

1 Cuprate superconducting materials above liquid nitrogen temperature from machine learning

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population

 - 1.2.1 The descriptor
 - 1.2.2 The feature space

- 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high-throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

1.1 Exploratory analysis of data

Contents ⚙

1.1 Exploratory analy:
▼ 1.2 Descriptors popul
1.2.1 The descriptoro
1.2.2 The feature s
▼ 1.3 Classification mod
1.3.1 Use pycaret to
1.3.2 Select the app
1.3.3 RFC
1.3.4 KNN
▼ 1.3.5 DT
1.3.5.1 Visualizat
1.3.6 ADC
▼ 1.3.7 RFC : Analysis
1.3.7.1 F1 score
1.3.7.2 confusion
1.3.7.3 Accuracy
1.3.7.4 The thres
1.3.7.5 Cost sens
▼ 1.4 Train the deep ne
1.4.1 Use S1, S2 de
1.4.2 Parameter op
1.4.3 The trained m
1.4.4 Residual anal
▼ 1.5 Virtual high throu
1.5.1 Virtual high-th
1.5.2 DNN predictio
▼ 1.6 Deep neural netw
1.6.1 PCA
1.6.2 T-SNE

In [2]:

```
1 import os
2 import sys
3 import time
4 import matplotlib
5 import pathlib
6 import graphviz
7 import itertools
8 import pycaret
9 import pyforest
10 import pandas as pd
11 import numpy as np
12 import matplotlib.pyplot as plt
13 from mpl_toolkits import mplot3d
14 import seaborn as sns
15 #from sklearn import *
16 import tensorflow as tf
17 from tensorflow import keras
18 from tensorflow.keras import layers
19 from sklearn.model_selection import train_test_split
20 from sklearn.preprocessing import StandardScaler
21 from sklearn.preprocessing import MinMaxScaler
22 from sklearn.metrics import f1_score
23 from sklearn.metrics import classification_report
24 from sklearn.metrics import confusion_matrix
25 from sklearn.metrics import precision_recall_curve
26 from sklearn.metrics import roc_auc_score
27 from sklearn.metrics import roc_curve
28 from sklearn.model_selection import ShuffleSplit
29 from sklearn.model_selection import cross_val_score
30 from sklearn.ensemble import AdaBoostClassifier
31 from sklearn.ensemble import RandomForestClassifier
32 from sklearn.neighbors import KNeighborsClassifier
33 font2 = {'family' : 'Times New Roman',
34          'weight' : 'normal',
35          'size' : 20,
36          }
37 def plot_roc_curve(fprs, tprs):
38     plt.figure(figsize=(8, 6), dpi=80)
39     plt.plot(fprs, tprs)
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
42 plt.plot([0, 1], linestyle='--')
43 plt.xticks(fontsize=13)
44 plt.yticks(fontsize=13)
45 font2 = {'family' : 'Times New Roman',
46 'weight' : 'normal',
47 'size' : 20,
48     }
49 plt.ylabel('TP rate', font2)
50 plt.xlabel('FP rate', font2)
51 plt.title('ROC', font2)
52 plt.savefig('ROC.jpg', dpi=300)
53 plt.show()
54
55 def plot_cnf_matrix(cnf_matrix, description):
56     class_names = [0, 1]
57     fig = plt.gcf()
58     fig.set_size_inches(15.5, 10.5)
59     matplotlib.rcParams['font.sans-serif']=['SimHei'] # 用黑体显示中文
60     matplotlib.rcParams['axes.unicode_minus']=False
61     fig, ax = plt.subplots()
62     tick_marks = np.arange(len(class_names))
63     plt.xticks(tick_marks, class_names)
64     plt.yticks(tick_marks, class_names)
65     sns.heatmap(pd.DataFrame(cnf_matrix), annot = True, cmap = 'OrRd', fmt = 'g')
66     ax.xaxis.set_label_position('top')
67     plt.tight_layout()
68     plt.title(description, y = 1.1, fontsize=16)
69     font2 = {'family' : 'Times New Roman',
70 'weight' : 'normal',
71 'size' : 20,
72     }
73     plt.ylabel('True0/1', font2)
74     plt.xlabel('Pred0/1', font2)
75     fig = plt.gcf()
76     fig.set_size_inches(5.5, 4.5)
77     plt.savefig('cnf_matrix.jpg', dpi=300)
78     plt.show()
79
80 def plot_confusion_matrix(cm,
81                         classes,
82                         title='Confusion matrix',
83                         cmap=plt.cm.Blues):
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
 - ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
 - ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
 - ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
84 plt.imshow(cm, interpolation='nearest', cmap=cmap)
85 plt.title(title)
86 plt.colorbar()
87 tick_marks = np.arange(len(classes))
88 plt.xticks(tick_marks, classes, rotation=0)
89 plt.yticks(tick_marks, classes)
90
91 thresh = cm.max() / 2.
92 for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
93     plt.text(j, i, cm[i, j],
94             horizontalalignment="center",
95             color="white" if cm[i, j] > thresh else "black")
96
97 plt.tight_layout()
98 plt.ylabel('True label')
99 plt.xlabel('Predicted label')
```

In [3]:

```
1 data = pd.read_csv("./data.csv")
2 S1 = pd.read_csv("./S1.csv")
3 S2 = pd.read_csv("./S2.csv")
4 df = S1
```

In [4]: 1 data.describe()

Out[4]:

Tc

count	12340.000000
mean	23.503293
std	28.565083
min	0.000500
25%	3.400000
50%	10.000000
75%	33.000000
max	143.000000

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 data
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [5]: 1 df

Out[5]:

	name	Tc	H	He	Li	Be	B	C	N	O	...	Mt	Ds	Rg	Cn	Nh
0	Zr74.1Ti3.9Ni22	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
1	Zr70.2Ti7.8Ni22	3.320	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
2	Zr66.3Ti11.7Ni22	3.560	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
3	Zr62.4Ti15.6Ni22	3.215	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
4	Zr55Cu30Al10Ni5	1.350	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
...
12335	Ag0.05Rh0.04Ti0.91	1.950	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
12336	Ag0.03Ti0.97	2.670	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
12337	Ag0.035Cd0.01Sn0.955	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
12338	Ag0.005Zn0.995	0.763	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0
12339	Ag0.002Al0.998	1.128	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	0

12340 rows × 120 columns

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature set
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appr
 - 1.3.3 RFC
 - 1.3.4 KNN
- ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
- ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sens
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter opti
 - 1.4.3 The trained m
 - 1.4.4 Residual anal
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-th
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

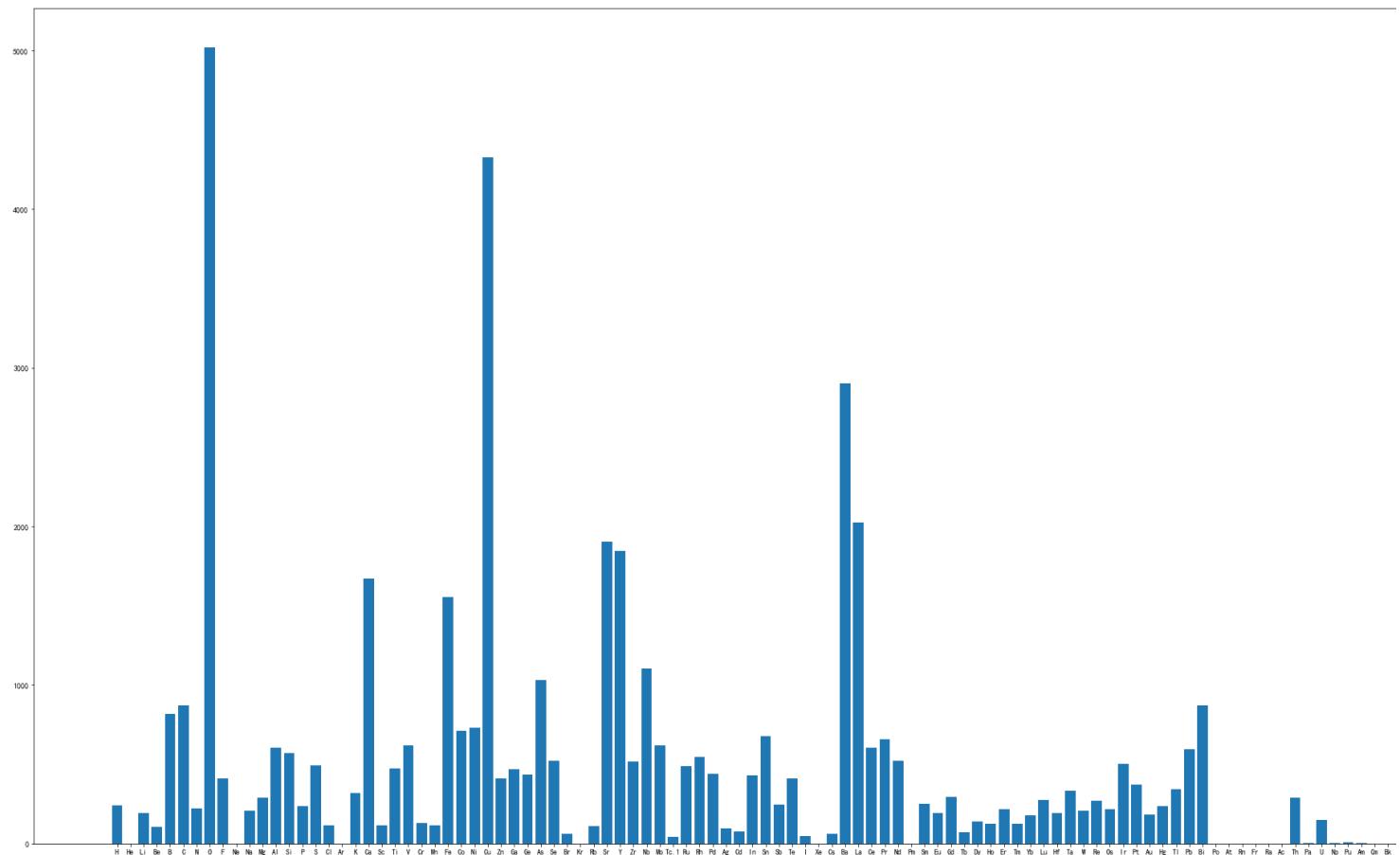
- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 data
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [8]:

```
1 distribution = []
2 elem = []
3 for i in df.columns[2:]:
4     ca = (df[i]==0).sum()
5     realele = 12340-ca
6     distribution.append(realele)
7     elem.append(i)
8 fengfu = pd.DataFrame(distribution)
9 elem = pd.DataFrame(elem)
10 elem["elem"] = elem
11 elem["distribution"] = distribution
12 x = elem["elem"]
13 y = elem["distribution"]
14 #seaborn.barplot(data=fengfu)
15 matplotlib.rcParams['font.sans-serif'] = ['SimHei']
16 matplotlib.rcParams['axes.unicode_minus'] = False
17 plt.figure(figsize=(40, 20))
18 plt.bar(x, y)
19 plt.savefig('distribution.jpg', dpi=300)
20 plt.show()
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

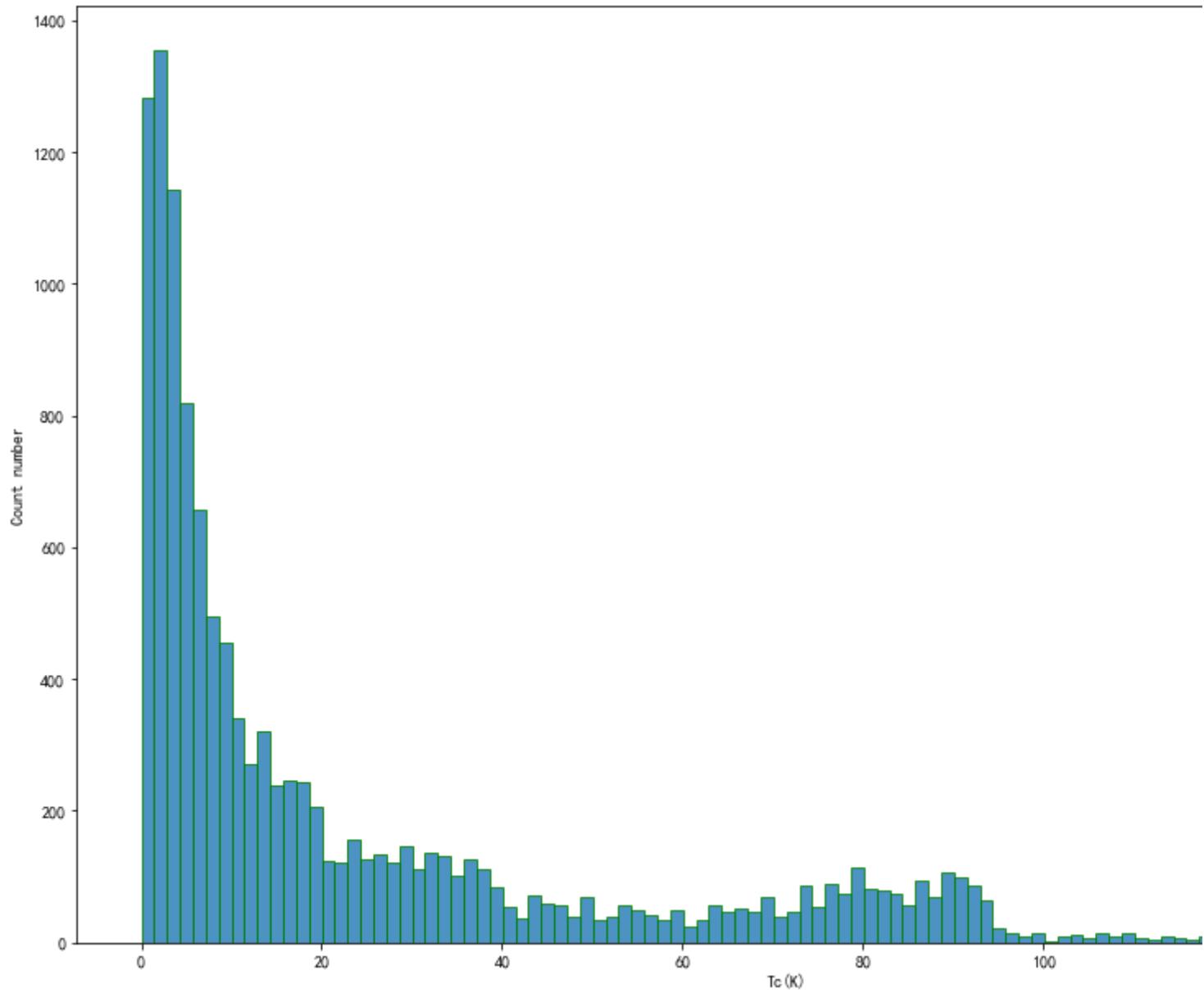


Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high-throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [9]:

