

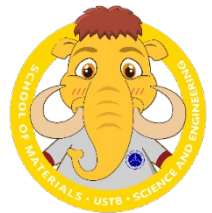


◆ 上机操作





3. 机器学习算法



◆ 小练习

```
[2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib
import matplotlib.pyplot as plt

#读取Excel数据
filefullpath = r"G:/DATA/Paper One-ForVisual.xlsx"
df = pd.read_excel(filefullpath)
#读取数据特征
df.describe()
```

```
[2]:
```

	ELEMENTS	T_MAX	T_MIN	T_AVE	RH_MIN	RH_MAX
count	522.000000	522.000000	522.000000	522.000000	522.000000	522.000000
mean	2.367471	34.321839	-0.923563	17.961494	39.214943	77.149430
std	1.991431	1.743940	6.744028	4.199953	11.876820	19.914310
min	0.030000	31.300000	-8.000000	14.200000	16.800000	72.800000



3. 机器学习算法



```
[4]: #检查重复值
print('未去重: ', df.shape)
print('去重: ', df.drop_duplicates().shape)
```

```
未去重: (522, 17)
去重: (511, 17)
```

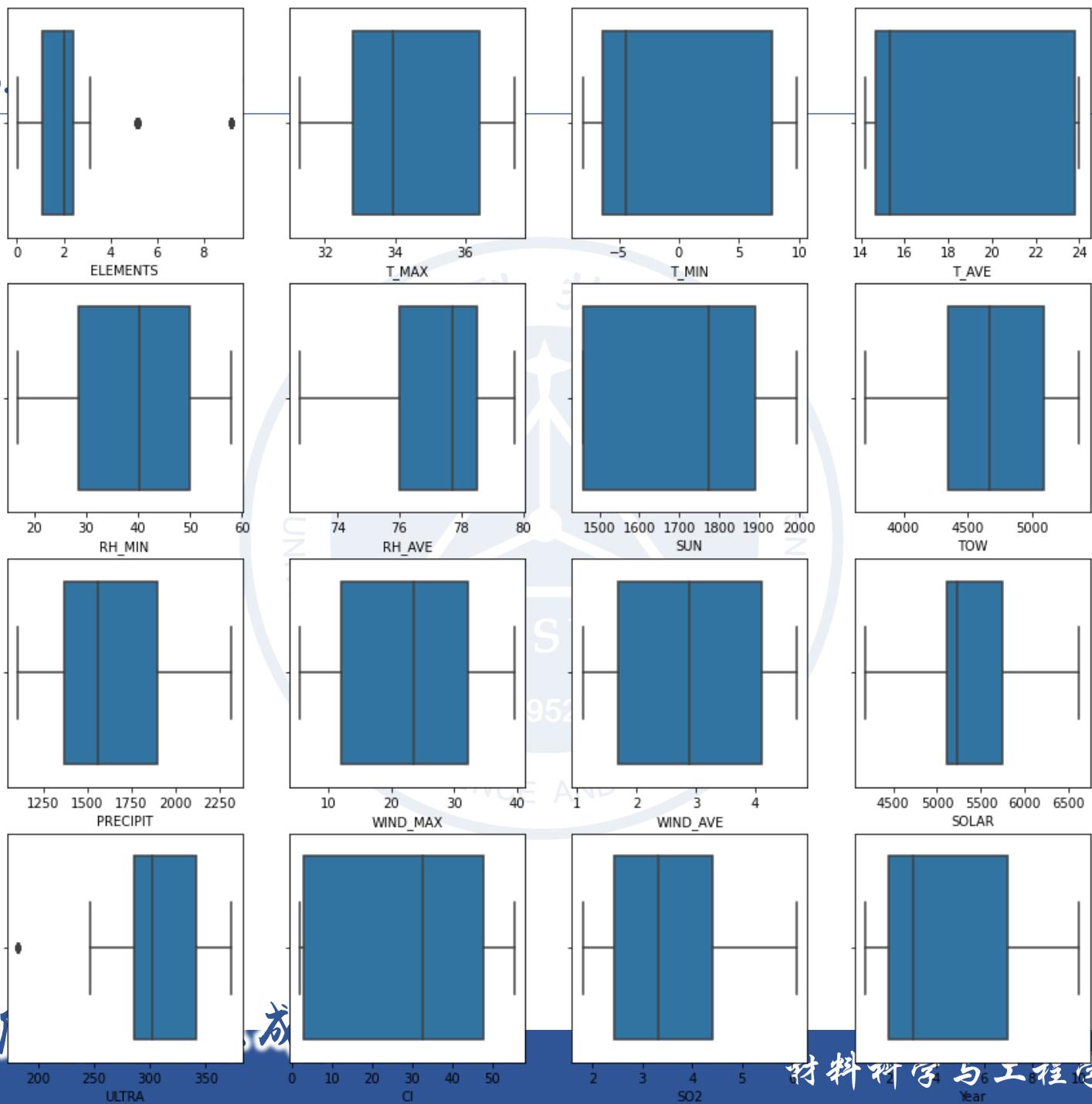
```
[5]: df.columns
```

```
[5]: Index(['ELEMENTS', 'T_MAX', 'T_MIN', 'T_AVE', 'RH_MIN', 'RH_AVE', 'SUN', 'TOW',
          'PRECIPIT', 'WIND_MAX', 'WIND_AVE', 'SOLAR', 'ULTRA', 'Cl', 'SO2',
