

FinTech HW2

tags: FinTech

年級：資工碩二

學號：R08922125

姓名：張皓鈞

INSTRUCTIONS

Dataset from kaggle Credit Card Fraud Detection. The features V1, V2, ... V28 are the principal components obtained with PCA, the only features which have not been transformed with PCA are 'Time' and 'Amount'. Feature 'Time' contains the seconds elapsed between each transaction and the first transaction in the dataset. Class value 1 in case of fraud and 0 otherwise.

Data Preprocessing

- train data=80%, test data=20%
- amount 屬性的範圍太大，需要做feature scale的調整所以將此做normalize

```
1  import pandas as pd
2  import numpy as np
3  import matplotlib.pyplot as plt
4  from sklearn.preprocessing import StandardScaler
5  from sklearn.model_selection import train_test_split
6  import keras
7
8  ## Load data and preprocessing
9  data = pd.read_csv('Data.csv')
10 scaler = StandardScaler()
11 data['NormalizedAmount'] = scaler.fit_transform(data['Amount'].values.reshape(-1))
12 data = data.drop(['Amount', 'Time'], axis = 1)
13 y = data['Class']
14 x = data.drop(['Class'], axis = 1)
15
16 from sklearn.model_selection import train_test_split
17 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 42)
18
```

Problem 1. Classification

(i) Construct DNN model and Gridsearch

- 用keras sequential建立四層的DNN model, 總共有 $V_1 \dots V_{28}$, *normalizedAmount* 29個feature去做training。
- 前三層的activation為ReLU, 最後一層output為sigmoid
- 中間dropout 50% · 以防止overfitting
- $\text{learn_rate} = [0.0001, 0.0003, 0.0005, 0.001]$
 $\text{batch_size} = [10, 60]$
 $\text{epochs} = [20, 70]$
 $\text{neurons} = [5, 10]$
做gridsearch找出32種參數組合中 · accuracy最高的一組參數。

```

1  from keras.models import Sequential
2  from keras.layers import Dense
3  from keras.layers import Dropout
4  from sklearn.model_selection import GridSearchCV
5  from keras.wrappers.scikit_learn import KerasClassifier
6
7  def create_model(learn_rate=0.001, neurons=10):
8      # create model
9      model = Sequential()
10     model.add(Dense(neurons, input_dim = 29, activation = 'relu'))
11     # model.add(Dense(24, activation = 'relu'))
12     model.add(Dense(neurons, activation = 'relu'))
13     Dropout(0.5),
14     model.add(Dense(neurons, activation = 'relu'))
15     model.add(Dense(1, activation = 'sigmoid'))
16
17     opt = keras.optimizers.Adam(learning_rate=learn_rate)
18     # Compile model
19     model.compile(loss='binary_crossentropy', optimizer=opt, metrics=['accuracy'])
20     return model
21
22     # Gridsearch Params
23     learn_rate = [0.0001, 0.0003 ,0.0005, 0.001]
24     batch_size = [10, 60]
25     epochs = [20, 70]
26     neurons = [5, 10]
27     param_grid = dict(batch_size=batch_size, epochs=epochs, learn_rate=learn_rate, neu
28
29     model = KerasClassifier(build_fn=create_model)
30     grid = GridSearchCV(estimator=model, param_grid=param_grid, n_jobs=-1, cv=None)
31     grid_result = grid.fit(x_train, y_train)
32
33     # summarize results
34     print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
35     means = grid_result.cv_results_['mean_test_score']
36     stds = grid_result.cv_results_['std_test_score']
37     params = grid_result.cv_results_['params']
38     for mean, stdev, param in zip(means, stds, params):
39         print("%f (%f) with: %r" % (mean, stdev, param))