```
1 #Superconducting transition temperature frequency diagram of all materials
2 fig = plt.gcf()
3 fig.set_size_inches(15.5, 10.5)
4 matplotlib.rcParams['font.sans-serif']=['SimHei']
5 matplotlib.rcParams['axes.unicode_minus']=False
6 plt.hist(data["Tc"], bins=100, color=None, edgecolor="green", alpha=0.8)
7 plt.title('Superconducting transition temperature frequency diagram')
8 plt.xlabel("Tc (K)")
9 plt.ylabel("Count number")
10 plt.show()
```



Contents ⚙️

- 1.1 Exploratory analysis
- 1.2 Descriptors population
- 1.2.1 The descriptor
- 1.2.2 The feature selection
- 1.3 Classification models
- 1.3.1 Use pycaret to
- 1.3.2 Select the appr
- 1.3.3 RFC
- 1.3.4 KNN
- 1.3.5 DT
- 1.3.5.1 Visualization
- 1.3.6 ADC
- 1.3.7 RFC : Analysis
- 1.3.7.1 F1 score
- 1.3.7.2 confusion
- 1.3.7.3 Accuracy
- 1.3.7.4 The thres
- 1.3.7.5 Cost sens
- 1.4 Train the deep net
- 1.4.1 Use S1, S2 de
- 1.4.2 Parameter opti
- 1.4.3 The trained m
- 1.4.4 Residual anal
- 1.5 Virtual high through
- 1.5.1 Virtual high-th
- 1.5.2 DNN prediction
- 1.6 Deep neural netw
- 1.6.1 PCA
- 1.6.2 T-SNE

Contents ⚙

1.1 Exploratory analysis
1.2 Descriptors population
1.2.1 The descriptor
1.2.2 The feature selection
1.3 Classification models
1.3.1 Use pycaret to
1.3.2 Select the appr
1.3.3 RFC
1.3.4 KNN
1.3.5 DT
1.3.5.1 Visualizat
1.3.6 ADC
1.3.7 RFC : Analysis
1.3.7.1 F1 score
1.3.7.2 confusion
1.3.7.3 Accuracy
1.3.7.4 The thresh
1.3.7.5 Cost sens
1.4 Train the deep neural network
1.4.1 Use S1, S2 de
1.4.2 Parameter optim
1.4.3 The trained model
1.4.4 Residual anal
1.5 Virtual high throughput
1.5.1 Virtual high-throu
1.5.2 DNN prediction
1.6 Deep neural network
1.6.1 PCA
1.6.2 T-SNE

In [10]:

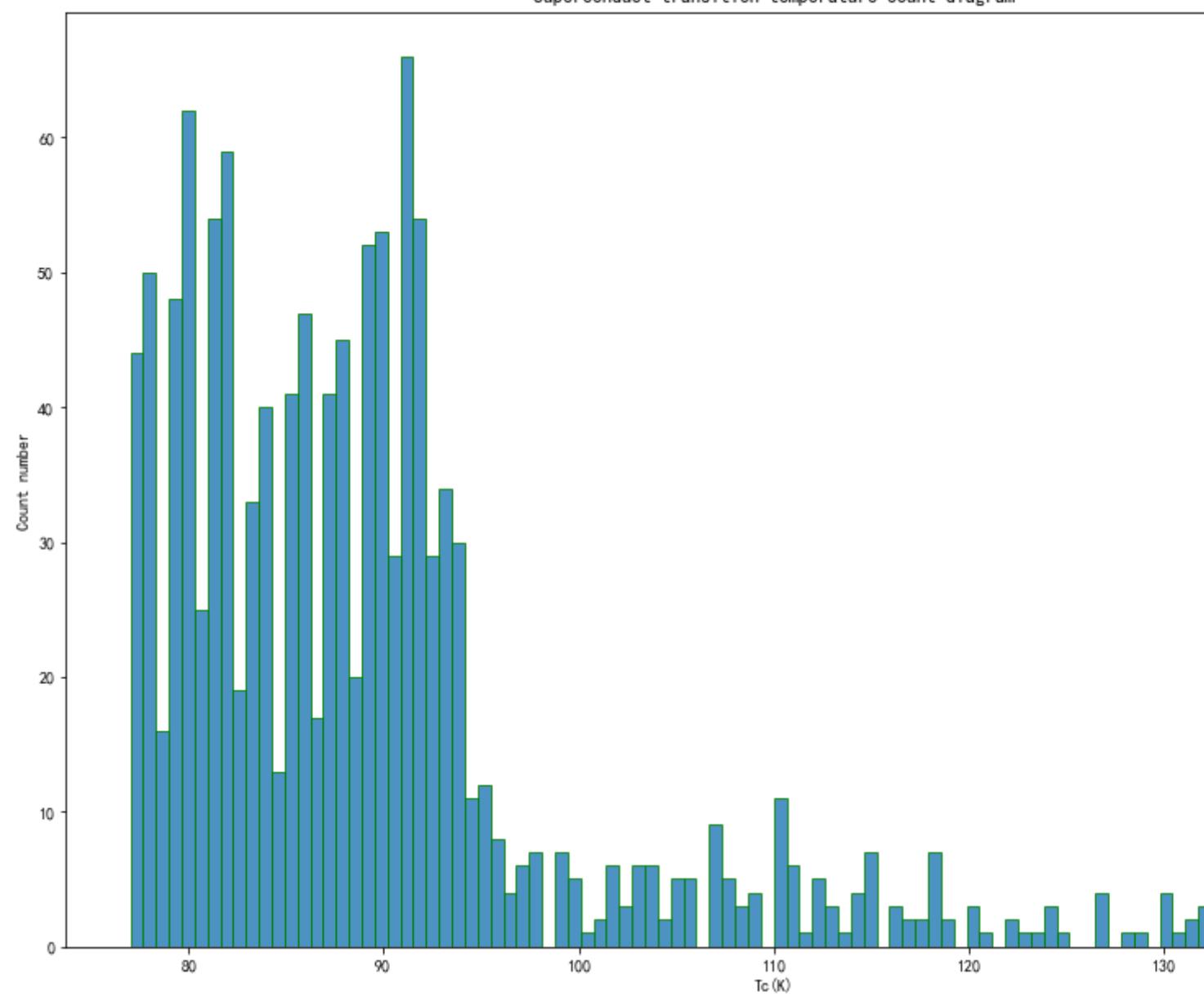
```
1 #Superconducting transition temperature frequency diagram of materials which Tc > 77K
2 fig = plt.gcf()
3 fig.set_size_inches(15.5, 10.5)
4 dataTc = data[data["Tc"]>=77]
5 dataTc
6 matplotlib.rcParams['font.sans-serif']=['SimHei']
7 matplotlib.rcParams['axes.unicode_minus']=False
8 #plt.hist(dataTc, bins=100, log=False, color=None)
9 plt.hist(dataTc["Tc"], bins=100, normed=0, edgecolor="green", alpha=0.8)
10 plt.xlabel("Tc (K)")
11 plt.ylabel("Count number")
12 plt.title('Superconduct transition temperature count diagram')
```

D:\py soft\conda\lib\site-packages\ipykernel_launcher.py:9: MatplotlibDeprecationWarning:
The 'normed' kwarg was deprecated in Matplotlib 2.1 and will be removed in 3.1. Use 'density' instead.
if __name__ == '__main__':

Out[10]: Text(0.5, 1.0, 'Superconduct transition temperature count diagram')

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high-throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE



1.2 Descriptors populate the data overview

1.2.1 The descriptor is generated by AFS1 and AFS2 software

```
In [11]: 1 S1 = pd.read_csv("./S1.csv")
          2 S2 = pd.read_csv("./S2.csv")
```

```
In [12]: 1 S1
```

Out[12]:

	name	Tc	H	He	Li	Be	B	C	N	O	...	Mt	Ds	Rg	Cn	Nh
0	Zr74.1Ti3.9Ni22	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
1	Zr70.2Ti7.8Ni22	3.320	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
2	Zr66.3Ti11.7Ni22	3.560	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
3	Zr62.4Ti15.6Ni22	3.215	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
4	Zr55Cu30Al10Ni5	1.350	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
...	
12335	Ag0.05Rh0.04Ti0.91	1.950	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
12336	Ag0.03Ti0.97	2.670	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
12337	Ag0.035Cd0.01Sn0.955	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
12338	Ag0.005Zn0.995	0.763	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	
12339	Ag0.002Al0.998	1.128	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	...	0	0	0	0	

12340 rows × 120 columns

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the app
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thres
 - 1.3.7.5 Cost sens
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 des
 - 1.4.2 Parameter optim
 - 1.4.3 The trained mod
 - 1.4.4 Residual anal
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throu
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [13]:

1 S2

Out[13]:

	name	Tc	Cell_Min_Number	Cell_Max_Number	Cell_range_Number	C
0	Hg0.66Pb0.34Ba2Ca1.98Cu2.9O8.4	143.000000	1.71253	6.87961	5.16708	
1	Hg0.75Ba2.07Ca2.07Cu3.11O8.208	135.800000	2.55429	7.15202	4.59773	
2	Hg0.75Ba2.07Ca2.07Cu3.11O8.187	135.400000	2.55761	7.16130	4.60369	
3	Hg1Ba2Ca2Cu3O8.29	135.000000	2.45549	6.87538	4.41989	
4	Hg1Ba2Ca2Cu3O8.27	135.000000	2.45851	6.88384	4.42532	
...
12335	Ag0.9Ga0.1	0.006500	3.10000	42.30000	39.20000	
12336	Pt1	0.004157	78.00000	78.00000	0.00000	
12337	La1.92Sr0.08Cu0.99Ni0.01O4	0.001000	0.04000	15.63430	15.59430	
12338	Gd0.9Pr0.95Ce0.15Cu1O4	0.001000	1.24286	8.22857	6.98571	
12339	Au0.978In0.022	0.000500	1.07800	77.26200	76.18400	

12340 rows × 122 columns

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population

 - 1.2.1 The descriptor
 - 1.2.2 The feature space

- 1.3 Classification models

 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appr
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT

 - 1.3.5.1 Visualizati

 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis

 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion ma
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresho
 - 1.3.7.5 Cost sensitiv

- 1.4 Train the deep neural network

 - 1.4.1 Use S1, S2 descript
 - 1.4.2 Parameter optimizat
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis

- 1.5 Virtual high throughput screening

 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions

- 1.6 Deep neural networks

 - 1.6.1 PCA
 - 1.6.2 T-SNE

1.3 Classification model training

1.3.1 Use pycaret to make a rough assessment of multimodels

In [14]:

```

1 CLS1 = pd.read_csv("./CLS1.csv")
2 CLS2 = pd.read_csv("./CLS2.csv")
3 dataset = CLS1
4 dataset

```

Out[14]:

	name	Tc	H	He	Li	Be	B	C	N	O	...	Mt	Ds	R
0	Hg0.66Pb0.34Ba2Ca1.98Cu2.9O8.4	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.515971	...	0	0	0
1	Hg0.75Ba2.07Ca2.07Cu3.11O8.208	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.506417	...	0	0	0
2	Hg0.75Ba2.07Ca2.07Cu3.11O8.187	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.505776	...	0	0	0
3	Hg1Ba2Ca2Cu3O8.29	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.508901	...	0	0	0
4	Hg1Ba2Ca2Cu3O8.27	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.508297	...	0	0	0
...
12335	Ag0.9Ga0.1	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.000000	...	0	0	0
12336	Pt1	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.000000	...	0	0	0
12337	La1.92Sr0.08Cu0.99Ni0.01O4	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.571429	...	0	0	0
12338	Gd0.9Pr0.95Ce0.15Cu1O4	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.571429	...	0	0	0
12339	Au0.978In0.022	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.000000	...	0	0	0

12340 rows × 120 columns

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appr
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sens
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter opti
 - 1.4.3 The trained m
 - 1.4.4 Residual anal
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-th
 - 1.5.2 DNN predictio
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor types
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to setup the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizations
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [15]:

```
1 data = dataset.sample(frac=0.8, random_state=2020).reset_index(drop=True)
2 data_unseen = dataset.drop(data.index).reset_index(drop=True)
3
4 print('Data for Modeling: ' + str(data.shape))
5 print('Unseen Data For Predictions: ' + str(data_unseen.shape))
6 from pycaret.classification import *
7 exp_clf101 = setup(data = data, target = 'Tc', session_id=123)
```

6	Categorical Features	37
7	Ordinal Features	False
8	High Cardinality Features	False
9	High Cardinality Method	None
10	Sampled Data	(9872, 120)
11	Transformed Train Set	(6910, 9991)
12	Transformed Test Set	(2962, 9991)
13	Numeric Imputer	mean
14	Categorical Imputer	constant
15	Normalize	False
16	Normalize Method	None
17	Transformation	False
18	Transformation Method	None

In []:

```
1 compare_models() ##Run to get model ranking
```

In []:

```
1 dataset = CLS2
2 data = dataset.sample(frac=0.8, random_state=2020).reset_index(drop=True)
3 data_unseen = dataset.drop(data.index).reset_index(drop=True)
4
5 print('Data for Modeling: ' + str(data.shape))
6 print('Unseen Data For Predictions: ' + str(data_unseen.shape))
7 from pycaret.classification import *
8 exp_clf101 = setup(data = data, target = 'Tc', session_id=123)
9 compare_models()
```

▼ 1.3.2 Select the appropriate model for classification (RFC,KNN,DT,ADC)

▼ 1.3.3 RFC

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
- ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

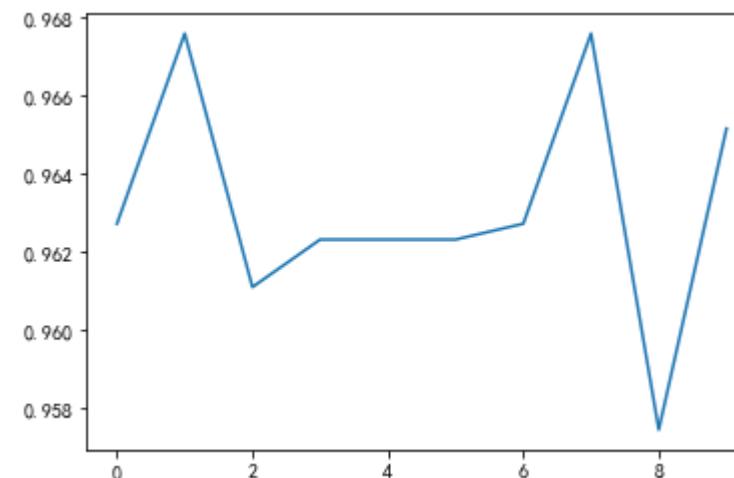
Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [17]:

```
1 # RFC
2 data = CLS1
3 X = data.iloc[:, 2:].values
4 y = data.iloc[:, 1].values
5 from sklearn.model_selection import train_test_split
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2019)
7 rfc = RandomForestClassifier(n_estimators = 200)
8 rfc.fit(X_train,y_train)
9 print("RFCscore",rfc.score(X_test, y_test.astype('int')))
10 cv = ShuffleSplit(n_splits=10, test_size=.2)
11 scores = cross_val_score(rfc, X, y, cv=cv)
12 print(scores)
13 print(scores.mean())
14 ll = np.linspace(0, 9, 10)
15 ll_x = ll
16 ll_y = scores
17 plt.plot(ll_x,ll_y)
18 plt.show()
```

RFCscore 0.9688006482982172
[0.96272285 0.96758509 0.96110211 0.96231767 0.96231767 0.96231767
0.96272285 0.96758509 0.95745543 0.96515397]
0.9631280388978931



1 3 4 KNN

In [18]:

```
1 # KNN
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.model_selection import GridSearchCV
4 param_search = [
5     {
6         "weights": ["uniform"],
7         "n_neighbors": [i for i in range(1, 11)]
8     },
9     {
10        "weights": ["distance"],
11        "n_neighbors": [i for i in range(1, 11)],
12        "p": [i for i in range(1, 6)]
13    }
14]
15 knn_clf = KNeighborsClassifier()
```

In [19]:

```
1 grid_search = GridSearchCV(knn_clf, param_search)
```

In [21]:

```
1 %%time
2 grid_search.fit(X_train, y_train)
3 grid_search.best_estimator_
```

Wall time: 24min 44s

Out[21]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
metric_params=None, n_jobs=None, n_neighbors=6, p=2,
weights='distance')

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizati
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sensi
 - ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter optim
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
 - ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput
 - 1.5.2 DNN prediction
 - ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizati
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitiv
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 data
 - 1.4.2 Parameter optimi
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

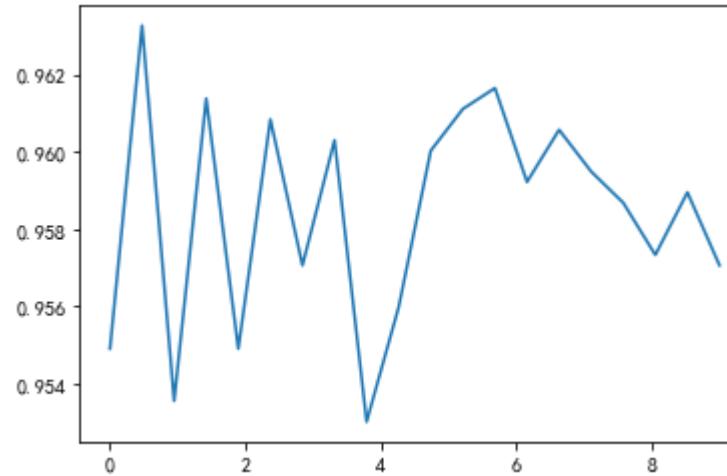
In [22]:

