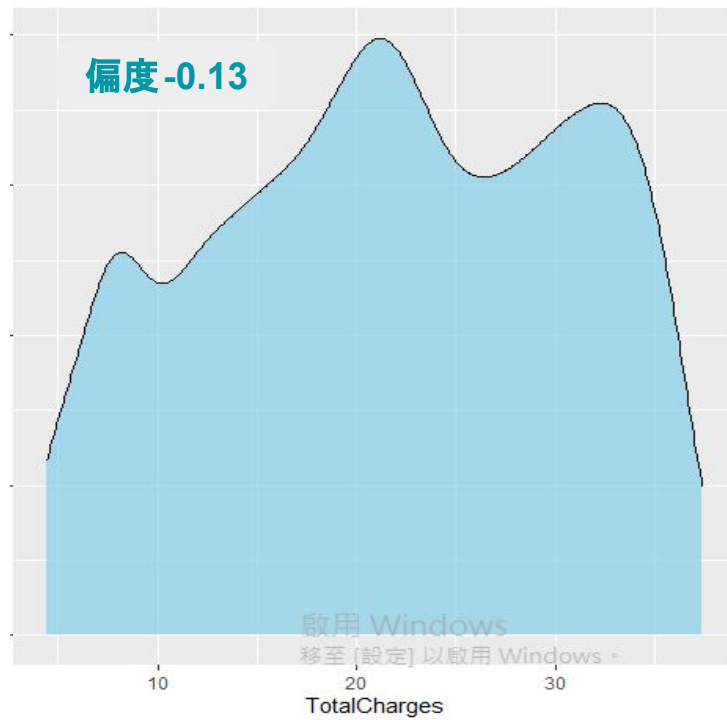
透過IQR法檢驗後，並未發現任何極端 值



步驟二:改善 TotalCharges 偏態問題



Boxcox轉換



步驟三：類別變數處理

現有問題	有部份變數會顯示顧客是否有使用該公司的某項附加服務 例如在 MultipleLines 中，共有 Yes、No、No PhoneService 三類	MultipleLines OnlineSecurity OnlineBackup DeviceProtection TechSupport StreamingTV/Movies
隱患	某些變數可由另一個變數完全表示，引發 共線性 、線性相依問題	
改善方式	將上述 MultipleLines 中的 No 及 No PhoneService 均視為同一類	
預期結果	降低共線性問題，使各個變數間彼此獨立，提升預測準確率	

步驟四：解決目標變數樣本不平衡

現有問題

客戶流失比(Churn)為1:3, 流失的客戶佔比較少

可能使模型 **隨機猜測** 仍保有高準確率, 但Recall值低

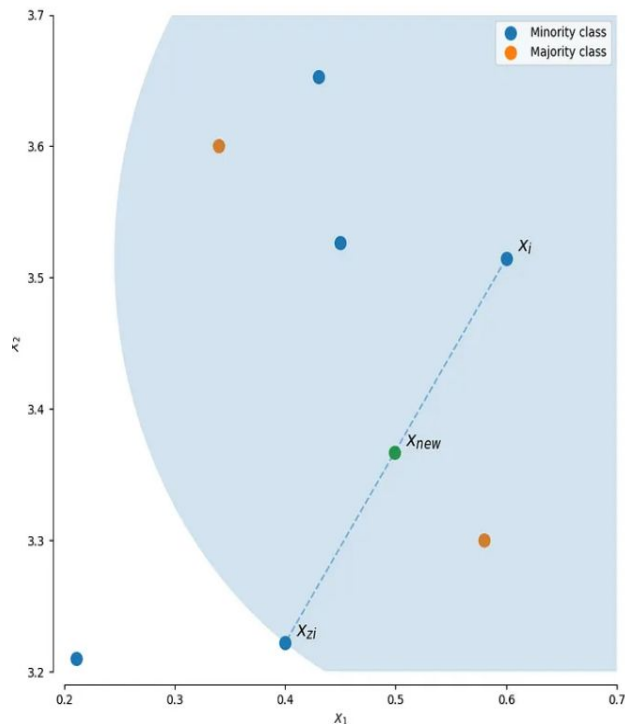
改善方式

因數據量有限, 採取 OVERSAMPLE方式, 降低模型 Bias

使用 **SMOTE-NC** 處理包含連續、類別變數的資料集

預期結果

透過OVERSAMPLE解決樣本不平衡可能造成的隱患,
使模型的預測效率以及 Recall值提高



AGENDA

1.