```

需要跑一點時間，已經特地降低grid數目、epochs跟neurons，準確度影響並不大。

- 得出最佳參數組合

Best: 0.952679

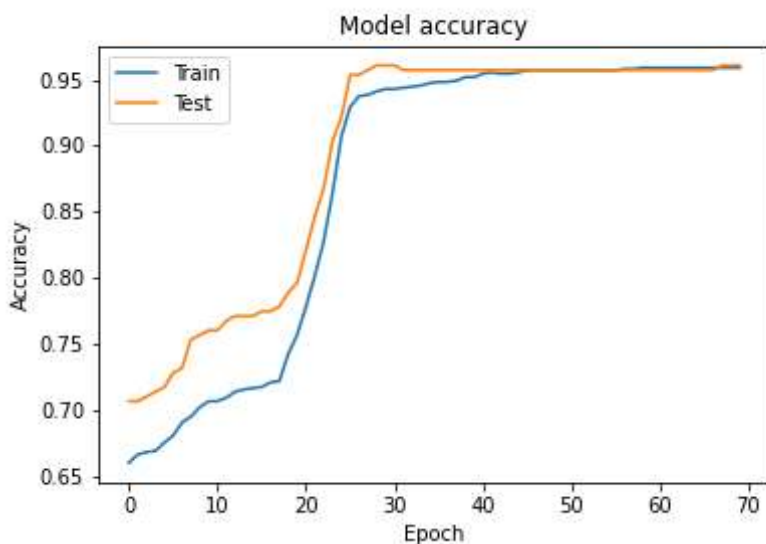
using {'batch_size': 10, 'epochs': 70, 'learn_rate': 0.0001, 'neurons': 10}

Choose best model from gridsearch and plot

```

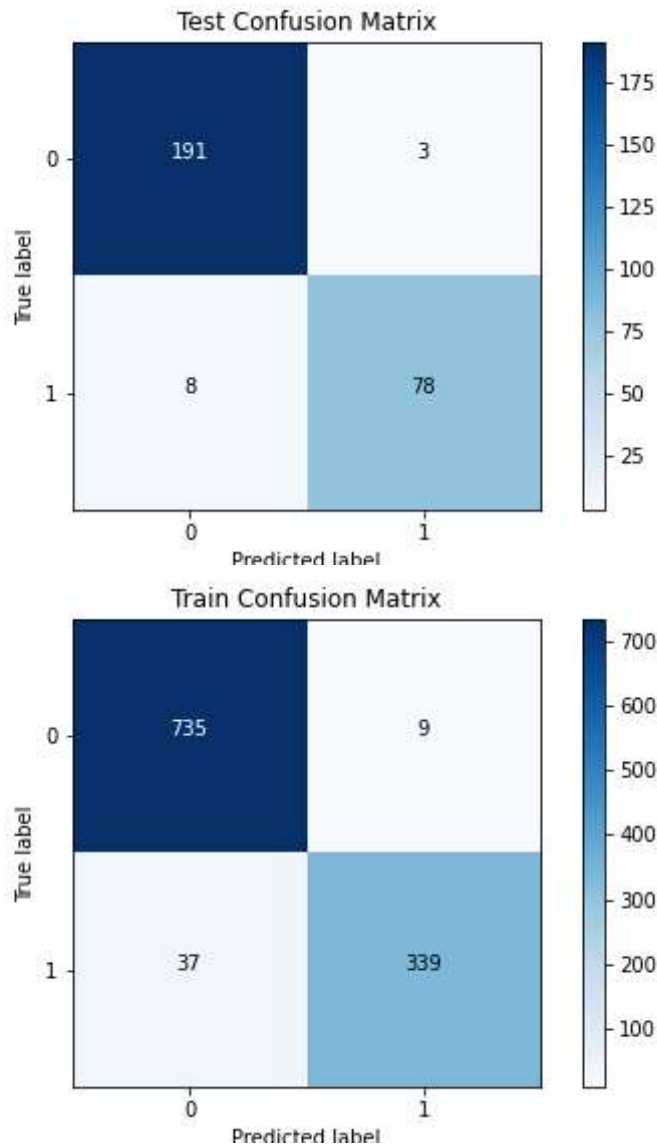
1  grid_best = grid_result.best_estimator_
2  history = grid_best.fit(x_train, y_train, validation_data=(x_test, y_test))
3
4  ## DNN Draw accuracy and loss
5  plt.plot(history.history['accuracy'])
6  plt.plot(history.history['val_accuracy'])
7  plt.title('Model accuracy')
8  plt.ylabel('Accuracy')
9  plt.xlabel('Epoch')
10 plt.legend(['Train', 'Test'], loc='upper left')
11 plt.show()
12
13 plt.plot(history.history['loss'])
14 plt.plot(history.history['val_loss'])
15 plt.title('Model loss')
16 plt.ylabel('loss')
17 plt.xlabel('Epoch')
18 plt.legend(['Train', 'Test'], loc='upper left')
19 plt.show()