```
1 # KNN
2 X = data.iloc[ :, 2: ].values
3 y = data.iloc[ :, 1 ].values
4 from sklearn.model_selection import train_test_split
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2020)
6 from sklearn.neighbors import KNeighborsClassifier
7 reg = KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
8                             metric_params=None, n_jobs=None, n_neighbors=6, p=2,
9                             weights='distance')
10 reg.fit(X_train, y_train.astype('int'))
11 print(reg.score(X_test, y_test.astype('int')))
12
13 from sklearn.model_selection import ShuffleSplit
14 from sklearn.model_selection import cross_val_score
15 cv = ShuffleSplit(n_splits=20, test_size=.3)
16 scores = cross_val_score(reg, X, y, cv=cv)
17 print(scores)
18 print(scores.mean())
19 ll = np.linspace(0, 9, 20)
20 ll_x = ll
21 ll_y = scores
22 plt.plot(ll_x, ll_y)
23 plt.show()
24
```

```
0.9554294975688817
[0.95488925 0.9632631 0.95353863 0.96137223 0.95488925 0.96083198
 0.95705024 0.96029173 0.95299838 0.95596975 0.96002161 0.96110211
 0.96164236 0.95921124 0.96056186 0.95948136 0.95867099 0.95732037
 0.95894111 0.95705024]
0.9584548892490545
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizati
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresho
 - 1.3.7.5 Cost sensitiv
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimiz
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE



1.3.5 DT

Contents ⚙

1.1 Exploratory analysis
1.2 Descriptors population
1.2.1 The descriptor population
1.2.2 The feature selection
1.3 Classification models
1.3.1 Use pycaret to select the best model
1.3.2 Select the appropriate model
1.3.3 RFC
1.3.4 KNN
1.3.5 DT
1.3.5.1 Visualization
1.3.6 ADC
1.3.7 RFC : Analysis
1.3.7.1 F1 score
1.3.7.2 confusion matrix
1.3.7.3 Accuracy
1.3.7.4 The threshold
1.3.7.5 Cost sensitivity
1.4 Train the deep neural network
1.4.1 Use S1, S2 descriptors
1.4.2 Parameter optimization
1.4.3 The trained model
1.4.4 Residual analysis
1.5 Virtual high throughput screening
1.5.1 Virtual high-throughput screening
1.5.2 DNN predictions
1.6 Deep neural network
1.6.1 PCA
1.6.2 T-SNE

In [24]:

```
1 # DT
2 data = CLS1
3 X = data.iloc[:, 2:].values
4 y = data.iloc[:, 1].values
5 from sklearn.model_selection import train_test_split
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2022)
7 data
```

Out[24]:

	name	Tc	H	He	Li	Be	B	C	N	O	F	Ne	Na	M
0	Hg0.66Pb0.34Ba2Ca1.98Cu2.9O8.4	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.515971	0.0	0	0.0	
1	Hg0.75Ba2.07Ca2.07Cu3.11O8.208	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.506417	0.0	0	0.0	
2	Hg0.75Ba2.07Ca2.07Cu3.11O8.187	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.505776	0.0	0	0.0	
3	Hg1Ba2Ca2Cu3O8.29	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.508901	0.0	0	0.0	
4	Hg1Ba2Ca2Cu3O8.27	0	0.0	0	0.0	0.0	0.0	0.0	0.0	0.508297	0.0	0	0.0	
...	
12335	Ag0.9Ga0.1	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	0	0.0	
12336	Pt1	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	0	0.0	
12337	La1.92Sr0.08Cu0.99Ni0.01O4	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.571429	0.0	0	0.0	
12338	Gd0.9Pr0.95Ce0.15Cu1O4	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.571429	0.0	0	0.0	
12339	Au0.978In0.022	1	0.0	0	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	0	0.0	

12340 rows × 120 columns

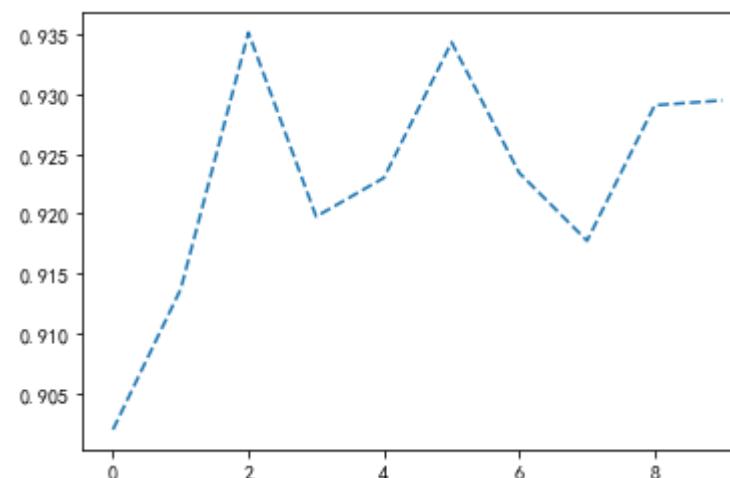
Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [25]:

```
1 clf = tree.DecisionTreeClassifier(max_depth=7, max_features=30, criterion="entropy") #初始化树模型
2 clf = clf.fit(X, y) #实例化训练集
3 score = clf.score(X, y) #返回预测的准确度
4 print("DT", score)
5 cv = ShuffleSplit(n_splits=10, test_size=.2)
6 scores = cross_val_score(clf, X, y, cv=cv)
7 print(scores)
8 print(scores.mean())
9 11 = np.linspace(0, 9, 10)
10 11_x = 11
11 11_y = scores
12 plt.plot(11_x, 11_y, "--")
13 plt.show()
```

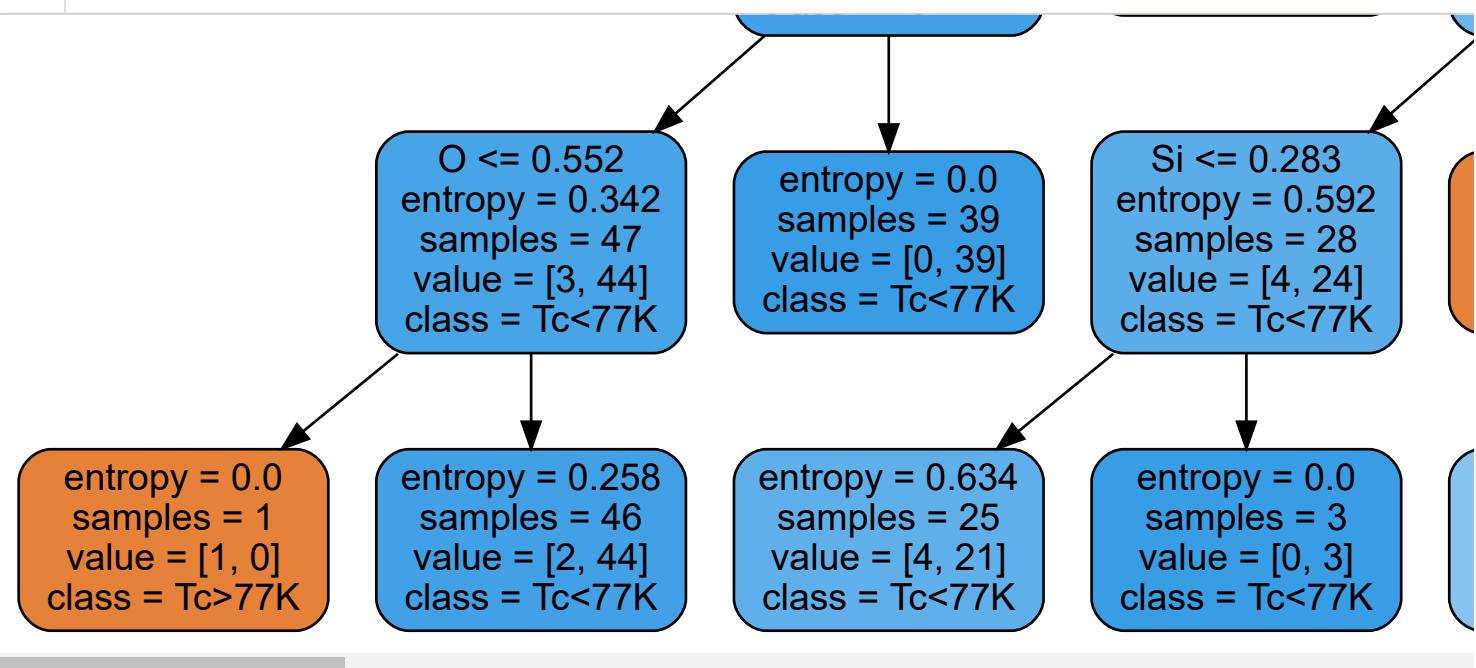
```
DT 0.9289303079416532
[0.90194489 0.9136953 0.93517018 0.9197731 0.92301459 0.93435981
 0.92341977 0.91774716 0.92909238 0.92949757]
0.9227714748784441
```



1.3.5.1 Visualization of decision trees

In [26]:

```
1 feature_name = data.iloc[:, 2:].columns
2 import graphviz
3 dot_data = tree.export_graphviz(clf
4                                         , out_file=None
5                                         , feature_names=feature_name
6                                         , class_names=["Tc>77K", "Tc<77K"]
7                                         , filled=True
8                                         , rounded=True
9                                         )
10 graph = graphviz.Source(dot_data)
11 graph
```



In []:

```
1 # graph.render(filename="tree") ##save DT graph as a pdf
```

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
- 1.2.1 The descriptor
- 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the approp
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sens
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter optim
 - 1.4.3 The trained mo
 - 1.4.4 Residual anal
- 1.5 Virtual high throughput
 - 1.5.1 Virtual high-th
 - 1.5.2 DNN prediction
- 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

1.3.6 ADC

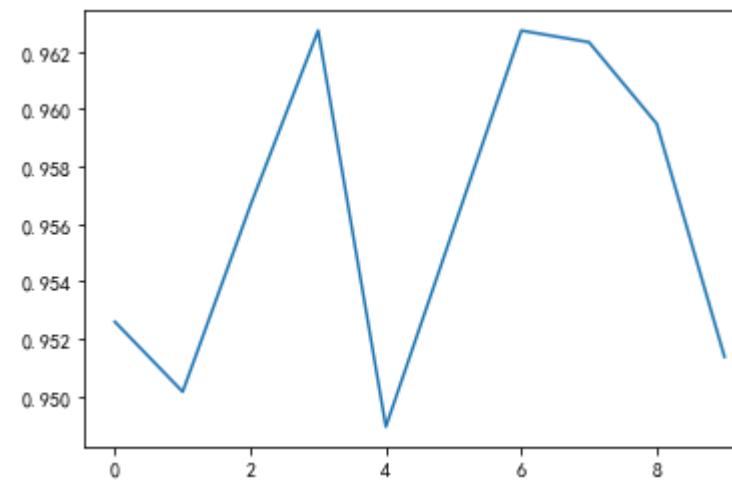
Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [27]:

```
1 # Ada
2 from sklearn.ensemble import AdaBoostClassifier
3 abc = AdaBoostClassifier(n_estimators = 900)
4 data = CLS1
5 X = data.iloc[:, 2:].values
6 y = data.iloc[:, 1].values
7 from sklearn.model_selection import train_test_split
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2019)
9 abc.fit(X_train,y_train)
10 print("ABCscore",abc.score(X_test, y_test.astype('int')))
11 cv = ShuffleSplit(n_splits=10, test_size=.2)
12 scores = cross_val_score(abc, X, y, cv=cv)
13 print(scores)
14 print(scores.mean())
15 ll = np.linspace(0, 9, 10)
16 ll_x = ll
17 ll_y = scores
18 plt.plot(ll_x,ll_y)
19 plt.show()
```

ABCscore 0.9639384116693679
[0.95259319 0.95016207 0.95664506 0.96272285 0.94894652 0.95583468
0.96272285 0.96231767 0.95948136 0.95137763]
0.9562803889789302



1.3.7 RFC : Analysis of multiple achievement indicators

In [28]:

```
1 # RFC
2 data = CLS1
3 X = data.iloc[:, 2:].values
4 y = data.iloc[:, 1].values
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2019)
6 rfc = RandomForestClassifier(n_estimators = 200)
7 rfc.fit(X_train, y_train)
8 rfc.score(X_test, y_test)
```

Out[28]: 0.9679902755267423

In [29]:

```
1 [*zip(feature_name, clf.feature_importances_)]
```

feature_name	clf.feature_importances_
'Ca'	0.0
'Tb'	0.0
'Dy'	0.0
'Ho'	0.0
'Er'	0.0
'Tm'	0.0
'Yb'	0.0
'Lu'	0.0
'Hf'	0.0
'Ta'	0.0
'W'	0.0
'Re'	0.0014127133231240914
'Os'	0.0
'Ir'	0.0
'Pt'	0.0
'Au'	0.0
'Hg'	0.022343358875741987
'Tl'	0.02037171072073104
'Pb'	0.0
'Bi'	0.0
'^ ^'	0.0

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptors
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 data
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high-throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

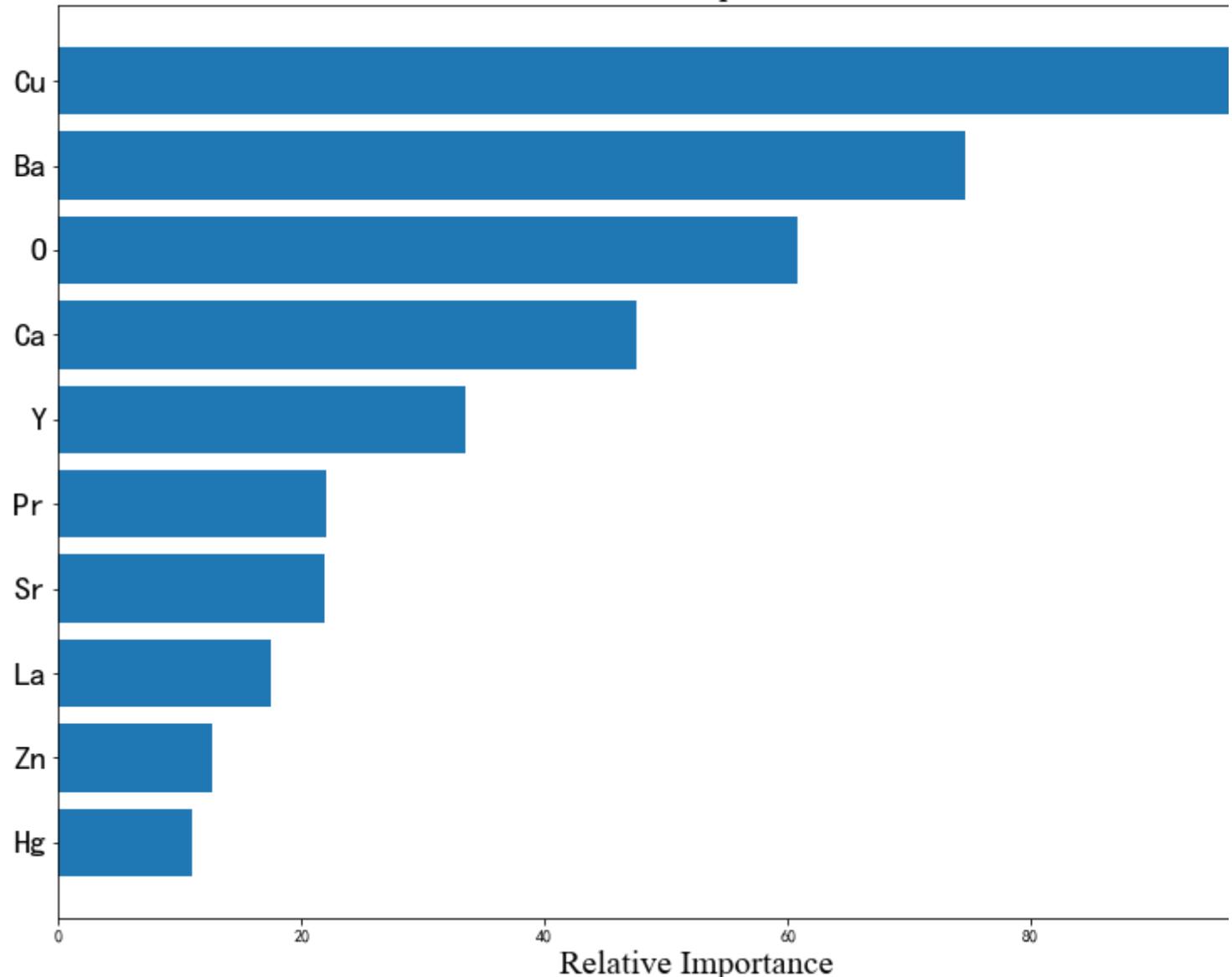
Contents ⚙

1.1 Exploratory analysis
▼ 1.2 Descriptors population
1.2.1 The descriptor population
1.2.2 The feature selection
▼ 1.3 Classification models
1.3.1 Use pycaret to select the best model
1.3.2 Select the appropriate classifier
1.3.3 RFC
1.3.4 KNN
▼ 1.3.5 DT
1.3.5.1 Visualizations
1.3.6 ADC
▼ 1.3.7 RFC : Analysis
1.3.7.1 F1 score
1.3.7.2 confusion matrix
1.3.7.3 Accuracy
1.3.7.4 The threshold
1.3.7.5 Cost sensitivity
▼ 1.4 Train the deep neural network
1.4.1 Use S1, S2 datasets
1.4.2 Parameter optimization
1.4.3 The trained model
1.4.4 Residual analysis
▼ 1.5 Virtual high throughput screening
1.5.1 Virtual high-throughput screening
1.5.2 DNN predictions
▼ 1.6 Deep neural network analysis
1.6.1 PCA
1.6.2 T-SNE