專案目標

4.

建造模型

2.

數據洞
察

5.

總結

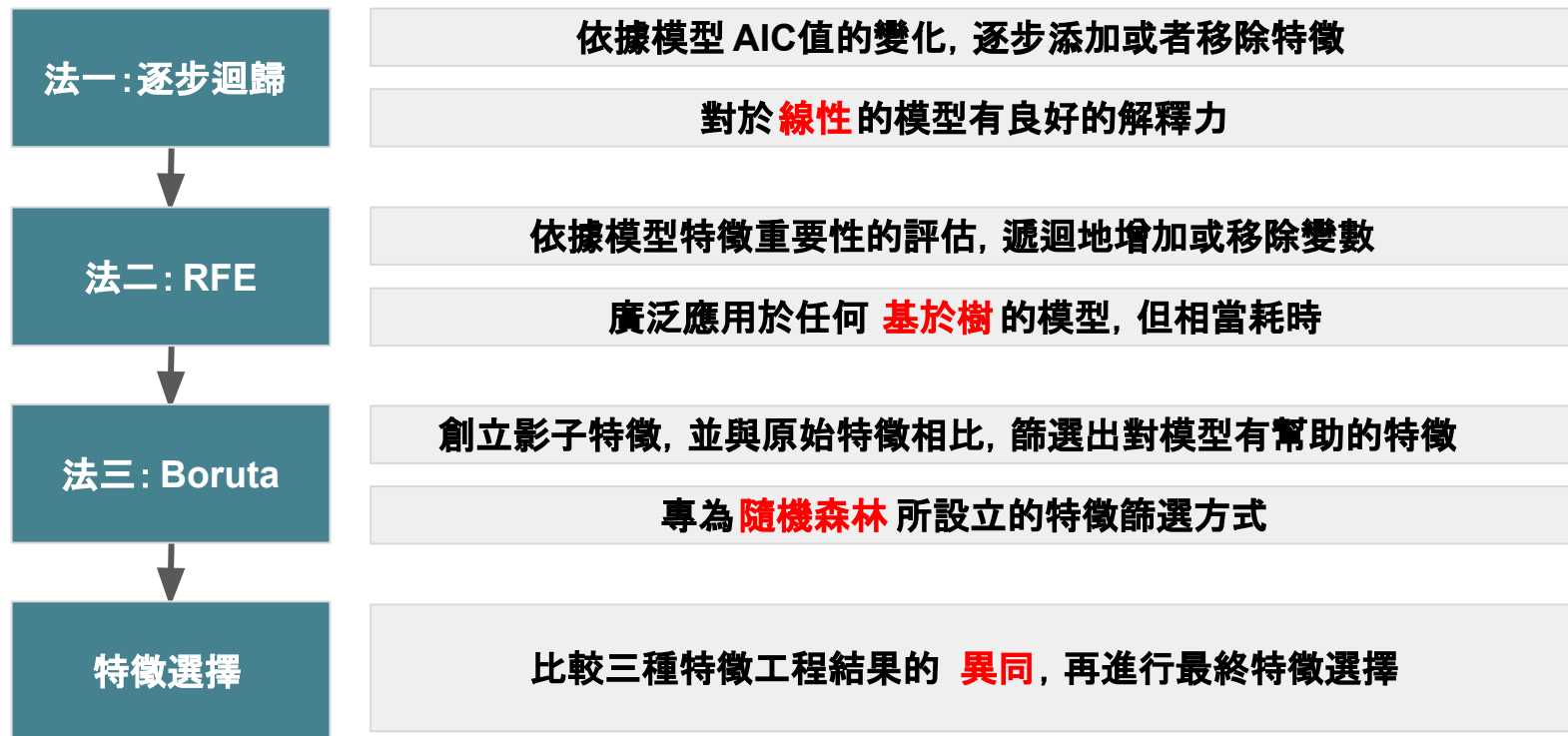
3.

