# Midterm

# Midterm Project Competition

N96094196 張維峻

工程科學系

# 0.摘要

InClass Prediction Competition

## Midterm Project Competition
This competition is a classification problem about material properties.

78 teams · 6 hours to go

Overview    Data    Notebooks    Discussion    Leaderboard    Rules    Team            My Submissions    **Submit Predictions**

### Overview

**Description**

Evaluation

This is the mid-term competition of the course of machine learning in engineering science in NCKU ES. This competition is a classification problem. Here we are going to build a machine learning model to classify material properties.

In this competition, we aim to design high-performance composites by arranging soft and stiff materials to fine-tune the mechanical properties of the entire composites and strengthen the ability of the material to resist crack propagation. In the data, there will be four categories there denote the superiority of the material properties. Your Machine Learning model should be able to distinguish them from each other and make predictions among the private testing dataset. The baseline of this competition is 60% accuracy, which means your model should be able to predict the results above 70% accuracy. The higher the accuracy the more scores you earn. Live long and happy coding!

### Overview

Description

**Evaluation**

#### Score

The evaluation metric for this competition is categorization accuracy. The categorization accuracy, commonly used in classification problem, measures accuracy using the number of correct predictions and total number of predictions. The categorization accuracy is given by:

$$\text{Score} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

#### Submission Format

Submission files should contain two columns: index and Label. The first column must be index, the second column is the predicted result. Each row is index and predicted label, please separate with commas. When uploading, please replace the original string label with a number from 1 to 4.

| index | Label |
|-------|-------|
| 0 | 2 |
| 1 | 1 |
| 2 | 4 |
| ⋮ | ⋮ |

| Original label | | Submitted label |
|---|---|---|
| Excellent | → | 1 |
| Good | → | 2 |
| Fair | → | 3 |
| bad | → | 4 |

# 1.選擇方法-SVM

選擇使用 SVM 來訓練分類器,讀取測試資料,把所以特徵都

丟進去訓練

```
[1]: import numpy as np
     import pandas as pd

[2]: # open datafile
     val_file = pd.read_csv('val.csv')
     train_file = pd.read_csv('train.csv')
     test_file = pd.read_csv('test.csv')

[3]: label=['bad', 'fair', 'good', 'Excellent']

[9]: train_data = train_file.iloc[:,1:-1]
     train_data = np.nan_to_num(train_data, 0)
     #print(train_data)
     train_data.shape

[9]: (16000, 70)

[10]: train_data = np.nan_to_num(train_data, 0)
      #print(train_data)
      train_label = train_file[['Label']].to_numpy()
      #print(train_label)
      train_label_num = [label.index(l) for l in train_label]
      #print(train_label_num)

[11]: val_data = val_file.iloc[:,1:-1]

[12]: val_data = np.nan_to_num(val_data,0)
      # print(val_data)
      val_label = val_file[['Label']].to_numpy()
      #print(val_label)
      val_label_num = [label.index(l) for l in val_label]
      #print(val_label_num)

[13]: from sklearn.svm import SVC              #1. choose "Support vector claaifier"
      model = SVC(kernel='rbf', gamma=0.012, C=150)  # 2. instantiate model
      model.fit(train_data, train_label.ravel())              # 3. fit model to data
      result = model.predict(val_data)               # 4. predict on new data

      from sklearn.metrics import accuracy_score
      score = accuracy_score(val_label, result)
      print(score)

      0.45806451612903226
```

模型準確度:0.45806

# 2.選擇特徵參數

先去除訓練集中,後面一些材料特性參數,以前面 8*8=64 個點構

成的圖去訓練模型

```python
1  import  numpy  as  np
2  import  pandas  as  pd
3  # open  datafile
4  val_file  =  pd.read_csv('val.csv')
5  train_file  =  pd.read_csv('train.csv')
6  test_file  =  pd.read_csv('test.csv')
7  label=['bad',  'fair',  'good',  'Excellent']
8  train_data  =  train_file.iloc[:,1:65]
9  train_data  =  np.nan_to_num(train_data,  0)
10 train_data.shape
11 train_data  =  np.nan_to_num(train_data,  0)
12 train_label  =  train_file[['Label']].to_numpy()
13 train_label_num  =  [label.index(l)  for  l  in  train_label]
14 val_data  =  val_file.iloc[:,1:65]
15 val_data  =  np.nan_to_num(val_data,0)
16 #  print(val_data)
17 val_label  =  val_file[['Label']].to_numpy()
18 #print(val_label)
19 val_label_num  =  [label.index(l)  for  l  in  val_label]
20 #print(val_label_num)
21 from  sklearn.svm  import  SVC         #1.  choose  "Support  vector  claaifier"
22 model  =  SVC(kernel='rbf',  gamma=0.1,  C=10)  # 2.  instantiate  model
23 model.fit(train_data,  train_label.ravel())       # 3.  fit  model  to  data
24 result  =  model.predict(val_data)                 # 4.  predict  on  new  data
25
26 from  sklearn.metrics  import  accuracy_score
27 score  =  accuracy_score(val_label,  result)
28 print(score)
```