In [30]:

```
1 clf = rfc
2 feature_importance = clf.feature_importances_
3 # make importances relative to max importance
4 feature_importance = 100.0 * (feature_importance / feature_importance.max())
5 sorted_idx = np.argsort(feature_importance)
6 pos = np.arange(sorted_idx.shape[0]) + .5
7 plt.subplot(1, 2, 2)
8 plt.barh(pos[-10:], feature_importance[sorted_idx[-10:]], align='center')
9 plt.yticks(pos[-10:], data.iloc[:, 2:].columns[sorted_idx[-10:]], fontsize=20)
10 plt.xlabel('Relative Importance', font2)
11 plt.title('Variable Importance', font2)
12 fig = plt.gcf()
13 fig.set_size_inches(30, 10)
14 #plt.savefig('Variable Importance', dpi=300, bbox_inches="tight")
15 plt.show()
```

Variable Importance



Contents ⚙️

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

1.3.7.1 F1 score

In [31]:

```
1 from sklearn.metrics import f1_score
2 y_pred = rfc.predict(X_test)
3 f1_score(y_test, y_pred)
```

Out[31]: 0.9823305748154776

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

1.3.7.2 confusion matrix

In [32]:

```
1 from sklearn.metrics import classification_report
2 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.85	0.81	0.83	239
1	0.98	0.99	0.98	2229
accuracy			0.97	2468
macro avg	0.92	0.90	0.91	2468
weighted avg	0.97	0.97	0.97	2468

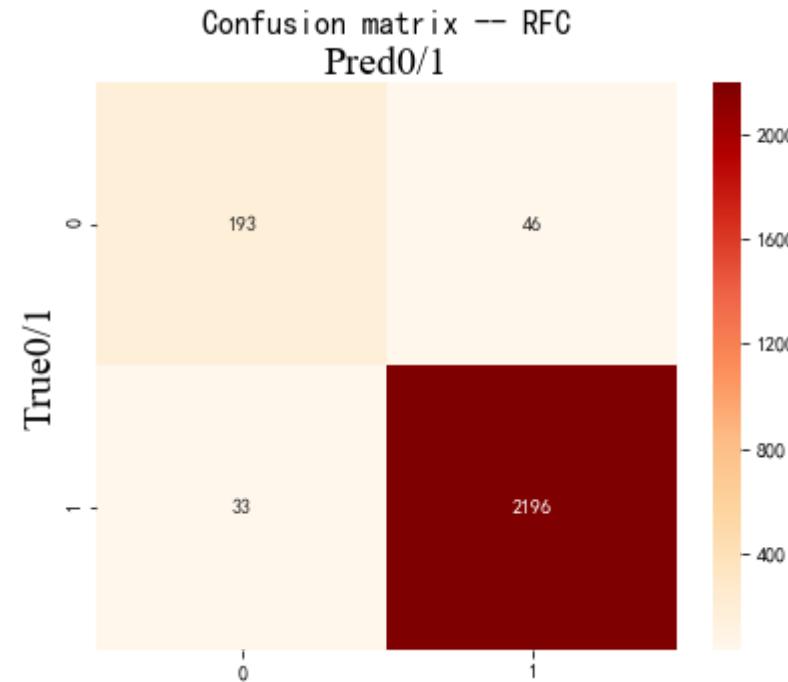
In [33]:

```
1 from sklearn.metrics import confusion_matrix
2 cnf_matrix = confusion_matrix(y_test, y_pred)
3 cnf_matrix
```

Out[33]: array([[193, 46],
 [33, 2196]], dtype=int64)

In [34]: 1 plot_cnf_matirx(cnf_matrix, 'Confusion matrix -- RFC')

<Figure size 1116x756 with 0 Axes>



Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

1.3.7.3 Accuracy and call-back curves

In [35]:

```
1 from sklearn.metrics import precision_recall_curve
2 from sklearn.metrics import roc_auc_score
3 rf_clf = rfc
4 y_probabilities_rf = rf_clf.predict_proba(X_test)[:, 1]
5
6 roc_auc_score(y_test, y_probabilities_rf)
```

Out[35]: 0.9901591985448568

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 data
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE

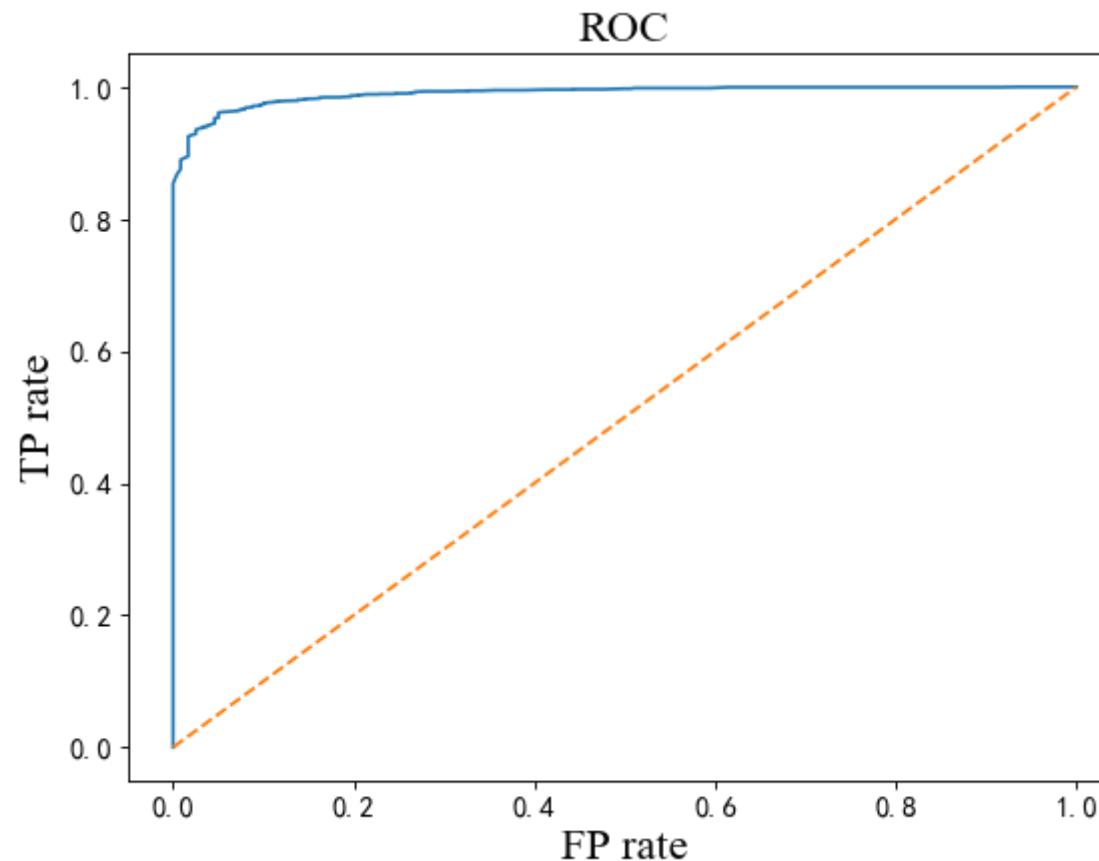
Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [36]:

```
1 from sklearn.metrics import roc_curve
2 fprs4, tprs4, thresholds4 = roc_curve(y_test, y_probabilities_rf)
3 fig = plt.gcf()
4 fig.set_size_inches(15.5, 10.5)
5
6 plot_roc_curve(fprs4, tprs4)
7
8 #plt.savefig('ROC.jpg', dpi=300)
```

⟨Figure size 1116x756 with 0 Axes⟩



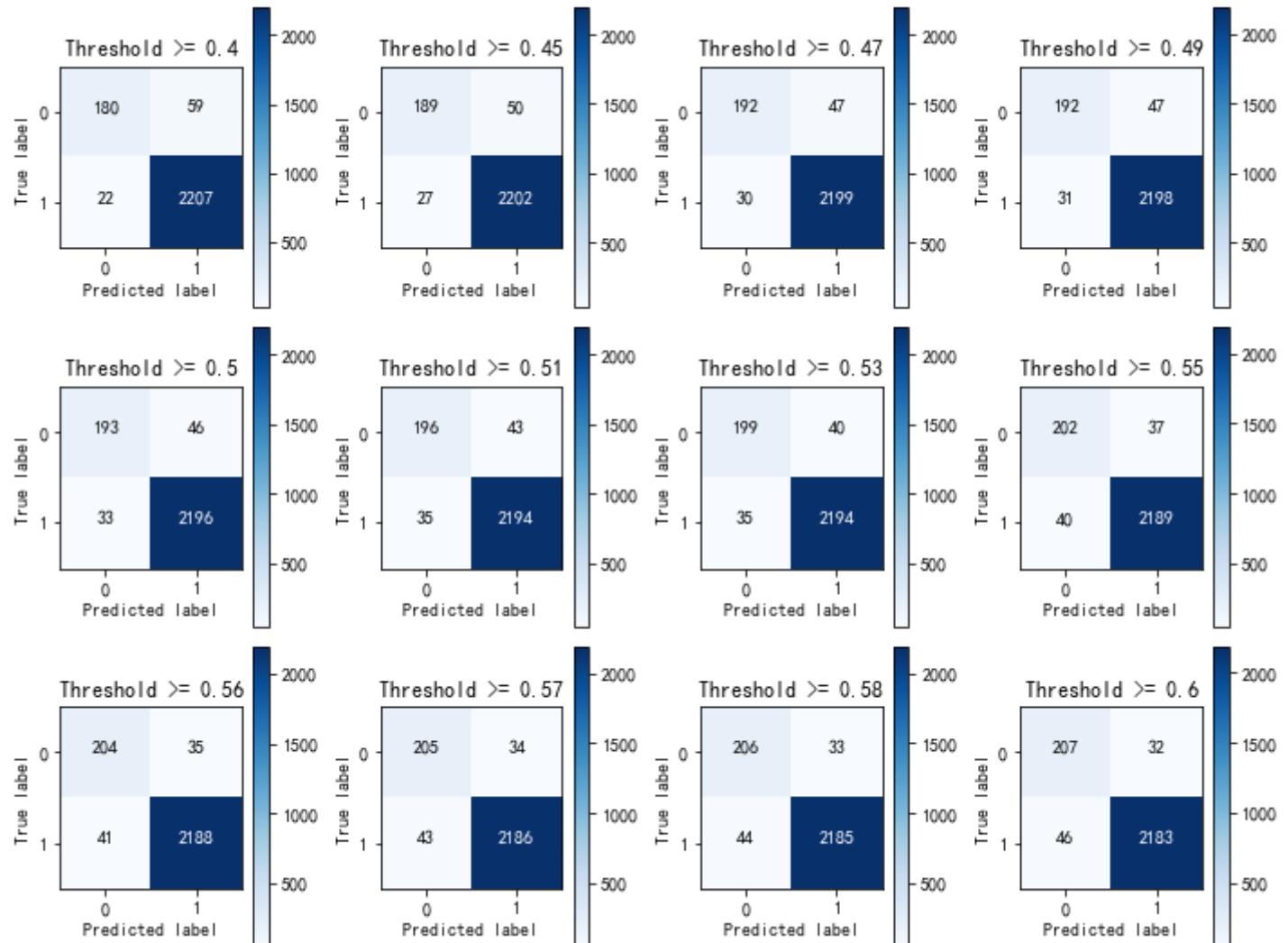
▼ 1.3.7.4 The threshold of model

Contents ☰⚙️

- 1.1 Exploratory analysis
 - ▼ 1.2 Descriptors population
 - 1.2.1 The descriptive statistics
 - 1.2.2 The feature selection
 - ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizations
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
 - ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
 - ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
 - ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate threshold
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitive learning
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE



1.3.7.5 Cost sensitive learning

Bayes formula: $P(A|B)=P(B|A)*P(A)/P(B)$

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appr
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizat
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sens
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter optim
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
In [38]: 1 len(data[data["Tc"]==0])
```

Out[38]: 1233

```
In [39]: 1 len(data[data["Tc"]==1])
```

Out[39]: 11107

```
In [43]: 1 log = []
2 for i in range(100, 900, 5):
3     # print(i/1000)
4     log.append(i/1000)
```

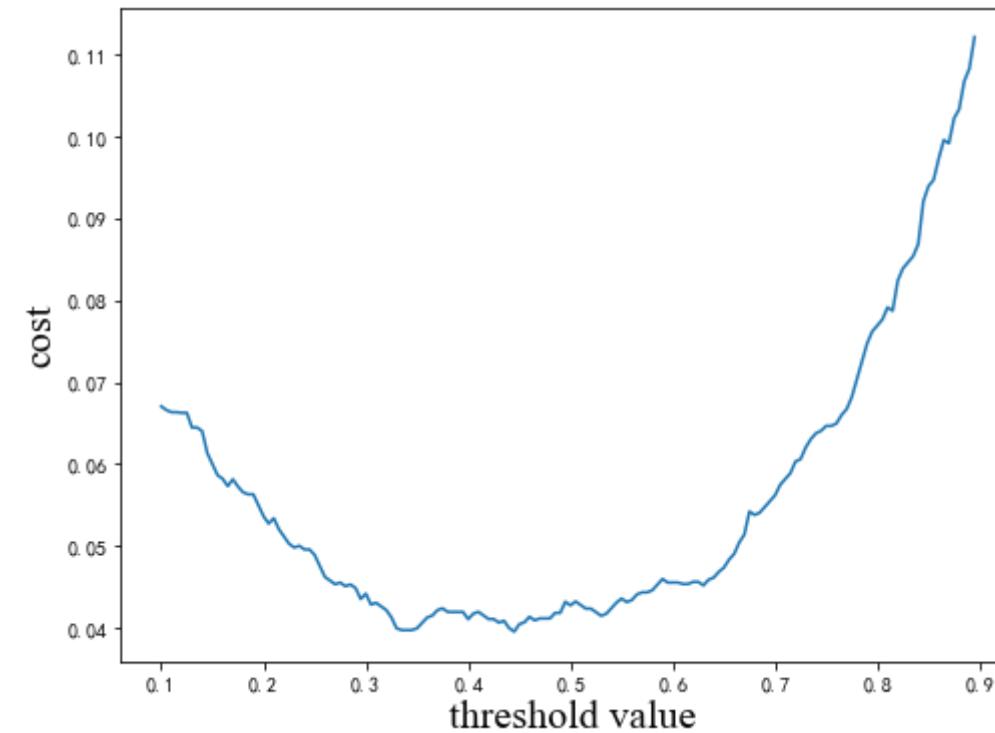
```
In [48]: 1 mg = []
2 thresholds = log
3 m = 1
4 for i in thresholds:
5     y_test_predictions_high_recall = y_probabilities_rf > i
6     cnf_matrix = confusion_matrix(y_test, y_test_predictions_high_recall)
7     np.set_printoptions(precision=2)
8     # print(cnf_matrix)
9     C01 = cnf_matrix[0][1]
10    C10 = cnf_matrix[1][0]
11    C00 = cnf_matrix[0][0]
12    C11 = cnf_matrix[1][1]
13    C = C01+C00+C10+C11
14    cost = (1*(C01)/((C11+C10)))+((1.5*C10)/((C11+C01)))
15    # cost = 1
16    # recall = cnf_matrix[1,1]/(cnf_matrix[1,0]+cnf_matrix[1,1])
17    # acc = (C00+C11)/C
18    # COST = (1/cost)*recall*acc
19    COST=cost
20    # print(COST)
21    mg.append(COST)
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate threshold value
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold value
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [49]:

