

學號:0411509 姓名:許家維

採用數據: 電網穩定性模擬數據

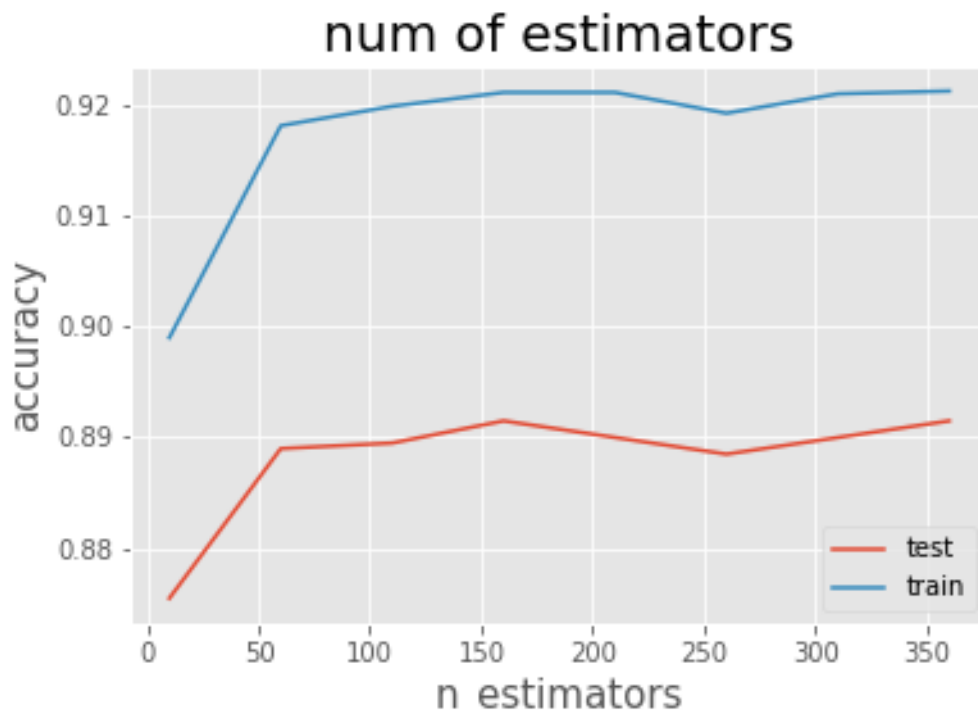
其中包含了(消耗功率、生產電力、反應時間等等), 並預測分類該電網是否穩定。

數據來源:

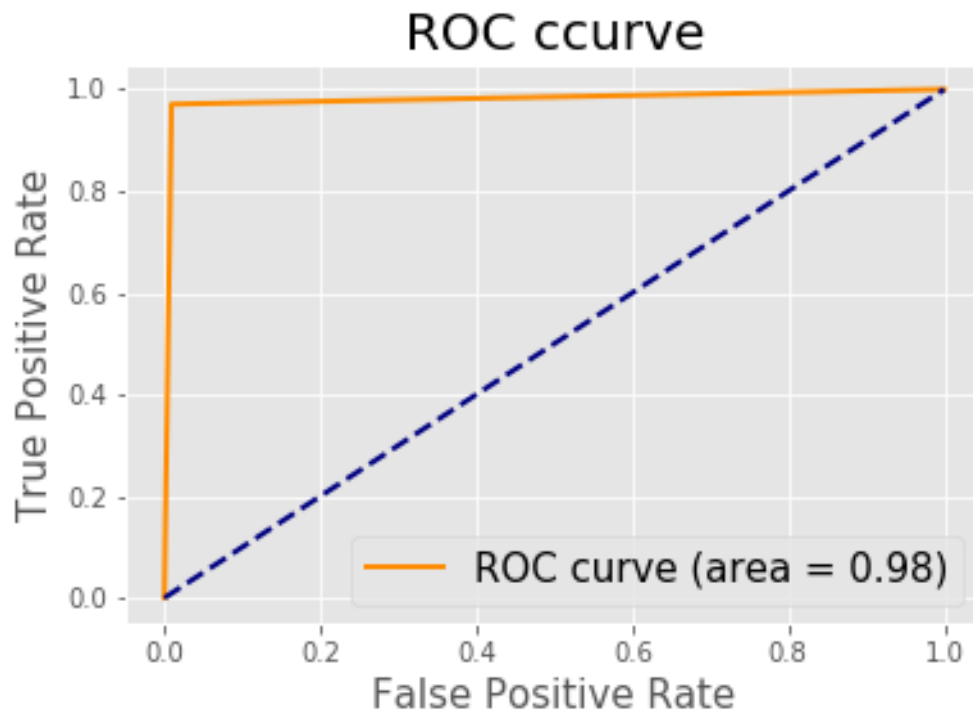
<http://archive.ics.uci.edu/ml/datasets/Electrical+Grid+Stability+Simulted+Data+#>