```



(ii) Plot Confusion matrices

- 用sklearn的
 - from sklearn.metrics import confusion_matrix
 - from sklearn.metrics import plot_confusion_matrix



```

1 def confusion(model,x,y):
2     y_predict = model.predict(x)
3     return confusion_matrix(y,y_predict)
4
5 cm_train_DNN = confusion(grid_best,x_train, y_train)
6 cm_test_DNN = confusion(grid_best,x_test, y_test)
7
8 cm_plot_labels = ['0','1']
9
10 plot_confusion_matrix(cm_train_DNN, cm_plot_labels, 'Train Confusion Matrix')
11 plot_confusion_matrix(cm_test_DNN, cm_plot_labels, 'Test Confusion Matrix')

```

(iii) Precision, Recall, F1-Score

- 雖然keras裡有現成的函式，但我想說上題都有Confusion Matrix了就帶入公式自己算。
- 題目是說對each class都要算上述的指標，所以在confusion matrix中，class 1與 class 0的 true positive等等數值會完全相反。

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive}$$

$$Recall = \frac{TruePositive}{TruePositive + FalseNegative}$$

$$F1 - score = 2 * \frac{Precision * Recall}{Precision + Recall}$$

```

1  def calculation(cm, class_ = 1):
2      TP = cm[1][1]
3      FP = cm[0][1]
4      FN = cm[1][0]
5      TN = cm[0][0]
6      if (class_ == 0):
7          TP = TN
8          tmp = FN
9          FN = FP
10         FP = tmp
11
12         precision = TP / (TP+FP)
13         recall = TP / (TP+FN)
14         F1_score = 2 * (precision*recall) / (precision+recall)
15         print("For Class "+ str(class_))
16         print("precision: ",precision)
17         print("recall: ", recall)
18         print("F1-score: ",F1_score)
19 #####
20 output:
21 -----DNN model training-----
22
23 For Class 1
24 precision:  0.9741379310344828
25 recall:    0.901595744680851
26 F1-score:  0.93646408839779
27
28 For Class 0
29 precision:  0.9520725388601037
30 recall:    0.9879032258064516
31 F1-score:  0.9696569920844327
32 -----testing-----
33
34 For Class 1
35 precision:  0.9629629629629629
36 recall:    0.9069767441860465
37 F1-score:  0.934131736526946
38
39 For Class 0
40 precision:  0.9597989949748744
41 recall:    0.9845360824742269
42 F1-score:  0.9720101781170485

```

(iv) Difference between Decision Tree and Random Forest

Decision Tree 是以樹狀為基礎的演算法，透過歸納規則將資料從樹根開始分類，一節一節尋找最佳分割點來將資料分成為小單位的集合，中間有時也會透過園丁修剪，而成為一顆樹形美麗的決策樹。

不過當訓練資料集內的數目太少，而變數太多時，分類的效果會變差，另外，決策樹在分類上屬於固定的路徑，沒辦法像類神經在分類過程有容錯能力，所以會有overfitting的狀況產生。

Random Forest是以隨機(重新抽樣)的方法種植出許多決策樹，樹的集合就是森林，接著從決策樹們的投票結果中選出票數最多的候選人作為本屆選舉結果。由於是基於隨機抽樣及多顆樹，效果會比決策樹好許多，但相對運算的時間也比較長。

(v) Decision Tree model

```

1  from sklearn.tree import DecisionTreeClassifier
2  tree = DecisionTreeClassifier(max_depth=3,min_samples_leaf=5,criterion='entropy')
3  tree.fit(x_train,y_train)
4
5  cm_train_tree = confusion(tree,x_train, y_train)
6  cm_test_tree = confusion(tree,x_test, y_test)
7
8  ## Precision, recall, F1-score from decision tree
9  print("-----Decision Tree Training-----")
10 print("acc_train=",tree.score(x_train,y_train))
11 calculation(cm_train_tree)
12 print("-----Test-----")
13 print("acc_test=",tree.score(x_test,y_test))
14 calculation(cm_test_tree)
15
16 #####
17 output:
18 -----Decision Tree traing-----
19
20 For Class 1
21 precision:  0.9938461538461538
22 recall:    0.8590425531914894
23 F1-score:  0.9215406562054208
24
25 For Class 0
26 precision:  0.9333333333333333
27 recall:    0.9973118279569892
28 F1-score:  0.9642625081221572
29 acc_train= 0.9508928571428571
30 -----Test-----
31
32 For Class 1
33 precision:  1.0
34 recall:    0.8837209302325582
35 F1-score:  0.9382716049382717
36
37 For Class 0
38 precision:  0.9509803921568627
39 recall:    1.0
40 F1-score:  0.9748743718592965
41 acc_test= 0.9642857142857143