```
1 mingan = np.array(mingan)
2 log = np.array(log)
3 plt.plot(log, mg)
4 font2 = {'family' : 'Times New Roman',
5 'weight' : 'normal',
6 'size' : 20,
7 }
8 plt.ylabel('cost', font2)
9 plt.xlabel('threshold value', font2)
10 fig = plt.gcf()
11 fig.set_size_inches(8, 6)
12 #plt.figure(figsize=(8, 6))
13 plt.savefig('cost.jpg', dpi=300)
14 plt.show()
```



Contents ⚙️

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizations
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

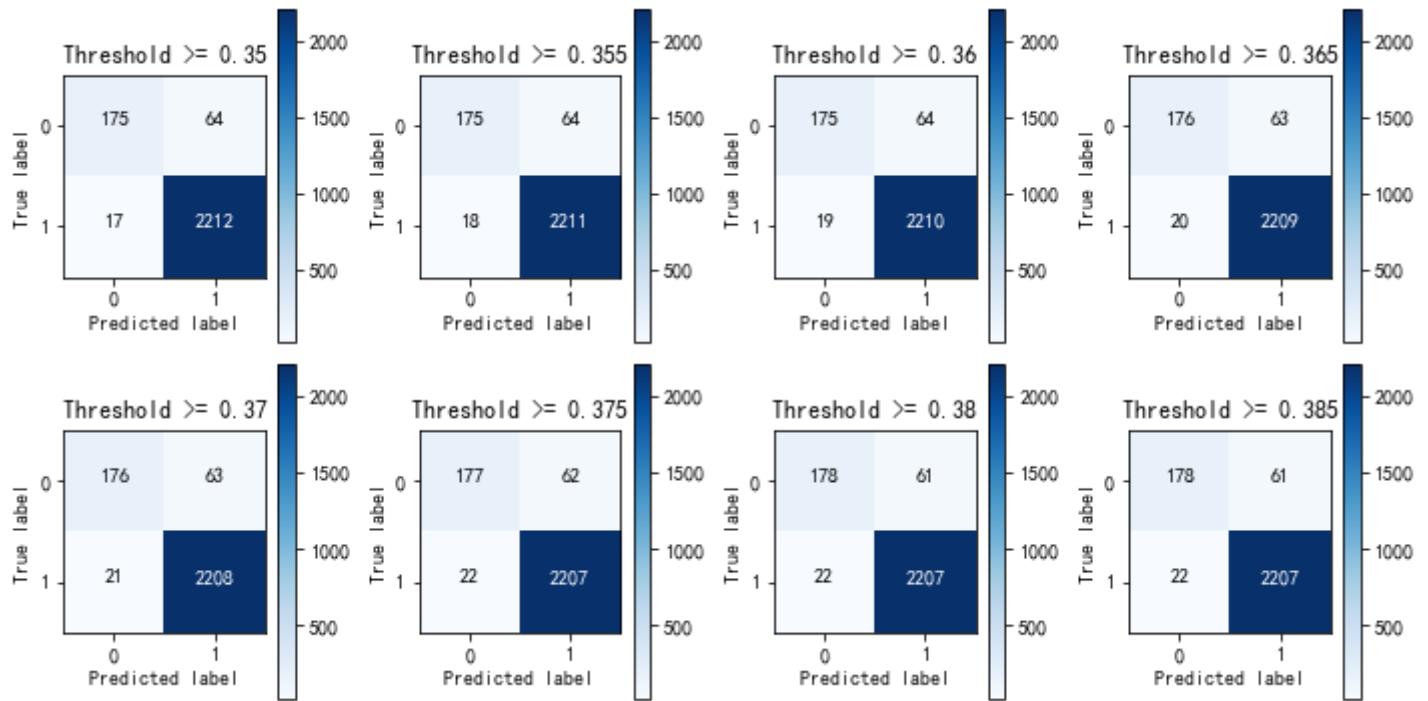
1.1 Exploratory analysis
1.2 Descriptors population
1.2.1 The descriptor population
1.2.2 The feature selection
1.3 Classification models
1.3.1 Use pycaret to select the best model
1.3.2 Select the appropriate model
1.3.3 RFC
1.3.4 KNN
1.3.5 DT
1.3.5.1 Visualization
1.3.6 ADC
1.3.7 RFC : Analysis
1.3.7.1 F1 score
1.3.7.2 confusion matrix
1.3.7.3 Accuracy
1.3.7.4 The threshold
1.3.7.5 Cost sensitivity
1.4 Train the deep neural network
1.4.1 Use S1, S2 descriptors
1.4.2 Parameter optimization
1.4.3 The trained model
1.4.4 Residual analysis
1.5 Virtual high throughput screening
1.5.1 Virtual high-throughput screening
1.5.2 DNN predictions
1.6 Deep neural network
1.6.1 PCA
1.6.2 T-SNE

In [50]:

```
1 thresholds = [0.350, 0.355, 0.360, 0.365, 0.370, 0.375, 0.380, 0.385]
2 plt.figure(figsize=(10, 10))
3
4 m = 1
5 for i in thresholds:
6     y_test_predictions_high_recall = y_probabilities_rf > i
7
8     plt.subplot(4, 4, m)
9     m += 1
10
11    cnf_matrix = confusion_matrix(y_test, y_test_predictions_high_recall)
12    np.set_printoptions(precision=2)
13
14    # print(i, "Recall: {}".format(cnf_matrix[0, 0]/(cnf_matrix[0, 1]+cnf_matrix[0, 0])))
15    # acc = rfc.score(X_test, y_test)
16    # print(acc)
17
18    class_names = [0, 1]
19    plot_confusion_matrix(cnf_matrix,
20                          classes=class_names,
21                          title='Threshold >= %s' % i)
22    # plt.savefig("laser", dpi=300)
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor S1
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to compare models
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE



1.4 Train the deep neural network

1.4.1 Use S1, S2 descriptors for training

In [51]: 1 df = S1

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [73]:

```
1 ## The basic framework of neural network
2 X = df.iloc[:, 2:].values
3 Y = df.iloc[:, 1].values
4 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.25)
5
6 sc = StandardScaler()
7 X_transform = sc.fit_transform(X_train)
8 sc2 = StandardScaler()
9 # y scale to [0, 1]
10
11 mms = MinMaxScaler()
12 Y_transform = mms.fit_transform(y_train.reshape(-1, 1))
13
14 class ProgressBar(keras.callbacks.Callback):
15     def on_epoch_end(self, epoch, logs):
16         self.draw_progress_bar(epoch + 1, EPOCHS)
17
18     def draw_progress_bar(self, cur, total, bar_len=50):
19         cur_len = int(cur / total * bar_len)
20         sys.stdout.write("\r")
21         sys.stdout.write("[{:<{}>} {}/{}]".format("=". * cur_len, bar_len, cur, total))
22         sys.stdout.flush()
23
24     def plot_history(history):
25         hist = pd.DataFrame(history.history)
26         hist['epoch'] = history.epoch
27         plt.figure()
28         plt.xlabel('epoch')
29         plt.ylabel('metric - MSE')
30         plt.plot(hist['epoch'], hist['mse'], label='tr')
31         plt.plot(hist['epoch'], hist['val_mse'], label='va')
32         plt.grid(True)
33         plt.legend()
34
35     plt.figure()
36     plt.xlabel('epoch')
37     plt.ylabel('metric - MAE')
38     plt.plot(hist['epoch'], hist['mae'], label='tr')
39     plt.plot(hist['epoch'], hist['val_mae'], label='va')
40     plt.grid(True)
41     plt.legend()
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
42
43 logdir = '.\callbacks'
44 output_model_file = os.path.join(logdir, "study.h5")
45 tensorboard = keras.callbacks.TensorBoard(logdir)
46 early_stop = keras.callbacks.EarlyStopping(monitor='val_loss', min_delta=1e-5, patience=50)
47
48 callbacks=[ProgressBar()
49             ,early_stop]
50
51 EPOCHS = 50000
52
53 model = tf.keras.Sequential()
54 model.add(tf.keras.layers.Flatten(input_shape=(X.shape[1],)))
55
56 model.add(tf.keras.layers.Dense(128, activation="relu"))
57 model.add(tf.keras.layers.Dense(256, activation="relu"))
58 model.add(tf.keras.layers.Dense(128, activation="relu"))
59 model.add(tf.keras.layers.Dense(128, activation="relu"))
60 model.add(tf.keras.layers.Dense(1, activation="relu"))
61 model.compile(optimizer=tf.keras.optimizers.Adam(0.0001),
62                 loss='mse',
63                 metrics=['mae', 'mse'])
64 )
```

In [74]:

```
1 model.summary()
```

Model: "sequential_9"

Layer (type)	Output Shape	Param #
flatten_9 (Flatten)	(None, 118)	0
dense_39 (Dense)	(None, 128)	15232
dense_40 (Dense)	(None, 256)	33024
dense_41 (Dense)	(None, 128)	32896
dense_42 (Dense)	(None, 128)	16512
dense_43 (Dense)	(None, 1)	129

Total params: 97,793

Trainable params: 97,793

Non-trainable params: 0

Contents ⚙️

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE

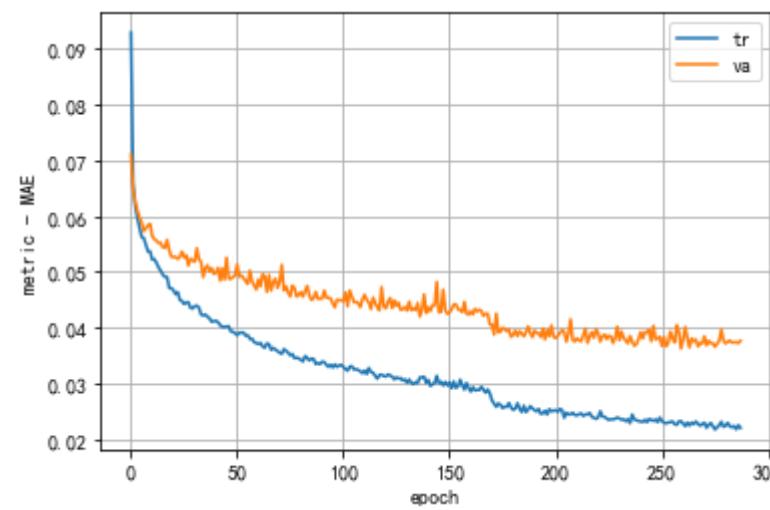
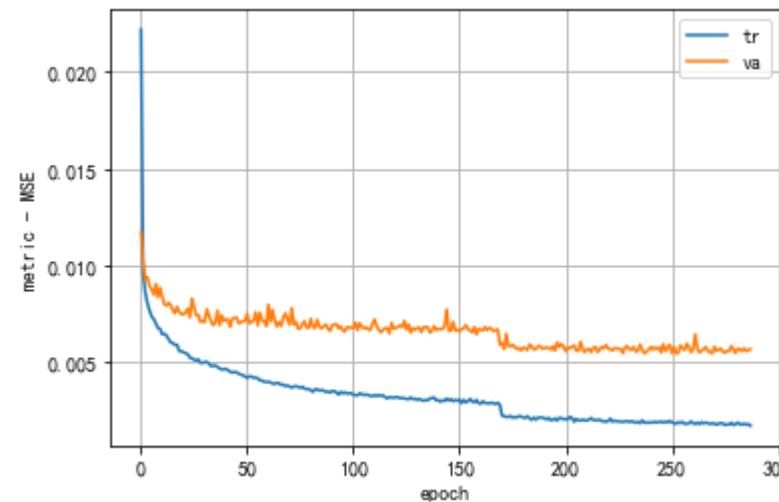
Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [75]:

```
1 history = model.fit(  
2     X_transform, Y_transform,  
3     epochs=EPOCHS, validation_split = 0.2, verbose=1,  
4     batch_size= 32, validation_freq= 1 , callbacks=callbacks)  
ae: 0.0375 - val_mse: 0.0057  
Epoch 283/50000  
7404/7404 [=====] - 0s 66us/sample - loss: 0.0018 - mae: 0.0222 - mse: 0.0018 -  
ae: 0.0377 - val_mse: 0.0056  
Epoch 284/50000  
7404/7404 [=====] - 0s 67us/sample - loss: 0.0017 - mae: 0.0222 - mse: 0.0017 -  
ae: 0.0373 - val_mse: 0.0055  
Epoch 285/50000  
7404/7404 [=====] - 1s 68us/sample - loss: 0.0018 - mae: 0.0223 - mse: 0.0018 -  
ae: 0.0374 - val_mse: 0.0057  
Epoch 286/50000  
7404/7404 [=====] - 1s 71us/sample - loss: 0.0017 - mae: 0.0217 - mse: 0.0017 -  
ae: 0.0374 - val_mse: 0.0055  
Epoch 287/50000  
7404/7404 [=====] - 1s 68us/sample - loss: 0.0018 - mae: 0.0225 - mse: 0.0018 -  
ae: 0.0372 - val_mse: 0.0056  
Epoch 288/50000  
7404/7404 [=====] - 1s 72us/sample - loss: 0.0017 - mae: 0.0220 - mse: 0.0017 -  
ae: 0.0376 - val_mse: 0.0056
```

In [76]: 1 plot_history(history)



Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [77]:

```
1 y_pred = model.predict(sc.transform(X_test))
2 y_pred= mms.inverse_transform(y_pred)
3 y_test = y_test.reshape(-1,1)
4 mae = y_pred-y_test
5 print("mae", np.abs(mae).reshape(1,-1).mean())
6
7 from matplotlib import pyplot as plt
8 plt.scatter(y_test,y_pred)
9 from sklearn.metrics import r2_score
10 r2_score(y_test,y_pred)
11 print("R2",r2_score(y_test,y_pred))
12 plt.title('DNN')
13 fig = plt.gcf()
14 fig.set_size_inches(5.5, 10.5)
15
16 ll = np.linspace(0,130,1200)
17 ll_x = ll
18 ll_y = ll
19 plt.plot(ll_x,ll_y)
20 plt.show()
```

mae 5.567640322635964
R2 0.8715308740832991

1.4.2 Parameter optimization (using additional library optunaas for a longer time)

Contents

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
1 import pyforest
2 import warnings
3 warnings.filterwarnings("ignore")
4 import pathlib
5 import matplotlib.pyplot as plt
6 import pandas as pd
7 import numpy as np
8 import seaborn as sns
9 import tensorflow as tf
10 from tensorflow import keras
11 from tensorflow.keras import layers
12 import os
13
14
15 df = pd.read_csv("S1.csv")
16
17 print(df.shape)
18
19 from sklearn.model_selection import train_test_split
20 X = df.iloc[:, 2:].values
21 Y = df.iloc[:, 1].values
22
23
24 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2)
25
26
27 from sklearn.preprocessing import StandardScaler
28 sc = StandardScaler()
29 X_transform = sc.fit_transform(X_train)
30 sc2 = StandardScaler()
31 # y scale to [0, 1]
32 from sklearn.preprocessing import MinMaxScaler
33 mms = MinMaxScaler()
34 #Y_transform = mms.fit_transform(y_train.reshape(-1, 1))
35 Y_transform = y_train
36
37 Y_transform.shape
38
```

Contents ⚙

1.1 Exploratory analysis
1.2 Descriptors population
1.2.1 The descriptor
1.2.2 The feature selection
1.3 Classification models
1.3.1 Use pycaret to
1.3.2 Select the appr
1.3.3 RFC
1.3.4 KNN
1.3.5 DT
1.3.5.1 Visualizat
1.3.6 ADC
1.3.7 RFC : Analysis
1.3.7.1 F1 score
1.3.7.2 confusion
1.3.7.3 Accuracy
1.3.7.4 The thresh
1.3.7.5 Cost sensi
1.4 Train the deep net
1.4.1 Use S1, S2 de
1.4.2 Parameter opti
1.4.3 The trained m
1.4.4 Residual anal
1.5 Virtual high through
1.5.1 Virtual high-th
1.5.2 DNN predictio
1.6 Deep neural network
1.6.1 PCA
1.6.2 T-SNE