0.755

模型準確度:0.755

# 3.SVC 一些調參

使用 GridSearchCV 和 Validation Curve 找到模型最佳的

gamma 和 C

```
Best parameters set found on development set:

{'C': 150, 'gamma': 0.012, 'kernel': 'rbf'}

Grid scores on development set:

0.484 (+/-0.002) for {'C': 0.01, 'gamma': 0.012, 'kernel': 'rbf'}
0.633 (+/-0.011) for {'C': 1, 'gamma': 0.012, 'kernel': 'rbf'}
0.742 (+/-0.017) for {'C': 150, 'gamma': 0.012, 'kernel': 'rbf'}

Detailed classification report:

The model is trained on the full development set.
The scores are computed on the full evaluation set.

              precision    recall  f1-score   support

   Excellent       1.00      1.00      1.00      1086
         bad       1.00      0.94      0.97       217
        fair       0.97      0.99      0.98       787
        good       0.99      0.99      0.99      1110

    accuracy                           0.99      3200
   macro avg       0.99      0.98      0.98      3200
weighted avg       0.99      0.99      0.99      3200
```
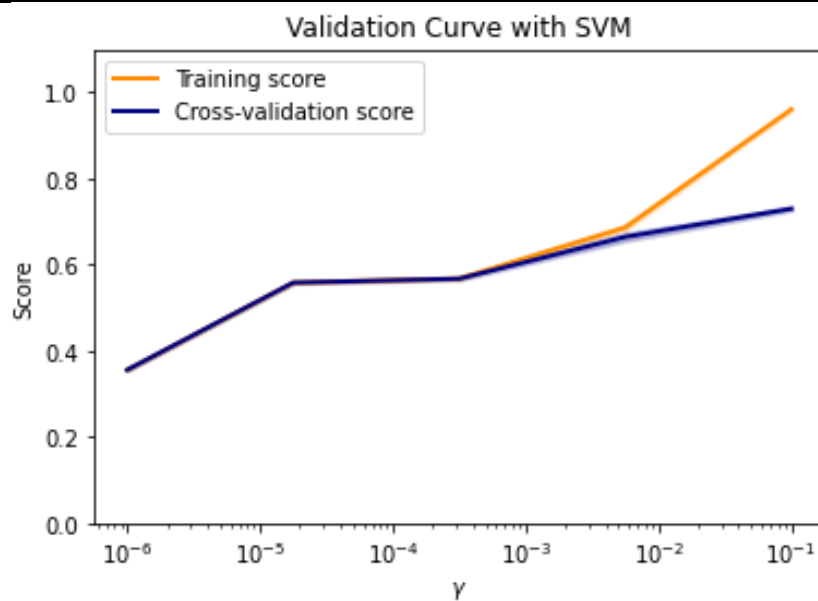
header_navigationMachine Learning Midterm
_N96094196_張維峻                                                                    6

```python
1 import numpy as np
2 import pandas as pd