分析過程

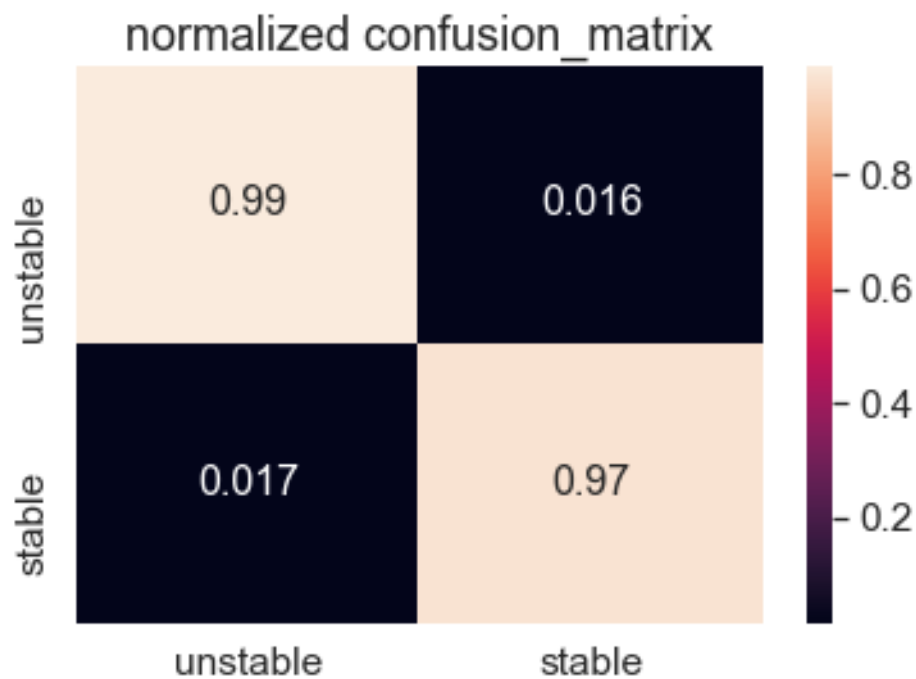
此次數據採用 Random Forest 分類, 其中是利用 ensembl 多個決策樹達成分類預測的, 其中若其中樹木數量是一個很重要的超參數, 太多容易造成過擬和 (overfitting), 過低則會有欠擬和 (under fitting), 透過逐一分析, 可發現在 $n=150$ 下有較良好的泛化能力並不失去預測準確度, 在不失去 testing data 準確度下提高 training data 準確度。



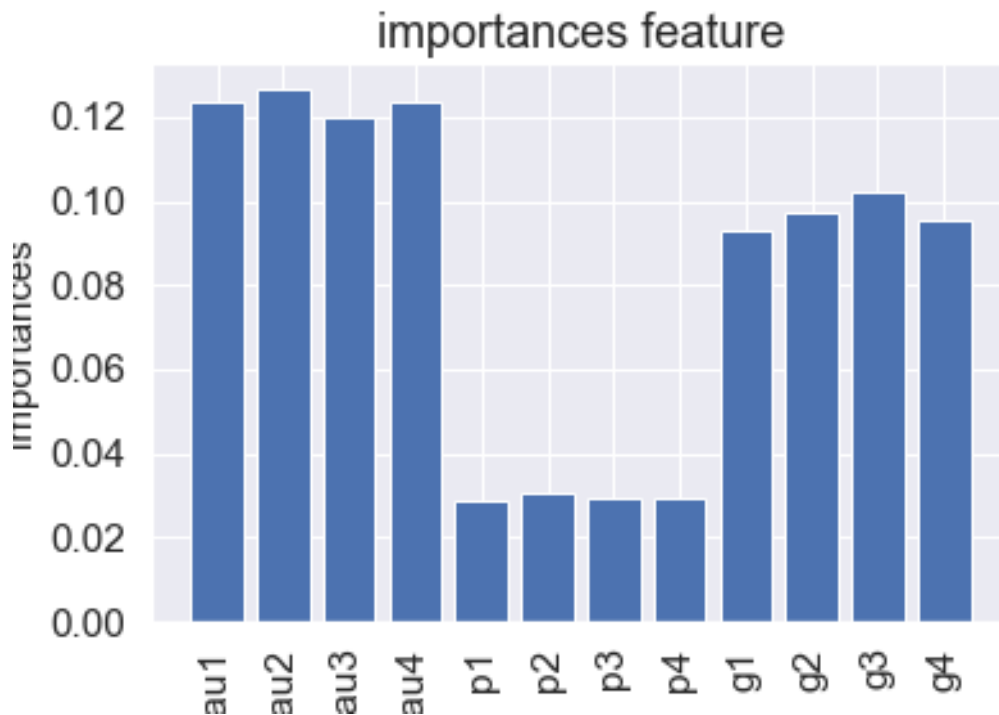
為了評量模型預測能力，我使用 Receiver operating characteristic(ROC)並計算 AUC，發現該隨機森林分類能力不錯，曲線下面積高達 0.98。