```

(v) Random Forest model

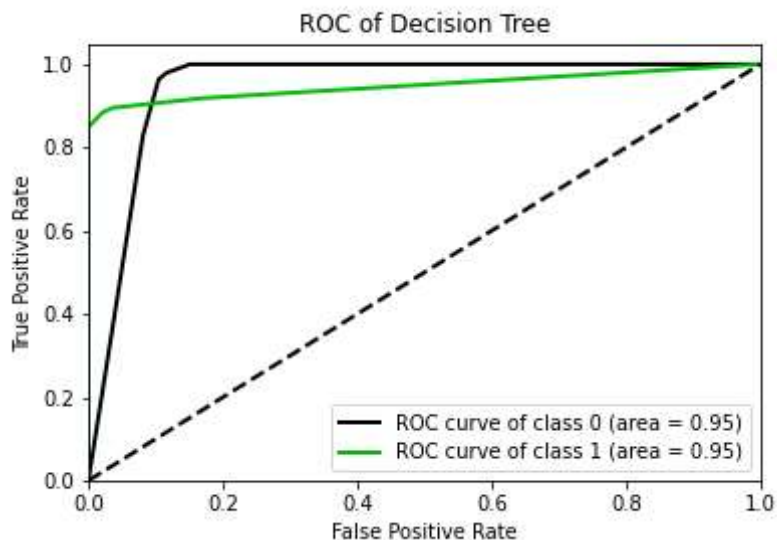
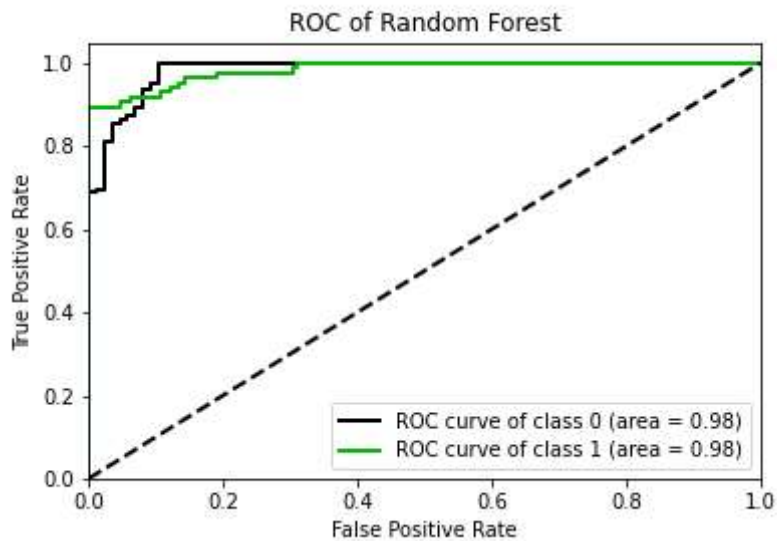
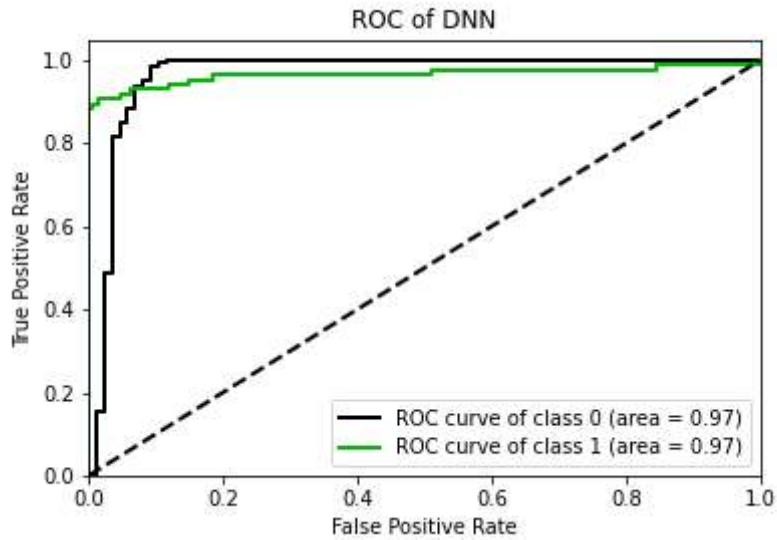
```

1  ## Random Forest model
2  from sklearn.ensemble import RandomForestClassifier
3  forest = RandomForestClassifier(n_estimators=100, max_depth=4)
4  forest.fit(x_train, y_train)
5
6  cm_train_forest = confusion(forest, x_train, y_train)
7  cm_test_forest = confusion(forest, x_test, y_test)
8  print("-----Random Forest Training-----")
9  calculation(cm_train_forest, class_=1)
10 calculation(cm_train_forest, class_=0)
11 print("acc_train=", forest.score(x_train, y_train))
12
13 print("-----Test-----")
14 calculation(cm_test_forest, class_=1)
15 calculation(cm_test_forest, class_=0)
16 print("acc_test=", forest.score(x_test, y_test))
17 #####
18 output:
19 -----Random Forest Training-----
20
21 For Class 1
22 precision:  0.9939577039274925
23 recall:    0.875
24 F1-score:  0.9306930693069307
25
26 For Class 0
27 precision:  0.9404309252217997
28 recall:    0.9973118279569892
29 F1-score:  0.9680365296803654
30 acc_train= 0.95625
31 -----Test-----
32
33 For Class 1
34 precision:  1.0
35 recall:    0.8837209302325582
36 F1-score:  0.9382716049382717
37
38 For Class 0
39 precision:  0.9509803921568627
40 recall:    1.0
41 F1-score:  0.9748743718592965
42 acc_test= 0.9642857142857143
43