3 # open datafile
4 val_file = pd.read_csv('val.csv')
5 train_file = pd.read_csv('train.csv')
6 test_file = pd.read_csv('test.csv')
7 label=['bad', 'fair', 'good', 'Excellent']
8 train_data = train_file.iloc[:,1:65]
9 train_data = np.nan_to_num(train_data, 0)
10 train_data.shape
11 train_data = np.nan_to_num(train_data, 0)
12 train_label = train_file[['Label']].to_numpy()
13 train_label_num = [label.index(l) for l in train_label]
14 val_data = val_file.iloc[:,1:65]
15 val_data = np.nan_to_num(val_data,0)
16 # print(val_data)
17 val_label = val_file[['Label']].to_numpy()
18 #print(val_label)
19 val_label_num = [label.index(l) for l in val_label]
20 #print(val_label_num)
21 from sklearn.svm import SVC          #1. choose "Support vector claaifier"
22 model = SVC(kernel='rbf', gamma=0.0123, C=150) # 2. instantiate model
23 model.fit(train_data, train_label.ravel())      # 3. fit model to data
24 result = model.predict(val_data)                # 4. predict on new data
25
26 from sklearn.metrics import accuracy_score
27 score = accuracy_score(val_label, result)
28 print(score)
```
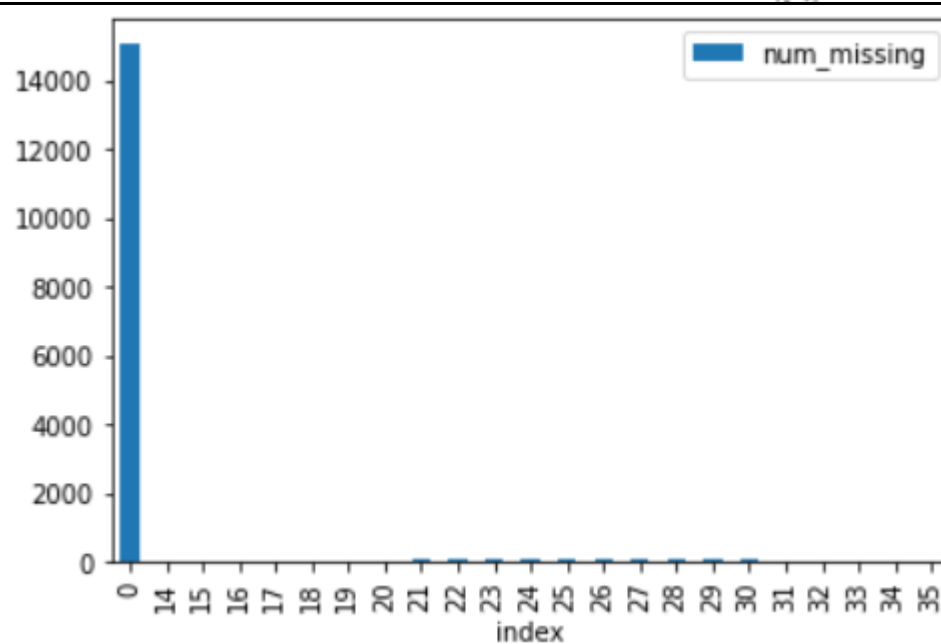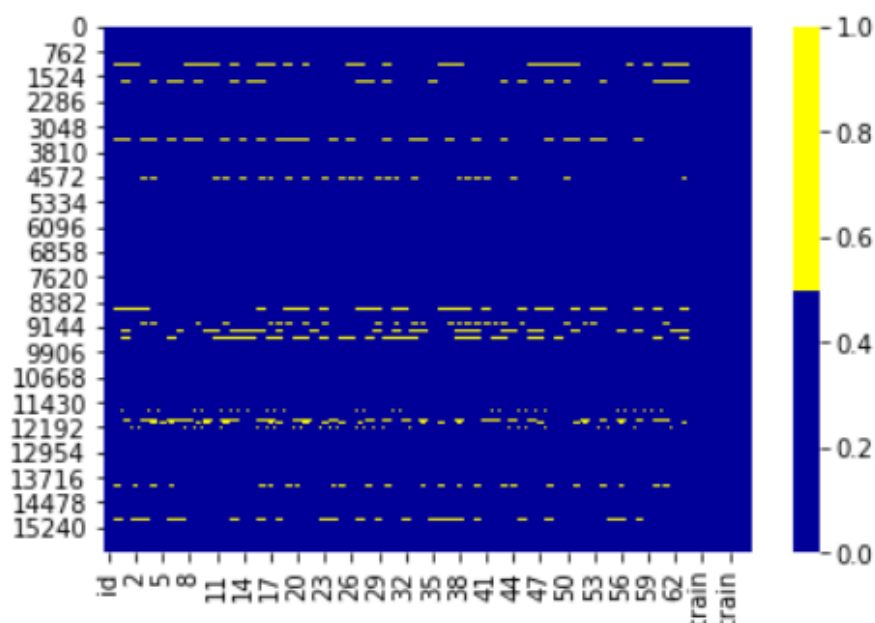
0.8479032258064516