數據處理

6.

附錄

使用三種特徵工程方法，對比結果後進行特徵選擇



特徵工程方法一（逐步迴歸）

根據羅吉斯迴歸的 VIF 值，移除產生共線性的特徵

自變數	GVIF
MonthlyCharges	85.9
InternetServices	50.0
TotalCharges	20.7
Tenure	15.5
Others	<10

Monthly
Charges

自變數	GVIF
TotalCharges	20.2
Tenure	15.2
Others	<10

Total
Charges

模型不存在
共線性

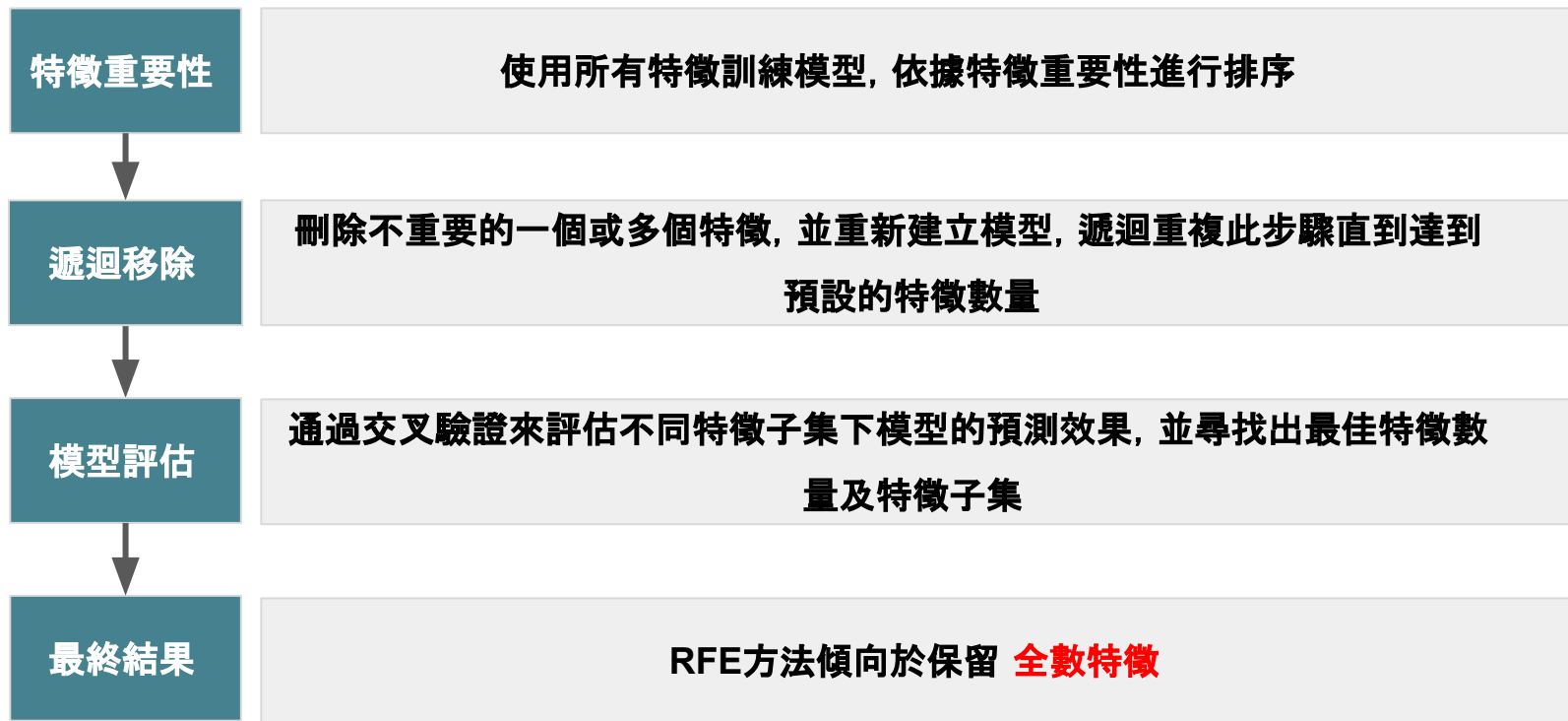
特徵工程方法一（逐步迴歸）

根據前一步驟遺留下來的變數，進行羅吉斯逐步迴歸



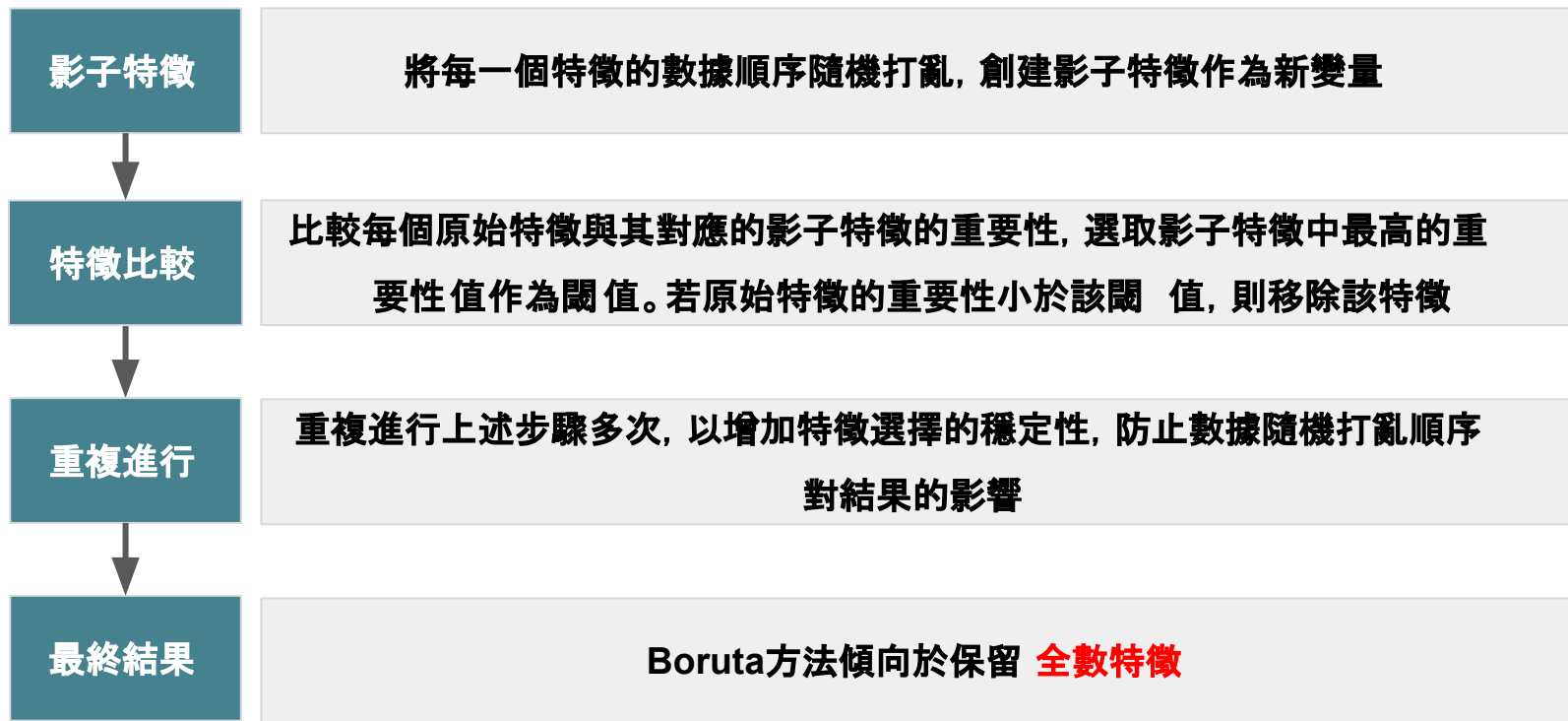
特徵工程方法二 (Recursive Feature Elimination)