```

(vi) ROC from DNN, Decision Tree, Random Forest

```
1  ## ROC Curve from DNN, Decision Tree and Random Forest
2  !pip install scikit-plot
3  import scikitplot as skplt
4
5  def roc_curve(model, x_test, y_test, model_name = "model_name"):
6      y_probs = model.predict_proba(x_test)
7      skplt.metrics.plot_roc(y_test, y_probs, plot_micro=False, plot_macro=False)
8      plt.title("ROC of "+ model_name)
9      plt.savefig("ROC of "+ model_name)
10     plt.show(block=False)
11
12     #Plot
13     roc_curve(grid_best, x_test, y_test, model_name = "DNN")
14     roc_curve(tree, x_test, y_test, model_name = "Decision Tree")
15     roc_curve(forest, x_test, y_test, model_name = "Random Forest")
```



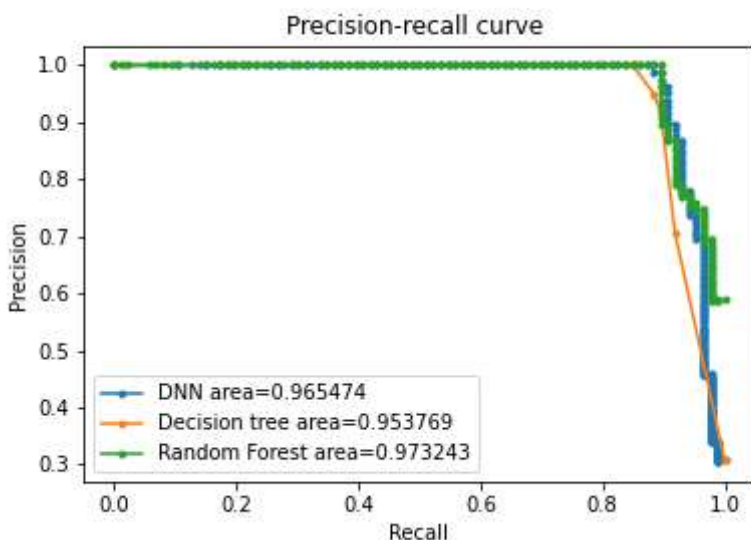
可以看出Random forest的效果最好(可能也跟我的參數有關) · 而DNN model在調參過後效果也不錯。