模型準確度:0.84790

# 4.Data Clean- Missing data

發現訓練集中某些欄位有缺值,因此刪除沒有值的特定項

熱圖,柱狀圖顯示出哪些欄位有缺值

```
1 #data  clean  ->Missing  data
2 import  seaborn  as  sns
3
4 cols  =  train_file.columns[:]
5 # specify  the  colours  -  yellow  is  missing.  blue  is  not  missing.
6 colours  =  ['#000099',  '#ffff00']
7 sns.heatmap(train_file[cols].isnull(),  cmap=sns.color_palette(colours))
```

用熱點圖表示 有缺值的欄位

```
1 # first  create  missing  indicator  for  features  with  missing  data
2 for  col  in  train_file.columns:
3       missing  =  train_file[col].isnull()
4       num_missing  =  np.sum(missing)
5
6       if  num_missing  >  0:
7             print('created  missing  indicator  for:  {}'.format(col))
8             train_file['{}_ismissing'.format(col)]  =  missing
9
10
11 # then  based  on  the  indicator,  plot  the  histogram  of  missing  values
12 ismissing_cols = [col  for  col  in  train_file.columns  if  'ismissing'  in  col]
13 train_file['num_missing']  =  train_file[ismissing_cols].sum(axis=1)
14
15 train_file['num_missing'].value_counts().reset_index().sort_values(by='index').plot.bar(x='index',  y='num_missing'
16
```

用柱狀圖表示 並且標出有缺值的欄位

```
1 print('train_data  shape  =  ',train_data.shape)
2 print('train_clndata  shape  ='  ,train_clndata.shape)
3
4 train_clnlabel  =  df_less_missing_rows[['Label']].to_numpy()
5 train_clnlabel_num  =  [label.index(1)  for  1  in  train_clnlabel]
6
```

```
train_data shape =  (16000, 64)
train_clndata shape = (15077, 63)
```

移除有缺值得欄位,對比原先訓練集 ,新的訓練及欄位減少一些

```
1 from  sklearn.svm  import  SVC          #1.  choose  "Support  vector  claaifier"
2 model  =  SVC(kernel='rbf',  gamma=0.0123,  C=150)   # 2.  instantiate  model
3 #model.fit(train_data,  train_label.ravel())       # 3.  fit  model  to  data
4 model.fit(train_clndata,train_clnlabel.ravel())
5 result  =  model.predict(val_data)                  # 4.  predict  on  new  data
6
7 from  sklearn.metrics  import  accuracy_score
8 score  =  accuracy_score(val_label,  result)
9 print(score)
```
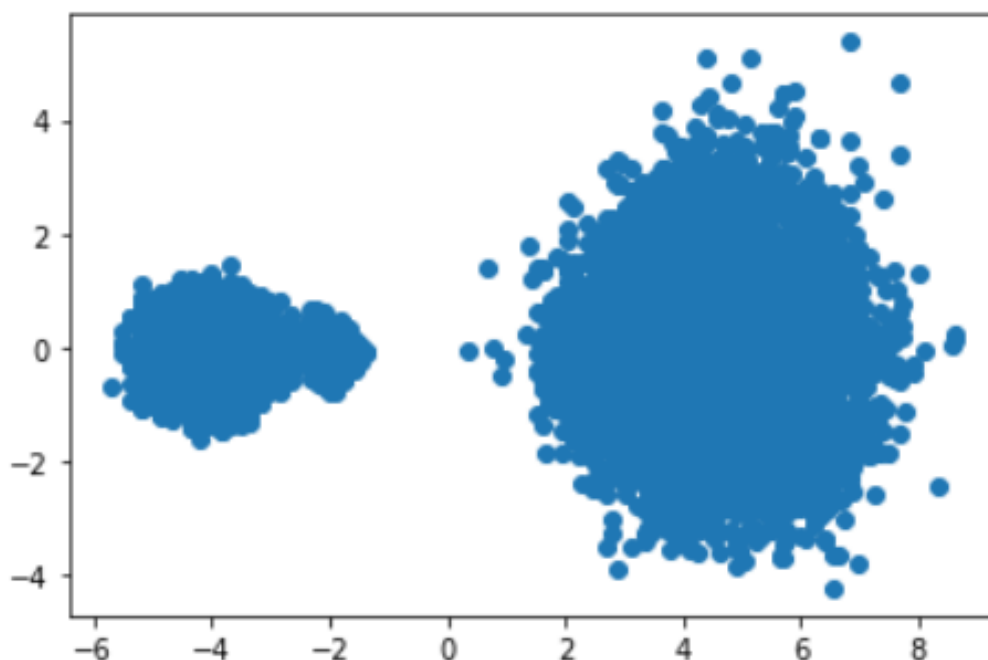
```
0.8608064516129033
```

模型準確度:0.86080

# 其他嘗試

使用 PCA 把 8*8 的圖形降低維度,在 2D 空間中表現

```
1 #64維度PCA降維成2D
2 from sklearn.decomposition import PCA
3 import matplotlib.pyplot as plt
4
5 n_components = 2
6 random_state = 6666
7
8 pca = PCA(n_components=n_components,
9                 random_state=random_state)
10 L = pca.fit_transform(Z)
11
12 plt.scatter(L[:,0],L[:,1])
13 # plt.axis('1');
```



過濾掉一些離群值,或許預測準度會上升