依據隨機森林模型特徵的 **重要性**，遞迴的移除不重要的變數



特徵工程方法三 (Boruta)

在隨機森林模型中以創造 **影子特徵** 的方式，與原始特徵進行比較並保留重要的特徵



使用逐步迴歸的結果選擇特徵，將特徵數量由 19降至 13

方法異同

逐步迴歸：在移除不顯著的變數後，模型保留了 13個特徵

RFE及Boruta：均傾向於保留全部的特徵

可能隱患

複雜的模型可能導致 過度擬合 問題，這使得模型在面對新數據時預測能力可能變差

保留所有變數會增加模型計算的時間成本，尤其對大數據集來說極其耗時

最終抉擇

選擇逐步迴歸的特徵選擇方法，因為它在統計上保持 顯著性，同時計算時間較為節省

儘管RFE降至13個變數僅損失約 1%的準確度，考慮到以上問題，因此決定以此作為平衡的選擇

建造隨機森林模型，並結合交叉驗證尋找最佳超參數

訓練模型	結果
ntree	500
mtry	6
splitrule	gini
min.node.size	1
Precision	89.5%
Recall	97.3%

	Yes	No
Yes	4020	111
No	474	3657

測試模型	結果
Accuracy	82.5%
Precision	80.0%
Recall	86.6%
Sensitivity	86.6%
Specificity	78.3%
F1 Score	83.2%

	Yes	No
Yes	894	138
No	224	808

建造XGboost模型，並結合交叉驗證尋找最佳超參數

訓練模型	結果
nrounds	150
max_depth	6
eta	0.3
gamma	0.1
colsample bytree	0.8
Precision	91.2%
Recall	90.0%

	Yes	No
Yes	3700	431
No	297	3834

測試模型	結果
Accuracy	83.9%
Precision	84.2%
Recall	83.3%
Sensitivity	83.6%
Specificity	84.3%
F1 Score	83.9%

	Yes	No
Yes	863	169
No	162	870

AGENDA

1.

專案目標

4.

建造模型

2.

數據洞
察

5.

總結

3.

數據處理

6.