(vi) Precision-recall curve from DNN, Decision Tree, Random Forest

```

1  ## Precision-Recall curve from DNN, Decision Tree and Random Forest
2  from sklearn.metrics import precision_recall_curve
3  from sklearn.metrics import auc
4
5  def plot_precisionrecall_curve(model,x_test,y_test, model_name="model"):
6      y_probs = model.predict_proba(x_test)
7      # y_probs_tree = tree.predict_proba(x_test)
8      # y_probs_forest = forest.predict_proba(x_test)
9
10     precision, recall, _ = precision_recall_curve(y_test, y_probs[:,1])
11     # precision_tree, recall_tree, _ = precision_recall_curve(y_test, y_probs_tree)
12     # precision_forest, recall_forest, _ = precision_recall_curve(y_test, y_probs_
13     plt.plot(recall, precision, marker='.', label=model_name + ' area=%f'%auc(reca
14
15     # Plot
16     plot_precisionrecall_curve(grid_best,x_test,y_test,model_name="DNN")
17     plot_precisionrecall_curve(tree,x_test,y_test,model_name="Decision tree")
18     plot_precisionrecall_curve(forest,x_test,y_test,model_name="Random Forest")
19     plt.title("Precision-recall curve ")
20     plt.xlabel('Recall')
21     plt.ylabel('Precision')
22     plt.legend()
23     plt.savefig("Precision-recall curve")
24     plt.show()

```



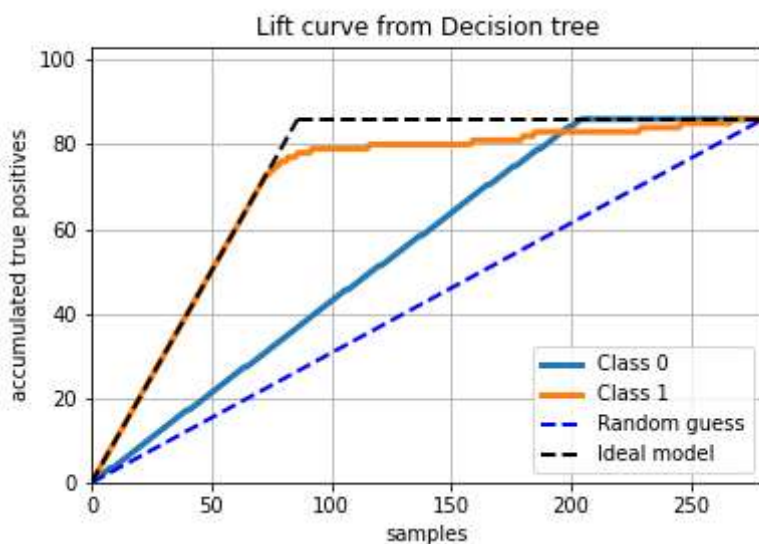
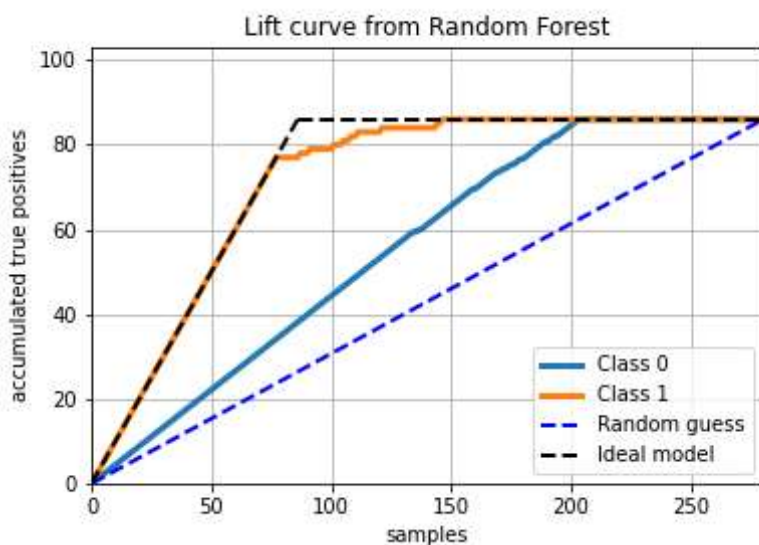
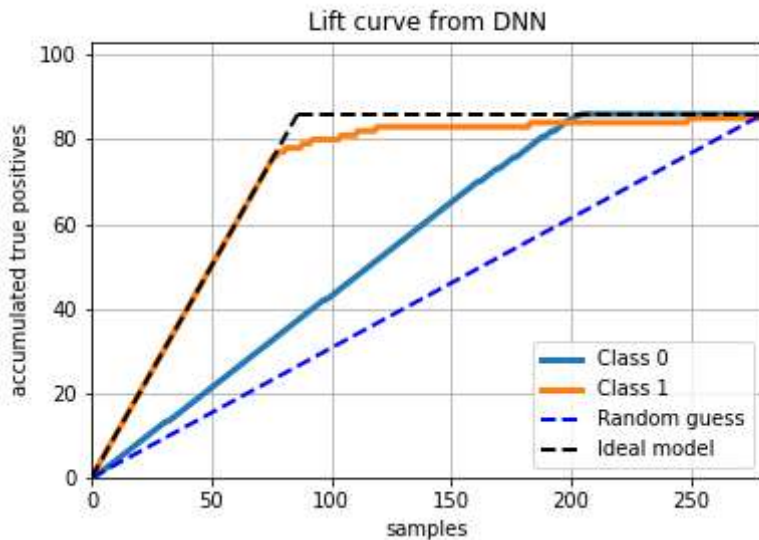
Problem 2.(a) Lift curve