          'Year', 'Vcorr'],
          dtype='object')
```

```
[6]: #检查异常值
#箱线图
fig, axes = plt.subplots(nrows=4, ncols=4, figsize=(15, 15))
#绘制箱线图
sns.boxplot(x="ELEMENTS", data=df, ax=axes[0][0])
sns.boxplot(x='T_MAX', data=df, ax=axes[0][1])
sns.boxplot(x='T_MIN', data=df, ax=axes[0][2])
sns.boxplot(x='T_AVE', data=df, ax=axes[0][3])
sns.boxplot(x='RH_MIN', data=df, ax=axes[1][0])
sns.boxplot(x="RH_AVE", data=df, ax=axes[1][1])
sns.boxplot(x='SUN', data=df, ax=axes[1][2])
sns.boxplot(x='TOW', data=df, ax=axes[1][3])
sns.boxplot(x='PRECIPIT', data=df, ax=axes[2][0])
sns.boxplot(x='WIND_MAX', data=df, ax=axes[2][1])
sns.boxplot(x='WIND_AVE', data=df, ax=axes[2][2])
sns.boxplot(x='SOLAR', data=df, ax=axes[2][3])
sns.boxplot(x='ULTRA', data=df, ax=axes[3][0])
sns.boxplot(x='Cl', data=df, ax=axes[3][1])
sns.boxplot(x='SO2', data=df, ax=axes[3][2])
sns.boxplot(x='Year', data=df, ax=axes[3][3])
plt.show()
```