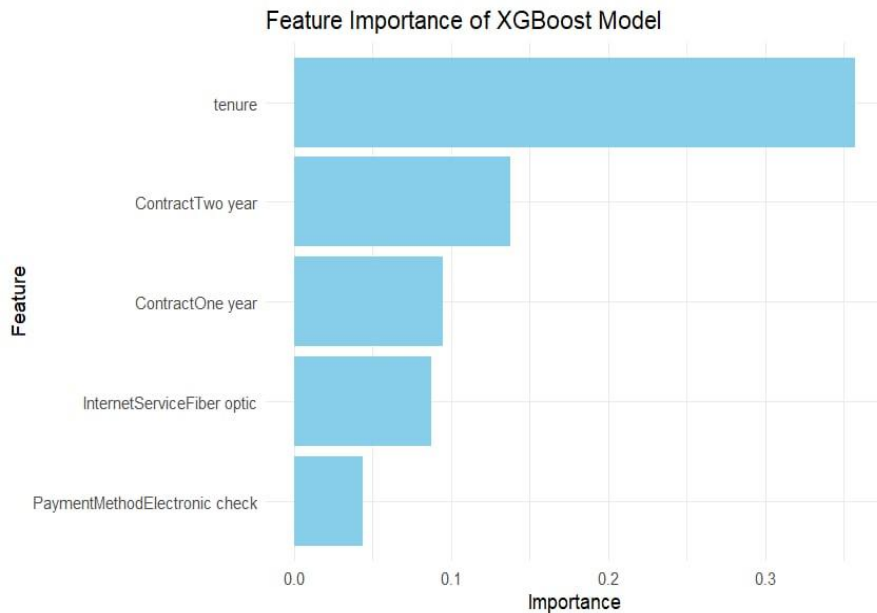
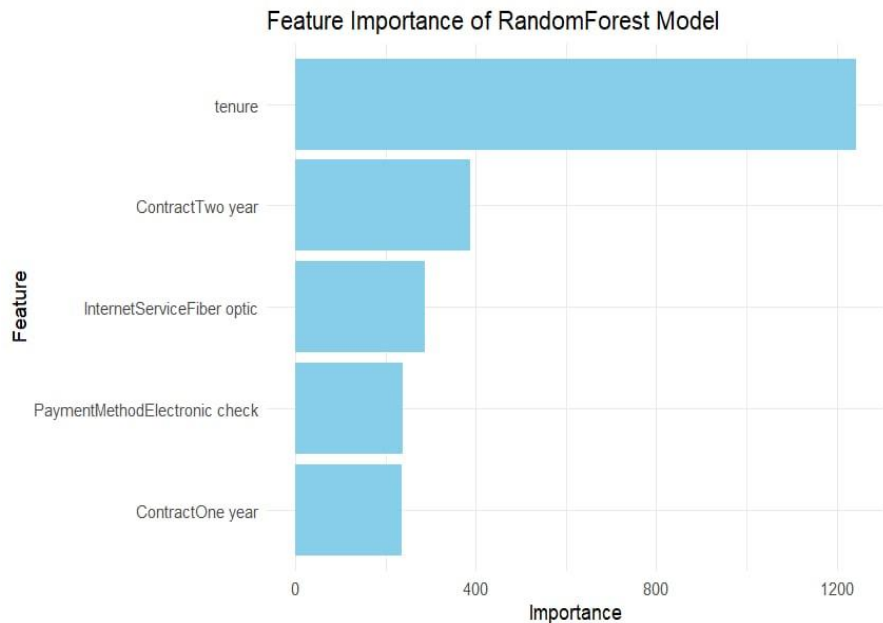
附錄

總結 – 模型比較 (Null Model未進行 Oversample處理, 且使用全變數進行羅吉斯迴歸)

隨機森林與 XGboost皆是優良模型 – 提升原始模型的 精確率及 召回率

Null Model			RandomForest			XGboost		
	Yes	No		Yes	No		Yes	No
Yes	216	157	Yes	894	138	Yes	863	169
No	111	921	No	224	808	No	162	870
Precision	67.1%		Precision	80.0%		Precision	84.2%	
Recall	57.9%		Recall	86.6%		Recall	83.6%	
F1 Score	62.2%		F1 Score	83.2%		F1 Score	83.9%	
Accuracy	80.9%		Accuracy	82.5%		Accuracy	83.9%	

預測顧客是否流失的重要指標



預測顧客是否流失的重要指標

相關性	Random Forest	XGboost
負/負	Tenure	Tenure
負/負	Contract Two Year	Contract Two Year
正/負	InternetService Fiber optic	Contract One Year
正/正	PaymentMethod Electronic check	InternetService Fiber optic
負/正	Contract One Year	PaymentMethod Electronic check