```
39
40 import sys
41 class ProgressBar(keras.callbacks.Callback):
42     def on_epoch_end(self, epoch, logs):
43
44         self.draw_progress_bar(epoch + 1, EPOCHS)
45
46     def draw_progress_bar(self, cur, total, bar_len=50):
47         cur_len = int(cur / total * bar_len)
48         sys.stdout.write("\r")
49         sys.stdout.write("[{:<{}}] {}/{}.format(" = " * cur_len, bar_len, cur, total))
50         sys.stdout.flush()
51
52     early_stop = keras.callbacks.EarlyStopping(monitor='val_loss', min_delta=1e-6, patience=50)
53     callbacks=[ProgressBar(), early_stop]
54
55     EPOCHS = 10000
56
57
58     def objective(trial):
59         kwargs = {}
60         jilu = []
61         # Categorical parameter
62         optimizer = trial.suggest_categorical('optimizer', ["RMSprop", "Adam", "SGD"])
63         if optimizer == "RMSprop":
64             kwargs["learning_rate"] = trial.suggest_float(
65                 "rmsprop_learning_rate", 1e-5, 1e-1, log=True
66             )
67             kwargs["decay"] = trial.suggest_float("rmsprop_decay", 0.85, 0.99)
68             kwargs["momentum"] = trial.suggest_float("rmsprop_momentum", 1e-5, 1e-1, log=True)
69         elif optimizer == "Adam":
70             kwargs["learning_rate"] = trial.suggest_float("adam_learning_rate", 1e-5, 1e-1, log=True)
71         elif optimizer == "SGD":
72             kwargs["learning_rate"] = trial.suggest_float(
73                 "sgd_opt_learning_rate", 1e-5, 1e-1, log=True
74             )
75             kwargs["momentum"] = trial.suggest_float("sgd_opt_momentum", 1e-5, 1e-1, log=True)
76
77         # Int parameter
78         sj_n = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16]
79         sj_n[0] = trial.suggest_int('sj_n[0]', 32, 256)
80         sj_n[1] = trial.suggest_int('sj_n[1]', 32, 256)
```

Contents ⚙

1.1 Exploratory analysis	81	sj_n[2] = trial.suggest_int('sj_n[2]', 32, 256)
1.2 Descriptors population	82	sj_n[3] = trial.suggest_int('sj_n[3]', 32, 256)
1.2.1 The descriptor types	83	sj_n[4] = trial.suggest_int('sj_n[4]', 32, 256)
1.2.2 The feature selection	84	sj_n[5] = trial.suggest_int('sj_n[5]', 32, 256)
1.3 Classification models	85	sj_n[6] = trial.suggest_int('sj_n[6]', 32, 256)
1.3.1 Use pycaret to select the best model	86	sj_n[7] = trial.suggest_int('sj_n[7]', 32, 256)
1.3.2 Select the appropriate classifier	87	sj_n[8] = trial.suggest_int('sj_n[8]', 32, 256)
1.3.3 RFC	88	sj_n[9] = trial.suggest_int('sj_n[9]', 32, 256)
1.3.4 KNN	89	sj_n[10] = trial.suggest_int('sj_n[10]', 32, 256)
1.3.5 DT	90	sj_n[11] = trial.suggest_int('sj_n[11]', 32, 256)
1.3.5.1 Visualization	91	sj_n[12] = trial.suggest_int('sj_n[12]', 32, 256)
1.3.6 ADC	92	sj_n[13] = trial.suggest_int('sj_n[13]', 32, 256)
1.3.7 RFC : Analysis	93	sj_n[14] = trial.suggest_int('sj_n[14]', 32, 256)
1.3.7.1 F1 score	94	sj_n[15] = trial.suggest_int('sj_n[15]', 32, 256)
1.4 Train the deep neural network	95	
1.4.1 Use S1, S2 descriptors	96	
1.4.2 Parameter optimization	97	# Uniform parameter
1.4.3 The trained model	98	dropout_rate = trial.suggest_uniform('dropout_rate', 0.0, 0.9)
1.4.4 Residual analysis	99	
1.5 Virtual high throughput screening	100	# Loguniform parameter
1.5.1 Virtual high-throughput screening	101	learning_rate = trial.suggest_loguniform('learning_rate', 1e-5, 1e-2)
1.5.2 DNN predictions	102	
1.6 Deep neural network	103	
1.6.1 PCA	104	num_layers = trial.suggest_int('num_layers', 3, 15)
1.6.2 T-SNE	105	model = tf.keras.Sequential()
	106	model.add(layers.Flatten(input_shape=(X.shape[1],)))
	107	test_n = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16]
	108	test_n[0] = trial.suggest_int('test_n[0]', 1, 2)
	109	test_n[1] = trial.suggest_int('test_n[1]', 1, 2)
	110	test_n[2] = trial.suggest_int('test_n[2]', 1, 2)
	111	test_n[3] = trial.suggest_int('test_n[3]', 1, 2)
	112	test_n[4] = trial.suggest_int('test_n[4]', 1, 2)
	113	test_n[5] = trial.suggest_int('test_n[5]', 1, 2)
	114	test_n[6] = trial.suggest_int('test_n[6]', 1, 2)
	115	test_n[7] = trial.suggest_int('test_n[7]', 1, 2)
	116	test_n[8] = trial.suggest_int('test_n[8]', 1, 2)
	117	test_n[9] = trial.suggest_int('test_n[9]', 1, 2)
	118	test_n[10] = trial.suggest_int('test_n[10]', 1, 2)
	119	test_n[11] = trial.suggest_int('test_n[11]', 1, 2)
	120	
	121	
	122	

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
123     test_n[12] = trial.suggest_int('test_n[12]', 1, 2)
124     test_n[13] = trial.suggest_int('test_n[13]', 1, 2)
125     test_n[14] = trial.suggest_int('test_n[14]', 1, 2)
126     test_n[15] = trial.suggest_int('test_n[15]', 1, 2)
127
128
129     for i in range(num_layers):
130         model.add(layers.Dense(sj_n[i], activation='relu'))
131         if test_n[i] == 1:
132             model.add(layers.BatchNormalization())
133             model.add(layers.Dropout(dropout_rate))
134
135         model.add(layers.Dense(1, activation='relu'))
136
137     model.compile(optimizer=optimizer, loss='mse', metrics=['mae', 'mse'])
138     history = model.fit(X_transform, Y_transform, epochs=EPOCHS, validation_split=0.2, verbose=0,
139                         validation_freq=1, callbacks=callbacks)
140     y_pred = model.predict(sc.transform(X_test))
141     global y_test
142
143     y_test = y_test.reshape(-1, 1)
144
145     mae = y_pred - y_test
146
147
148     import numpy as np
149     from sklearn.metrics import r2_score
150     global r2
151     r2 = r2_score(y_test, y_pred)
152
153     MAE = np.abs(mae).reshape(1, -1).mean()
154     import time
155
156     ts = str(time.time())
157     saved_model_path = "./saved999_models/{}".format(str(r2)) + "R2" + "|||||MAE" + str(
158         np.abs(mae).reshape(1, -1).mean()) + "TM" + ts
159     if r2 >= 0.91:
160         tf.keras.experimental.export_saved_model(model, saved_model_path)
161         model.save('optmodel.h5')
162
163     print("R2", r2)
164     print("mae", str(np.abs(mae).reshape(1, -1).mean()))
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor S1
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
165     return r2
166
167
168 import optuna
169 study = optuna.create_study(direction='maximize')
170 study.optimize(objective, n_trials=10000)
171 print(study.best_params)
172 print(study.best_value)
```

1.4.3 The trained model, take DNN under S1 descriptor as an example

In [78]:

```
1 save_model = tf.keras.experimental.load_from_saved_model("./S1model")
2 #save_model.predict(X).shape
3 save_model.summary()
```

Model: "sequential_418"

Layer (type)	Output Shape	Param #
flatten_418 (Flatten)	(None, 118)	0
dense_2105 (Dense)	(None, 251)	29869
dense_2106 (Dense)	(None, 153)	38556
dense_2107 (Dense)	(None, 239)	36806
dense_2108 (Dense)	(None, 124)	29760
dense_2109 (Dense)	(None, 1)	125

Total params: 135,116

Trainable params: 135,116

Non-trainable params: 0

1.4.4 Residual analysis

Contents ⚙

1.1 Exploratory analysis
1.2 Descriptors population
1.2.1 The descriptor
1.2.2 The feature selection
1.3 Classification models
1.3.1 Use pycaret to
1.3.2 Select the approp
1.3.3 RFC
1.3.4 KNN
1.3.5 DT
1.3.5.1 Visualizati
1.3.6 ADC
1.3.7 RFC : Analysis
1.3.7.1 F1 score
1.3.7.2 confusion
1.3.7.3 Accuracy
1.3.7.4 The thresho
1.3.7.5 Cost sensitiv
1.4 Train the deep neural
1.4.1 Use S1, S2 de
1.4.2 Parameter optimi
1.4.3 The trained model
1.4.4 Residual analysis
1.5 Virtual high throughput
1.5.1 Virtual high-throughput
1.5.2 DNN prediction
1.6 Deep neural network
1.6.1 PCA
1.6.2 T-SNE

In [79]:

```
1 ## take S2model for a example
2 model = tf.keras.experimental.load_from_saved_model("./S2model")
3 df = S2
4 from sklearn.model_selection import train_test_split
5 X = df.iloc[:, 2:].values
6 Y = df.iloc[:, 1].values
7 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.1)
8
9
10 from sklearn.preprocessing import StandardScaler
11 sc = StandardScaler()
12 X_transform = sc.fit_transform(X_train)
13 sc2 = StandardScaler()
14
15 from sklearn.preprocessing import MinMaxScaler
16
17 y_pred = model.predict(sc.transform(X_test))
18
19 from matplotlib import pyplot as plt
20 plt.scatter(y_test, y_pred)
21 from sklearn.metrics import r2_score
22 r2_score(y_test, y_pred)
23 print(r2_score(y_test, y_pred))
24 from sklearn.metrics import mean_squared_error, mean_absolute_error
25
26
27 np.sqrt(mean_squared_error(y_test, y_pred))
28
29
30 print("RMSE", np.sqrt(mean_squared_error(y_test, y_pred)))
31
32
33 mean_absolute_error(y_test, y_pred)
34
35 print("MAE", mean_absolute_error(y_test, y_pred))
36
37 font2 = {'family': 'Times New Roman',
38          'weight': 'normal',
39          'size': 12,
40          }
41 plt.title('Deeplearning')
```

Contents ⚙

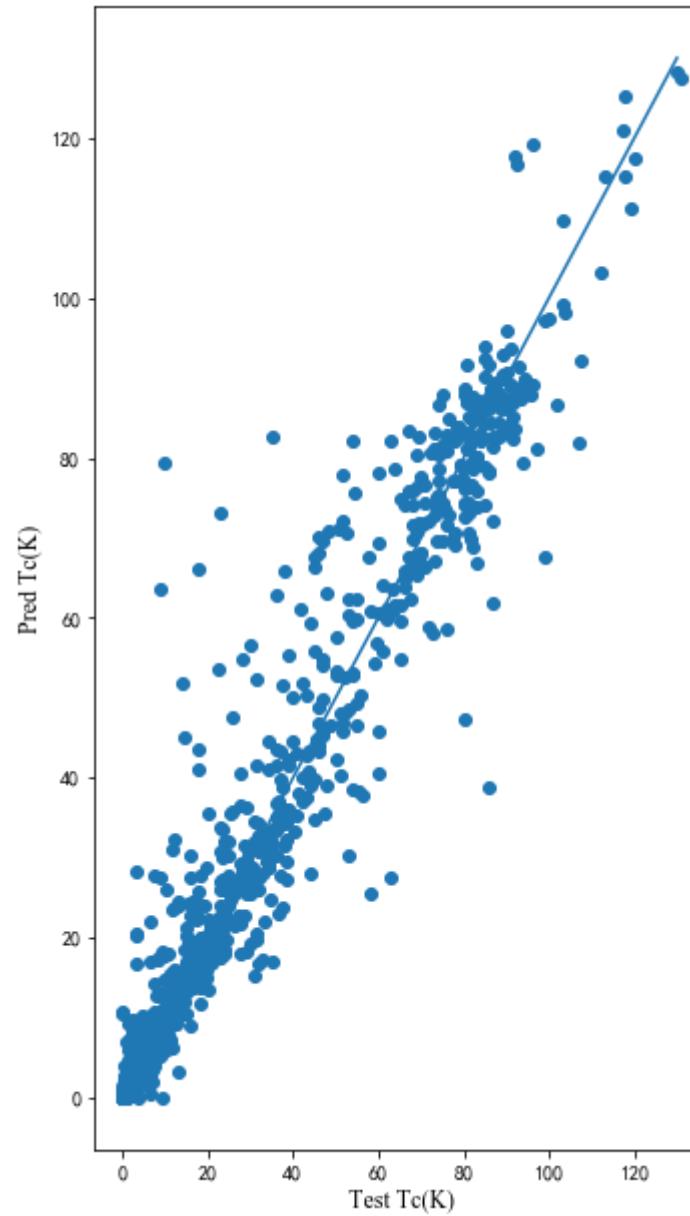
- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 data
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
42 fig = plt.gcf( )
43 fig.set_size_inches(5.5, 10.5)
44 plt.xlabel("Test Tc (K)", font2)
45
46 plt.ylabel("Pred Tc (K)", font2)
47
48 ll = np.linspace(0, 130, 1200)
49 ll_x = ll
50 ll_y = ll
51 plt.plot(ll_x, ll_y)
52 #plt.savefig('DNN回归.jpg', dpi=300)
53 plt.show()
```

0.9361781128924711
RMSE 7.345803863368909
MAE 3.731645200334191

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high-throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE



Contents ⚙

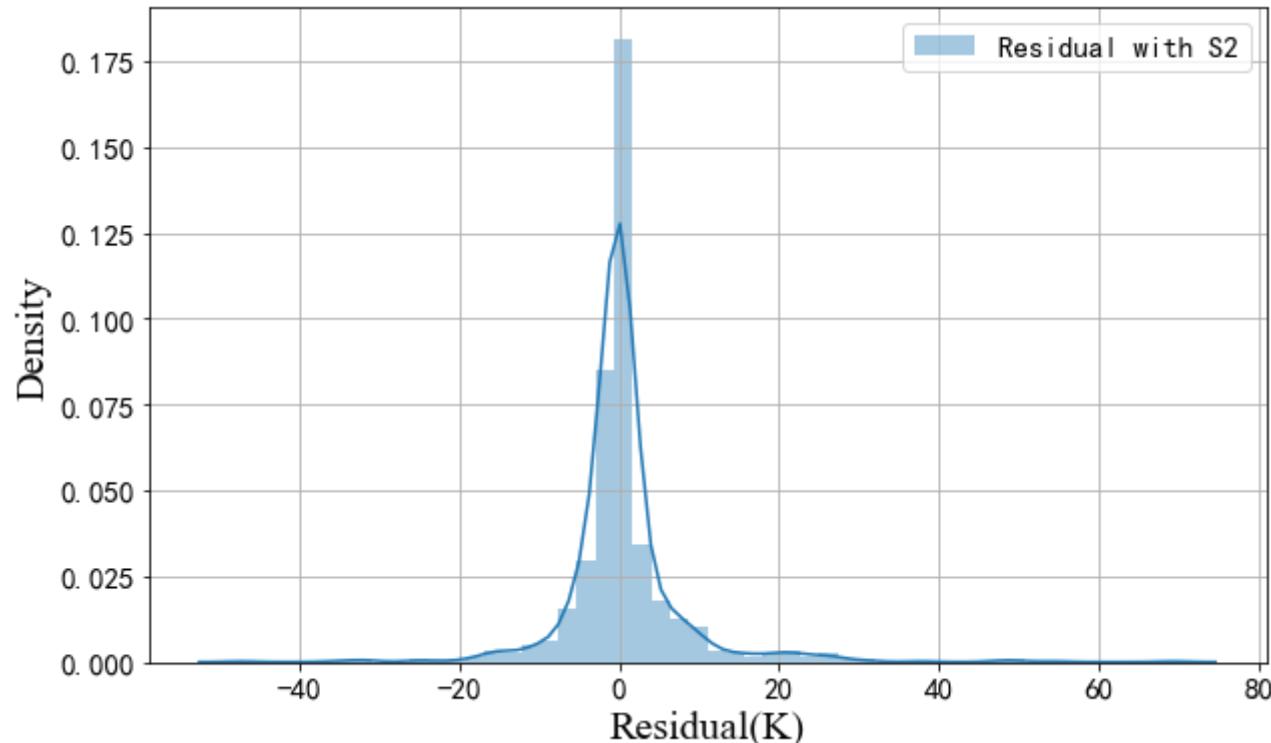
- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
 - ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
 - ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
 - ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [80]:

```
1 error = y_pred.T - y_test.T
2 error = pd.DataFrame(error)
3 fig = plt.figure(figsize = (10, 6))
4 #ax1 = fig.add_subplot(2, 1, 1) # 创建子图1
5 #ax1.scatter(s.index, s.values)
6 #plt.grid()
7 # 绘制数据分布图
8 font3 = {'family' : 'Times New Roman',
9          'weight' : 'normal',
10         'size' : 20,
11         }
12 ax = fig.add_subplot(1, 1, 1) # 创建子图2
13 sns.hist(bins=30, alpha = 0.5, ax = ax)
14 sns.distplot(error, ax = ax, bins = 50)
15 plt.legend(["Residual with S2"], fontsize =15)
16 plt.xlabel('Residual(K)', font3)
17 plt.ylabel('Density', font3)
18 ax.tick_params(axis='y', labelsize=15)
19 ax.tick_params(axis='x', labelsize=15)
20 #plt.title('Pred of VH')
21 #plt.savefig('误差分析.jpg', dpi=300)
22 plt.grid()
23 #plt.savefig("S2 residual", dpi=300, bbox_inches ="tight")
```

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput Settings
 - 1.5.1 Virtual high-throughput sample generation
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE



1.5 Virtual high throughput Settings

1.5.1 Virtual high-throughput sample generation (the following code is virtual sample distribution in HgCaBaPbCuO compounds)

```
1  ### It should be noted that the sample distribution generated by virtual high flux conforms to the distribution of training data as much as possible
2  vh = []
3  import numpy
4  import pandas as pd
5  x = numpy.arange(1, 300, 6)
6  x= x/1000
7  x = list(x)
8  print(x)
9  for i in x:
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

```
10     print(i)
11     vh = []
12     data =[ {'Hg':[[0, 0.26, 0.01]]}, {'Pb':[[0.09, 0.34, 0.01]]}, {'Ba':[[i, i, i]]}, {'Ca':[[0, 0.2
13     def findall(pos, sum, outcome):
14         if pos==len(data)-1:
15             if 1/sum >=0.7 and 1/sum<=2.4:
16                 outcome+="0"
17                 print(outcome)
18                 vh.append(outcome)
19             # -----output
20         else:
21             for (key, val) in data[pos].items():
22                 for vl in val:
23                     st = vl[0]
24                     if st==0:
25                         st=st+vl[2]
26                     en = vl[1]
27                     bu = vl[2]
28                     while(st<=en):
29                         findall(pos+1, sum+round(st, 6), outcome+str(key)+str(round(s
30                         st+=bu
31             findall(0, 0, '')
32
33             test=pd.DataFrame(vh)
34             print(test)
35             test.to_csv('./'+str(i)+'. Ba.csv', encoding='gbk', index=0, header=0)
```

1.5.2 DNN prediction of virtual samples

1 see it in DNN-VH.ipynb

1.6 Deep neural network combined with manifold learning

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [83]:

```
1 import tensorflow as tf
2 saved_model_path = "path_to\modell"
3 new_model = tf.keras.experimental.load_from_saved_model(saved_model_path)
4 new_model.summary()
5 model = new_model
6 a = []
7 for i in model.get_weights():
8
9     a.append(i[0])
10
11 from tensorflow.keras import Sequential
12 model2 = tf.keras.Sequential()
13 print(len(model.layers))
14 print(len(model.layers[:-1]))
15 for layer in model.layers[:-1]:
16     model2.add(layer)
17
18 model2.summary()
19 d = []
20 for i in model2.get_weights():
21     d.append(i[0])
22 print("-----")
```

Model: "sequential_418"

Layer (type)	Output Shape	Param #
flatten_418 (Flatten)	(None, 118)	0
dense_2105 (Dense)	(None, 251)	29869
dense_2106 (Dense)	(None, 153)	38556
dense_2107 (Dense)	(None, 239)	36806
dense_2108 (Dense)	(None, 124)	29760
dense_2109 (Dense)	(None, 1)	125

Total params: 135,116

Trainable params: 135,116

Non-trainable params: 0

6

5

Model: "sequential_11"

Contents ⚙️

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor types
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Layer (type)	Output Shape	Param #
flatten_418 (Flatten)	(None, 118)	0
dense_2105 (Dense)	(None, 251)	29869
dense_2106 (Dense)	(None, 153)	38556
dense_2107 (Dense)	(None, 239)	36806
dense_2108 (Dense)	(None, 124)	29760

Total params: 134,991

Trainable params: 134,991

Non-trainable params: 0

In [84]:

```
1 df = pd.read_csv("S1.csv")
2 df
```

Out[84]:

	name	Tc	H	He	Li	Be	B	C	N	O	F	Ne	Na	Mg	Al	Si
0	Zr74.1Ti3.9Ni22	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
1	Zr70.2Ti7.8Ni22	3.320	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
2	Zr66.3Ti11.7Ni22	3.560	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
3	Zr62.4Ti15.6Ni22	3.215	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
4	Zr55Cu30Al10Ni5	1.350	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.100	0
...
12335	Ag0.05Rh0.04Ti0.91	1.950	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12336	Ag0.03Ti0.97	2.670	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12337	Ag0.035Cd0.01Sn0.955	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12338	Ag0.005Zn0.995	0.763	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12339	Ag0.002Al0.998	1.128	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.998	0

12340 rows × 120 columns

In [86]:

```
1 y_pred = model2.predict(sc.transform(X))
2 len(y_pred)
```

Out[86]:

12340

In [87]:

```

1 y_pred = pd.DataFrame(y_pred)
2 newdf = pd.concat([df, y_pred], axis=1)
3 newdf

```

Out[87]:

	name	Tc	H	He	Li	Be	B	C	N	O	F	Ne	Na	Mg	Al	Si
0	Zr74.1Ti3.9Ni22	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
1	Zr70.2Ti7.8Ni22	3.320	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
2	Zr66.3Ti11.7Ni22	3.560	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
3	Zr62.4Ti15.6Ni22	3.215	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
4	Zr55Cu30Al10Ni5	1.350	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.100	0
...
12335	Ag0.05Rh0.04Ti0.91	1.950	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12336	Ag0.03Ti0.97	2.670	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12337	Ag0.035Cd0.01Sn0.955	3.650	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12338	Ag0.005Zn0.995	0.763	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.000	0
12339	Ag0.002Al0.998	1.128	0.0	0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0	0.0	0.0	0.998	0

12340 rows × 244 columns

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appr
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualizati
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion ma
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresho
 - 1.3.7.5 Cost sensit
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter optim
 - 1.4.3 The trained mod
 - 1.4.4 Residual anal
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throu
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

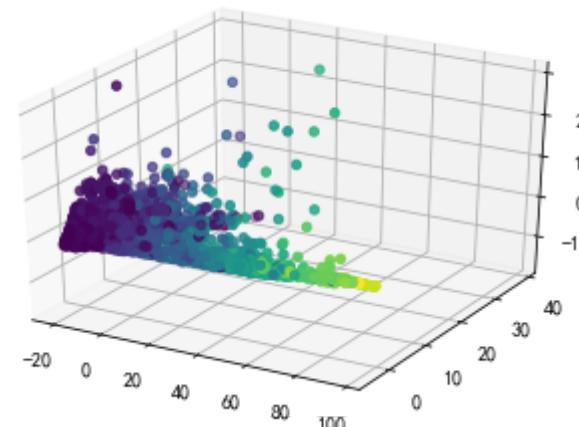
1.6.1 PCA

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sens
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter optim
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput
 - 1.5.2 DNN prediction
- 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [88]:

```
1 from sklearn.decomposition import PCA
2 ax = plt.axes(projection='3d')
3 x = y_pred
4 pca = PCA(n_components=3)
5 pca.fit(x)
6 x_pca = pca.transform(x)
7 newdf = np.array(newdf)
8 close = newdf[:, 1]
9 volume = newdf[:, 1]
10 import seaborn as sns
11 from matplotlib import pyplot as plt
12 plt.figure(figsize=(20, 10))
13 #plt.scatter(x_pca[:, 0], x_pca[:, 1], alpha=0.2, c=close)
14 ax.scatter3D(x_pca[:, 0], x_pca[:, 1], x_pca[:, 2], c=close)
15 plt.show()
16 plt.show()
```



Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizations
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

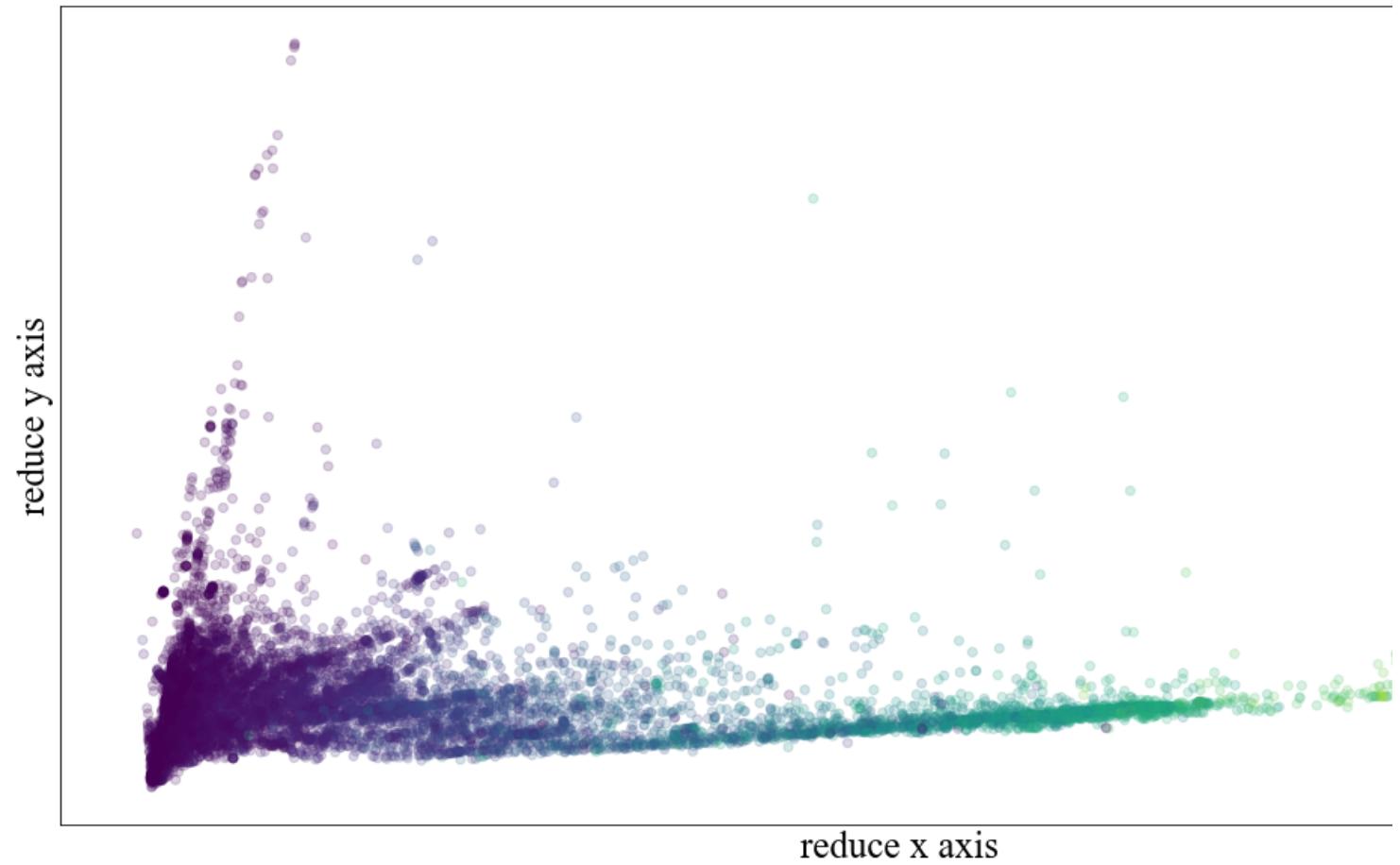
- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizati
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitiv
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimiz
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [121]:

```
1 font4 = {'family' : 'Times New Roman',
2      'weight' : 'normal',
3      'size' : 25,
4      }
5 from sklearn.decomposition import PCA
6 #ax = plt.axes(projection='3d')
7 x = y_pred
8 pca = PCA(n_components=2)
9 pca.fit(x)
10 x_pca = pca.transform(x)
11 newdf = np.array(newdf)
12 close = newdf[:, 1]
13 volume = newdf[:, 1]
14 plt.figure(figsize=(20, 10))
15 plt.scatter(x_pca[:, 0], x_pca[:, 1], alpha=0.2, c=close, label="PCA")
16 #ax.scatter3D(x_pca[:, 0], x_pca[:, 1], x_pca[:, 2], c=close)
17 plt.xlabel('reduce x axis', font4)
18 plt.ylabel('reduce y axis', font4)
19 plt.legend(loc='upper right', fontsize= 25)
20 plt.xticks([])
21 plt.yticks([])
22 plt.savefig("realPCA.png", dpi=300, bbox_inches="tight")
23 plt.show()
```

Contents ⚙

- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE



```
1 for i in range(2, 120, 1):  
2     newdf = np.array(newdf)  
3     close = newdf[:, i]  
4     print("elem Number", i-1)  
5  
6     import seaborn as sns  
7     from matplotlib import pyplot as plt  
8     plt.figure(figsize=(10, 5))  
9     plt.scatter(x_pca[:, 0], x_pca[:, 1], alpha=0.2 ,c =close )  
10    plt.show()
```



1.6.2 T-SNE

Contents ⚙️

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizations
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