根據助教給的hw2.pdf，其中對於lift curve的定義比較像是不同scale之下的cumulative gain curve，於是我修改原本scikitplot的plot_cumulative_gain使其符合作業中的敘述。

```

1  ## Lift curve (modified by cumulative gain curve)
2  from scikitplot.helpers import cumulative_gain_curve
3  def plot_cumulative_gain(model, x_test ,y_true, model_name="model_name",
4                          ax=None, figsize=None, title_fontsize="large",
5                          text_fontsize="medium"):
6      y_true = np.array(y_true)
7      y_probas = model.predict_proba(x_test)
8      y_probas = np.array(y_probas)
9
10     classes = np.unique(y_true)
11     if len(classes) != 2:
12         raise ValueError('Cannot calculate Cumulative Gains for data with '
13                           '{} category/ies'.format(len(classes)))
14
15     # Compute Cumulative Gain Curves
16     total = len(y_true)
17     true_positive = y_true.sum()
18
19     percentages1, gains1 = cumulative_gain_curve(y_true, y_probas[:, 0],
20                                                  classes[0])
21     percentages2, gains2 = cumulative_gain_curve(y_true, y_probas[:, 1],
22                                                  classes[1])
23
24     # print(gains1)
25     if ax is None:
26         fig, ax = plt.subplots(1, 1, figsize=figsize)
27
28     ax.set_title("Lift curve from "+model_name, fontsize=title_fontsize)
29
30     ax.plot(percentages1*total, gains1*true_positive, lw=3, label='Class {}'.for
31 ax.plot(percentages2*total, gains2*true_positive, lw=3, label='Class {}'.for
32
33     ax.set_xlim([0.0, total])
34     ax.set_ylim([0.0, true_positive*1.2])
35
36     ax.plot([0, total], [0, true_positive], 'b--', lw=2, label='Random guess')
37     ax.plot([0, true_positive], [0, true_positive], 'k--', lw=2, label='Ideal mc
38 ax.plot([true_positive, total], [true_positive, true_positive], 'k--', lw=2)
39
40     ax.set_xlabel('samples', fontsize=text_fontsize)
41     ax.set_ylabel('accumulated true positives', fontsize=text_fontsize)
42     ax.tick_params(labelsize=text_fontsize)
43     ax.grid('on')
44     ax.legend(loc='lower right', fontsize=text_fontsize)
45     plt.savefig("Lift curve from "+model_name)
46
47     return ax
48
49 plot_cumulative_gain(grid_best,x_test,y_test, model_name='DNN')
50 plot_cumulative_gain(tree,x_test,y_test, model_name='Decision tree')
51 plot_cumulative_gain(forest,x_test,y_test, model_name='Random Forest')

```



- 可以看出 Class 1 是比 Class 0 具有代表性的，這也是因為此dataset就是要預測是否為Class 1: Fraud的關係。
- 這三張圖看起來 Random Forest與 DNN 的預測效果比 Decision Tree好，與上面的ROC及 Precision-recall都有一樣的結果。