羅吉斯逐步回歸幫助鑑定變數相關性

若係數為正則代表正相關

若係數為負則代表負相關

預測顧客是否流失的重要指標 – FAMD分析



總結 – Recap

分析總結

預測效率

預測未來用戶流失的準確率約為 **83%**，且能夠以 **86%** 的準確度 **識別** 出流失顧客

對於高風險流失顧客，可 **提前採取** 專業策略以防止顧客流失

Tenure/Charges

使用期間長/總花費越高的顧客 **流失率較低**，可能與顧客忠誠度、轉換成本、顧客慣性有關

Contract

簽訂兩年或一年合約的顧客 **流失率較低**，可能與契約約束、優惠折扣及服務穩定性有關

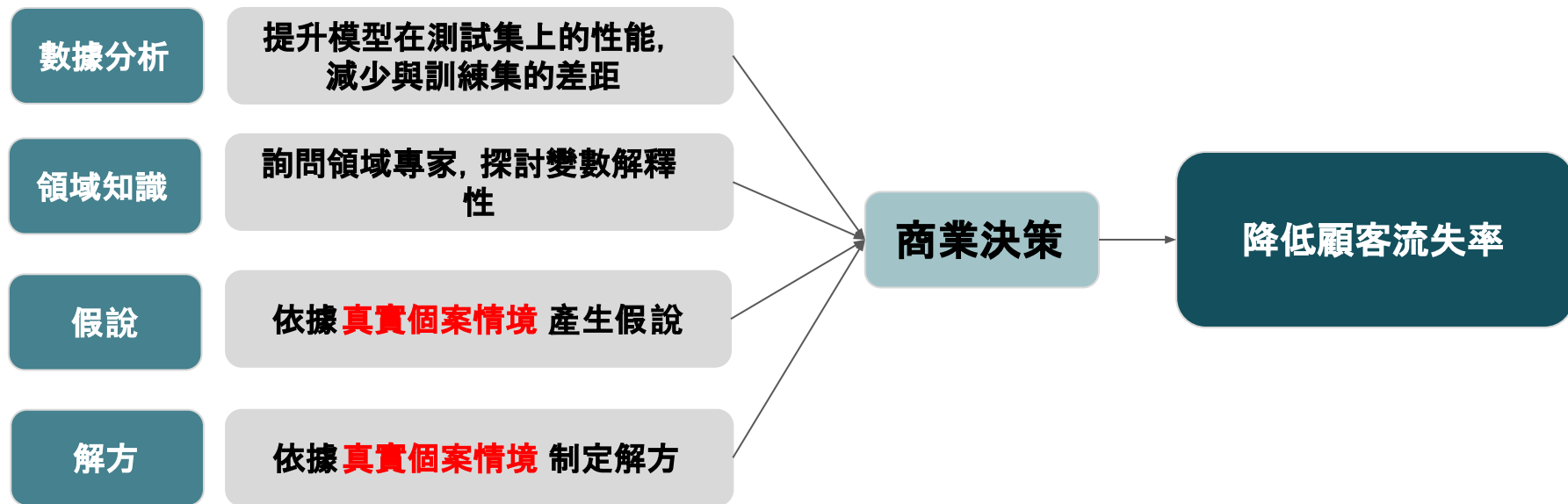
InternetService

使用光纖網路的顧客 **流失率較高**，可能與服務問題、成本高昂及市場競爭激烈有關

PaymentMethod

使用電子支付的顧客 **流失率較高**，可能與支付過程不流暢及代收費用高昂有關

Next Step - 結合領域知識及商業決策



AGENDA

1.

專案目標

4.

成效分析

2.

數據洞
察

5.

總結

3.

策略建議

6.