並把混淆矩陣(confusion matrix)可視化，發現雖然分類穩定的效果稍弱，但準確性也高達 97%

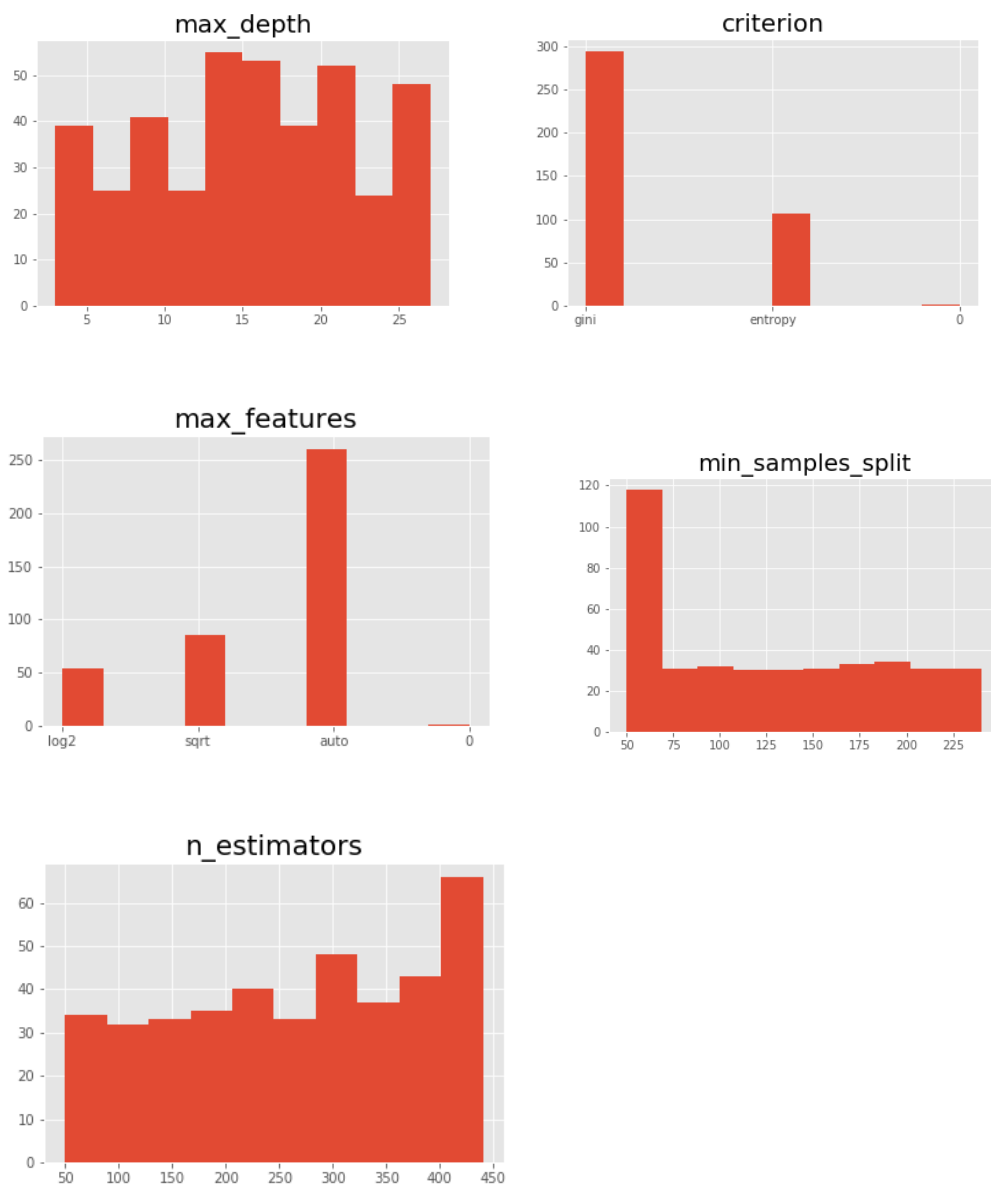


其中特過此隨機森林的模型也可以得知何種參數對於分類的影響性，透過分析葉子與根的熵(entropy)，即可進一步對數據的解讀，可以觀察到 tau(reaction time of participant)和 g(coefficient proportional to price elasticity)對於分類電網是否穩定很重要。



為了更加挖掘隨機森林潛力，我透過 Bayesian Hyperparam Optimization 自動調整超參數，包含了 criterio、n_estimators、max_depth 等，其透過高斯過程自動找尋最佳超參數，畢竟由於超參數眾多無法使用棋盤法一一尋找。

以下是我疊代 400 次後呈現的超參數分佈：



整理出最好的超參數：

criterion: gini

max_depth: 13

max_features: auto

min_samples_split: 50

n_estimators: 430

總結

透過隨機森林的算法進行分類，並計算出分類結果的**精準度**、**召回率** (recall)、**f1-score** 可以得到以下總表，可以顯示該分類方法十分優異，對於 class0(unstable)和 class1(stable)的分類都有著**極高的正確率**。

	precision	recall	f1-score	support
0.0	0.98	0.99	0.99	6380
1.0	0.98	0.97	0.98	3620
micro avg	0.98	0.98	0.98	10000
macro avg	0.98	0.98	0.98	10000
weighted avg	0.98	0.98	0.98	10000

以下是該隨機森林下的其中一個**分支的可視化**