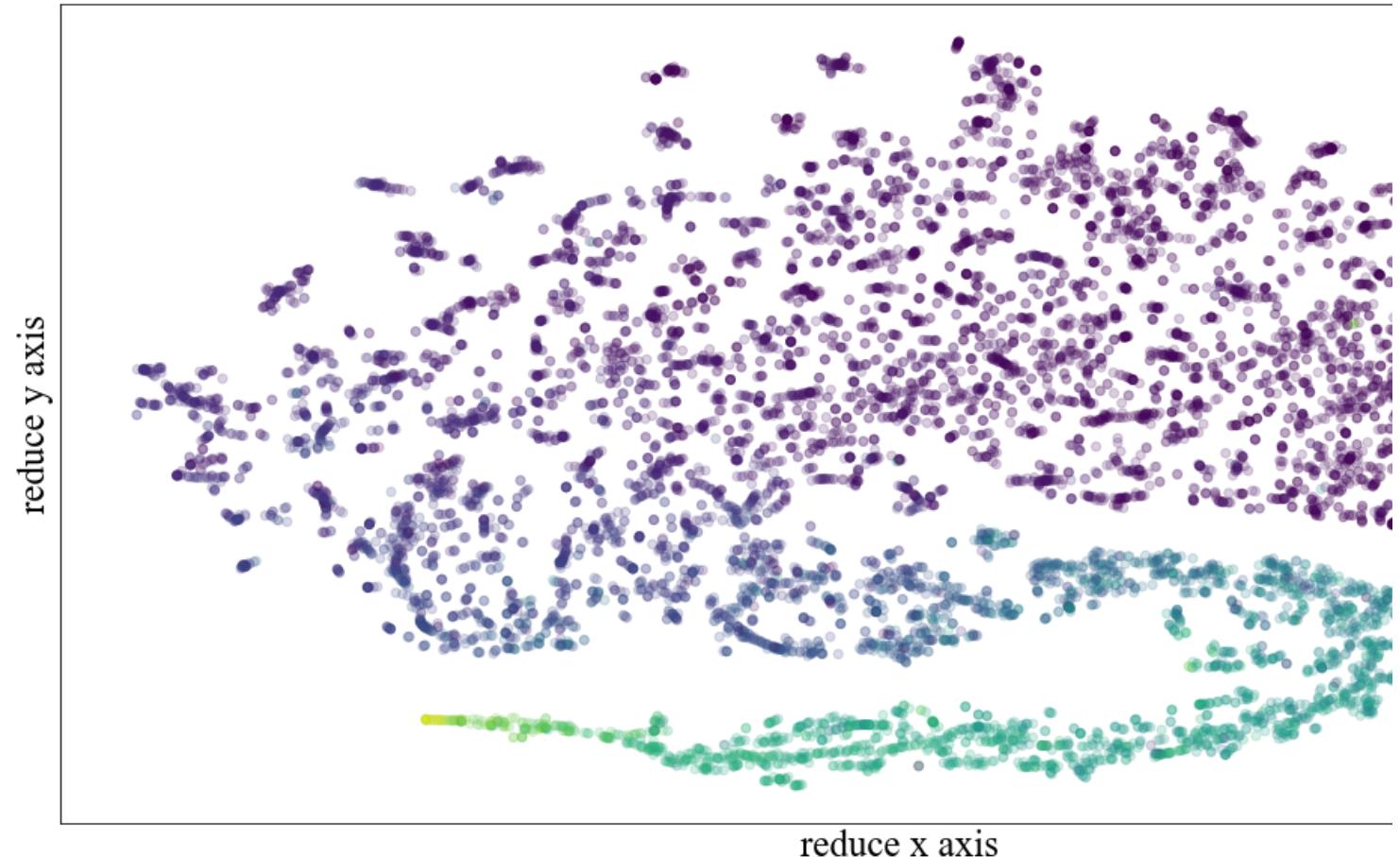
- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualizati
 - 1.3.6 ADC
- ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sensi
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 data
 - 1.4.2 Parameter optimi
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [127]:

```
1 from sklearn.manifold import TSNE
2 #ax = plt.axes(projection='3d')
3 x = y_pred
4 X_embedded = TSNE(n_components=2).fit_transform(x)
5 #X_embedded
6 plt.figure(figsize=(20, 10))
7 plt.scatter(X_embedded[:, 0], X_embedded[:, 1], alpha=0.2, c=close, label = "t-SNE")
8 #ax.scatter3D(X_embedded[:, 0], X_embedded[:, 1], X_embedded[:, 2], c=close)
9 plt.xlabel('reduce x axis', font4)
10 plt.ylabel('reduce y axis', font4)
11 plt.legend(loc='upper right', fontsize= 25)
12 plt.xticks([])
13 plt.yticks([])
14 plt.savefig("real1T-SNE.png", dpi=300, bbox_inches="tight")
15 plt.show()
16 plt.show()
```

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE



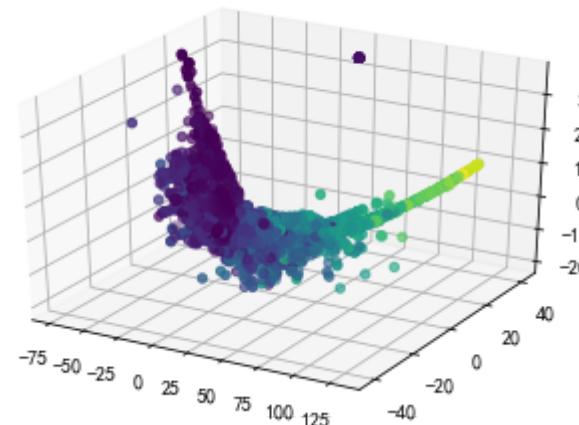
1.6.3 Isomap

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature selection
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to select the best model
 - 1.3.2 Select the appropriate classifier
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
 - ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
 - ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
 - ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

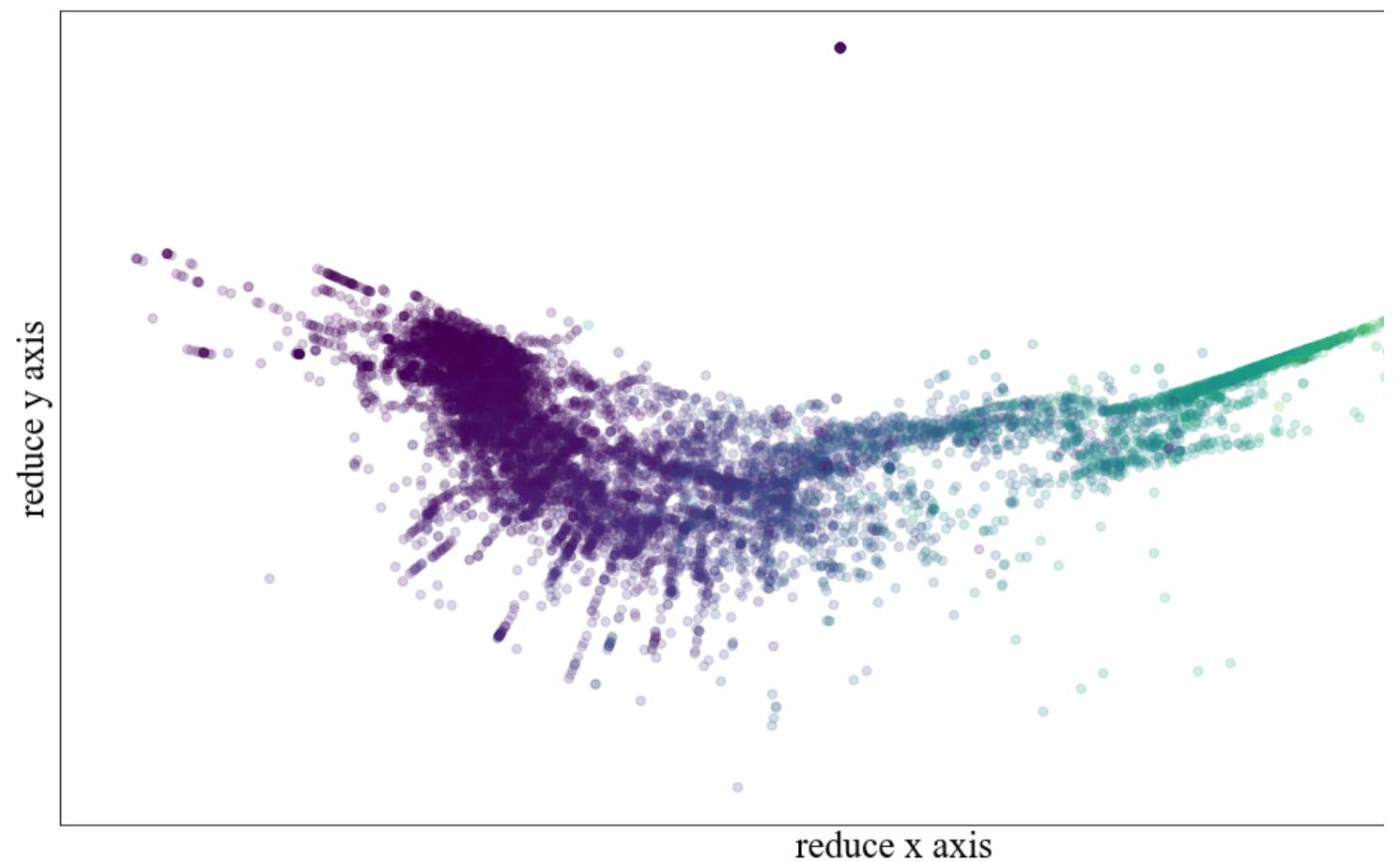
In [129]:

```
1 from sklearn import manifold
2 ax = plt.axes(projection='3d')
3 x = y_pred
4 X_embedded = manifold.Isomap(n_components=3).fit_transform(x)
5 import seaborn as sns
6 from matplotlib import pyplot as plt
7 plt.figure(figsize=(20, 10))
8 plt.scatter(X_embedded[:, 0], X_embedded[:, 1], alpha=0.2, c=close, label="Isomap")
9 ax.scatter3D(X_embedded[:, 0], X_embedded[:, 1], X_embedded[:, 2], c=close)
10 plt.xlabel('reduce x axis', font4)
11 plt.ylabel('reduce y axis', font4)
12 plt.legend(loc='upper right', fontsize= 25)
13 plt.xticks([])
14 plt.yticks([])
15 plt.savefig("reallIsomap.png", dpi=300, bbox_inches="tight")
16 plt.show()
```



Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 descriptors
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE



1.6.4 MDS

Contents ⚙️

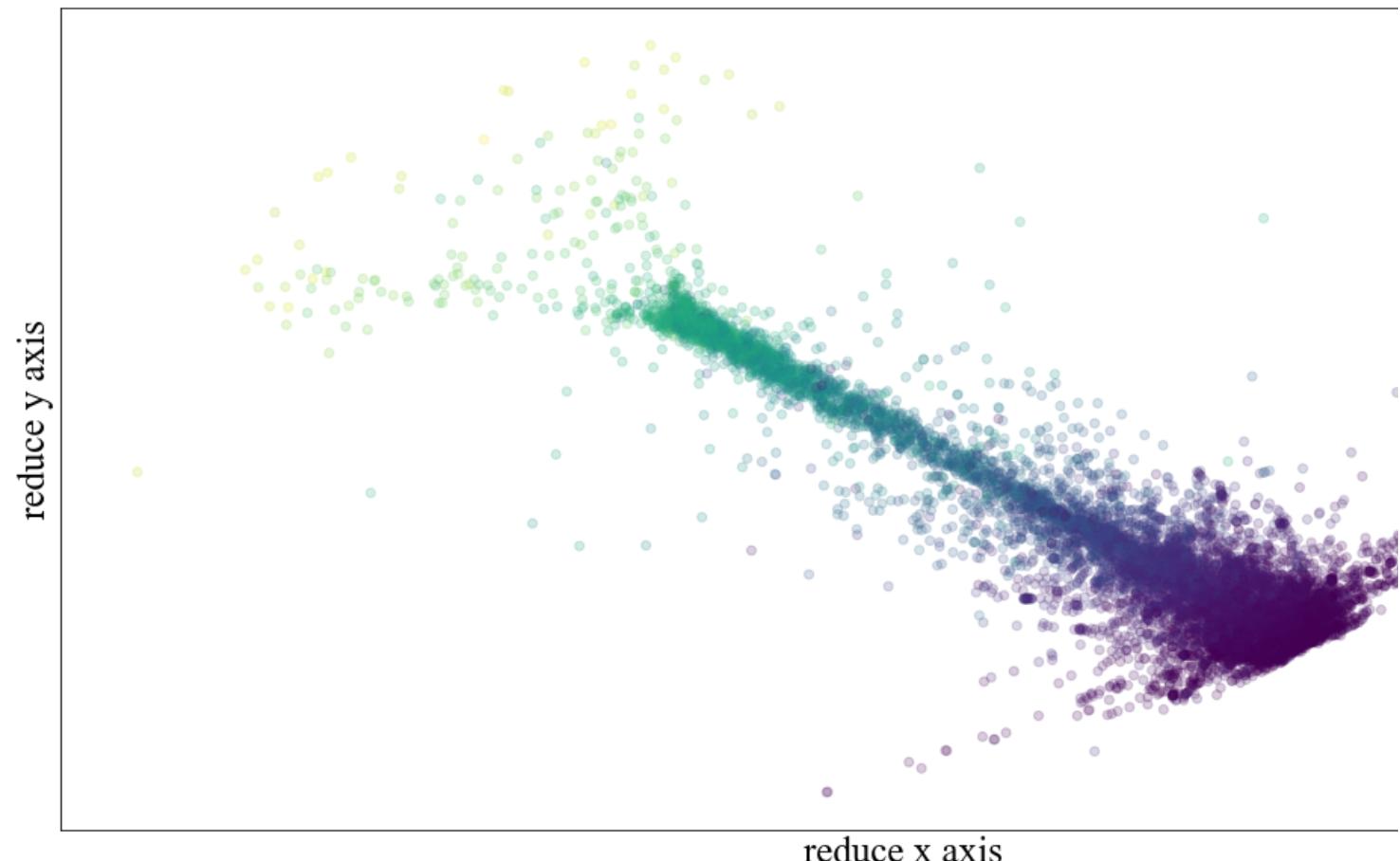
- 1.1 Exploratory analysis
- 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - 1.3.5 DT
 - 1.3.5.1 Visualizations
 - 1.3.6 ADC
 - 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- 1.6 Deep neural networks
 - 1.6.1 PCA
 - 1.6.2 T-SNE

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to
 - 1.3.2 Select the appro
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The thresh
 - 1.3.7.5 Cost sens
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 de
 - 1.4.2 Parameter optim
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput
 - 1.5.1 Virtual high-throughput
 - 1.5.2 DNN prediction
- ▼ 1.6 Deep neural network
 - 1.6.1 PCA
 - 1.6.2 T-SNE

In [130]:

```
1 from sklearn import manifold
2
3 x = y_pred
4 X_embedded = manifold.MDS(n_components = 2, max_iter=100, n_init=1).fit_transform(x)
5 plt.figure(figsize=(20, 10))
6 plt.scatter(X_embedded[:, 0], X_embedded[:, 1], alpha=0.2 ,c =close, label="MDS")
7 plt.xlabel('reduce x axis', font4)
8 plt.ylabel('reduce y axis', font4)
9 plt.legend(loc='upper right', fontsize= 25)
10 plt.xticks([])
11 plt.yticks([])
12 plt.savefig("real1MDS.png", dpi=300, bbox_inches="tight")
13 plt.show()
```



END

Contents ⚙

- 1.1 Exploratory analysis
- ▼ 1.2 Descriptors population
 - 1.2.1 The descriptor
 - 1.2.2 The feature space
- ▼ 1.3 Classification models
 - 1.3.1 Use pycaret to build a model
 - 1.3.2 Select the appropriate model
 - 1.3.3 RFC
 - 1.3.4 KNN
 - ▼ 1.3.5 DT
 - 1.3.5.1 Visualization
 - 1.3.6 ADC
 - ▼ 1.3.7 RFC : Analysis
 - 1.3.7.1 F1 score
 - 1.3.7.2 confusion matrix
 - 1.3.7.3 Accuracy
 - 1.3.7.4 The threshold
 - 1.3.7.5 Cost sensitivity
- ▼ 1.4 Train the deep neural network
 - 1.4.1 Use S1, S2 datasets
 - 1.4.2 Parameter optimization
 - 1.4.3 The trained model
 - 1.4.4 Residual analysis
- ▼ 1.5 Virtual high throughput screening
 - 1.5.1 Virtual high-throughput screening
 - 1.5.2 DNN predictions
- ▼ 1.6 Deep neural network analysis
 - 1.6.1 PCA
 - 1.6.2 T-SNE