附錄

羅吉斯逐步迴歸係數係數

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	0.570976	0.128203	4.454	8.44e-06	***
DependentsYes	-0.373902	0.070747	-5.285	1.26e-07	***
tenure	-0.032374	0.002024	-15.998	< 2e-16	***
PhoneServiceYes	-0.347918	0.110425	-3.151	0.001629	**
MultipleLinesYes	0.252633	0.072303	3.494	0.000476	***
`InternetServiceFiber optic`	0.938935	0.082032	11.446	< 2e-16	***
InternetServiceNo	-0.766592	0.107698	-7.118	1.10e-12	***
OnlineSecurityYes	-0.524045	0.075367	-6.953	3.57e-12	***
OnlineBackupYes	-0.134981	0.068992	-1.956	0.050409	.
TechSupportYes	-0.467485	0.076298	-6.127	8.95e-10	***
StreamingTVYes	0.239726	0.074792	3.205	0.001349	**
StreamingMoviesYes	0.331994	0.074307	4.468	7.90e-06	***
`ContractOne year`	-0.803014	0.090763	-8.847	< 2e-16	***
`ContractTwo year`	-1.463869	0.140579	-10.413	< 2e-16	***
PaperlessBillingYes	0.393372	0.065184	6.035	1.59e-09	***
`PaymentMethodCredit card (automatic)`	0.075733	0.099828	0.759	0.448066	
`PaymentMethodElectronic check`	0.498962	0.083925	5.945	2.76e-09	***
`PaymentMethodMailed check`	0.073753	0.099276	0.743	0.457538	

6.1 SMOTE-NC

While our SMOTE approach currently does not handle data sets with all nominal features, it was generalized to handle mixed datasets of continuous and nominal features. We call this approach Synthetic Minority Over-sampling TEchnique-Nominal Continuous [SMOTE-NC]. We tested this approach on the Adult dataset from the UCI repository. The SMOTE-NC algorithm is described below.

1. Median computation: Compute the median of standard deviations of all continuous features for the minority class. If the nominal features differ between a sample and its potential nearest neighbors, then this median is included in the Euclidean distance computation. We use median to penalize the difference of nominal features by an amount that is related to the typical difference in continuous feature values.
2. Nearest neighbor computation: Compute the Euclidean distance between the feature vector for which k-nearest neighbors are being identified (minority class sample) and the other feature vectors (minority class samples) using the continuous feature space. For every differing nominal feature between the considered feature vector and its potential nearest-neighbor, include the median of the standard deviations previously computed, in the Euclidean distance computation. Table 2 demonstrates an example.

F1 = 1 2 3 A B C [Let this be the sample for which we are computing nearest neighbors]

F2 = 4 6 5 A D E

F3 = 3 5 6 A B K

So, Euclidean Distance between F2 and F1 would be:

$$\text{Eucl} = \sqrt{(4-1)^2 + (6-2)^2 + (5-3)^2 + \text{Med}^2 + \text{Med}^2}$$

Med is the median of the standard deviations of continuous features of the minority class.

The median term is included twice for feature numbers 5: B→D and 6: C→E, which differ for the two feature vectors: F1 and F2.

3. Populate the synthetic sample: The continuous features of the new synthetic minority class sample are created using the same approach of SMOTE as described earlier. The nominal feature is given the value occurring in the majority of the k-nearest neighbors.

The SMOTE-NC experiments reported here are set up the same as those with SMOTE, except for the fact that we examine one dataset only. SMOTE-NC with the Adult dataset differs from our typical result: it performs worse than plain under-sampling based on AUC, as shown in Figures 26 and 27. We extracted only continuous features to separate the effect of SMOTE and SMOTE-NC on this dataset, and to determine whether this oddity was due to our handling of nominal features. As shown in Figure 28, even SMOTE with only continuous features applied to the Adult dataset, does not achieve any better performance than plain under-sampling. Some of the minority class continuous features have a very high variance, so, the synthetic generation of minority class samples could be overlapping with the majority class space, thus leading to more false positives than plain under-sampling. This hypothesis is also supported by the decreased AUC measure as we SMOTE at degrees greater than 50%. The higher degrees of SMOTE lead to more minority class samples in the dataset, and thus a greater overlap with the majority class decision space.

<https://www3.nd.edu/~dial/publications/chawla2002smote.p>

Thank you for watching