3.



学厚

成

材料科学与工程学院 4





3. 机器学习算法

◆ 另一种箱线图

```
7]: df1 = df.iloc[:, 0:8].copy()
df2 = df.iloc[:, 8:17].copy()
#划分
fig, axes = plt.subplots(1,8,figsize=(16,8)) #这个函数是设置子图参数的
color = dict(boxes='DarkGreen', whiskers='DarkOrange', medians='DarkBlue', caps='Red')
# boxes表示箱体, whisker表示触须线
# medians表示中位数, caps表示最大与最小值界限
df1.plot(kind='box',ax=axes, subplots=True, title='Different boxplots', color=color, sym='r+')
# sym参数表示异常值标记的方式
axes[0].set_ylabel('%')
axes[1].set_ylabel('°C')
axes[2].set_ylabel('°C')
axes[3].set_ylabel('°C')
axes[4].set_ylabel('%')
axes[5].set_ylabel('%')
axes[6].set_ylabel('h')
axes[7].set_ylabel('h')
fig.subplots_adjust(wspace=1,hspace=1) # 调整子图之间的间距
plt.show()
```

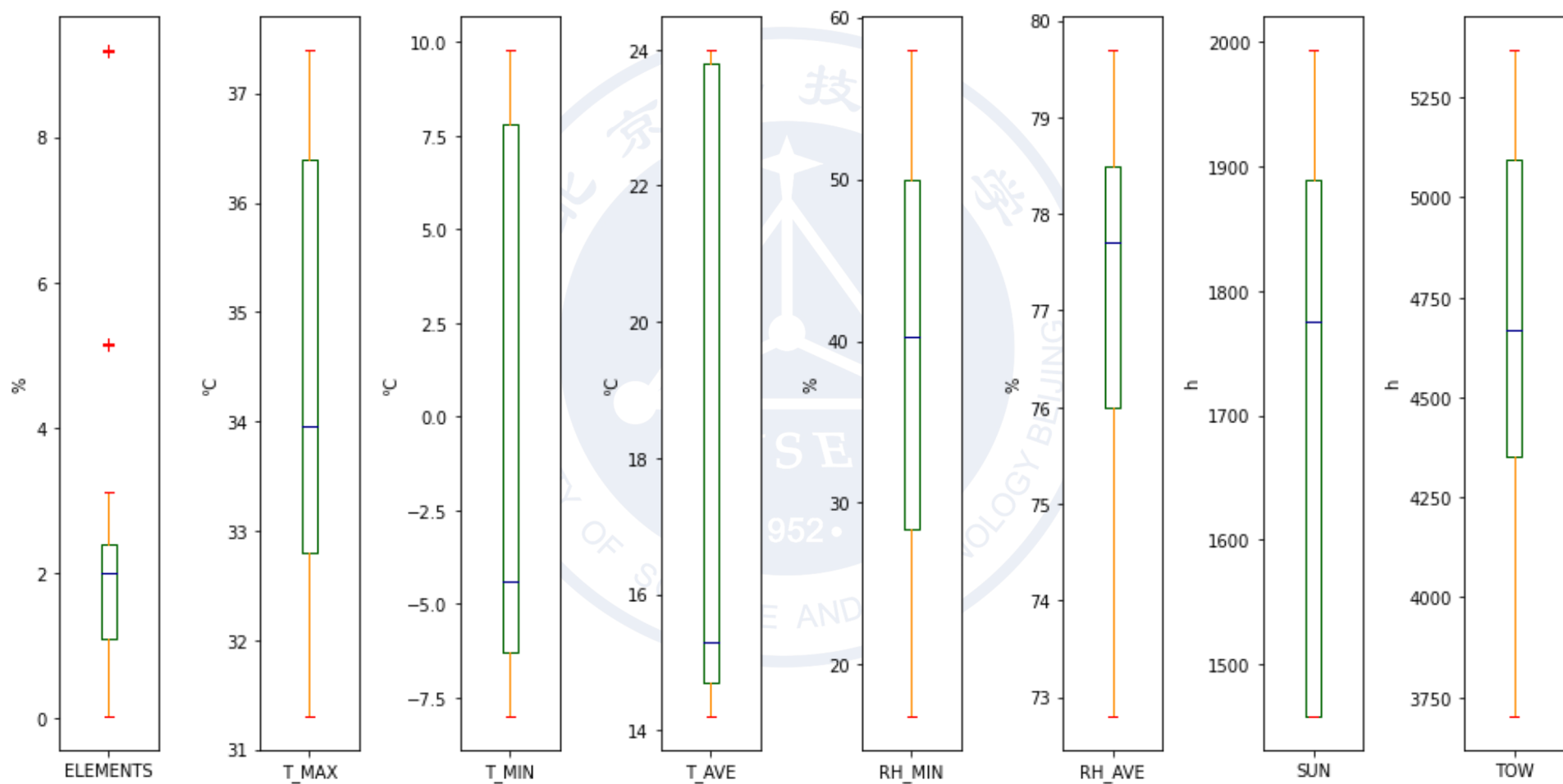




3. 机器学习算法



Different boxplots





◆绘制出变量 ‘PRECIPIT’, ‘WIND_MAX’,
‘WIND_AVE’, ‘SOLAR’, ‘ULTRA’, ‘CI’,
‘SO2’, ‘Year’, ‘Vcorr’的箱线图

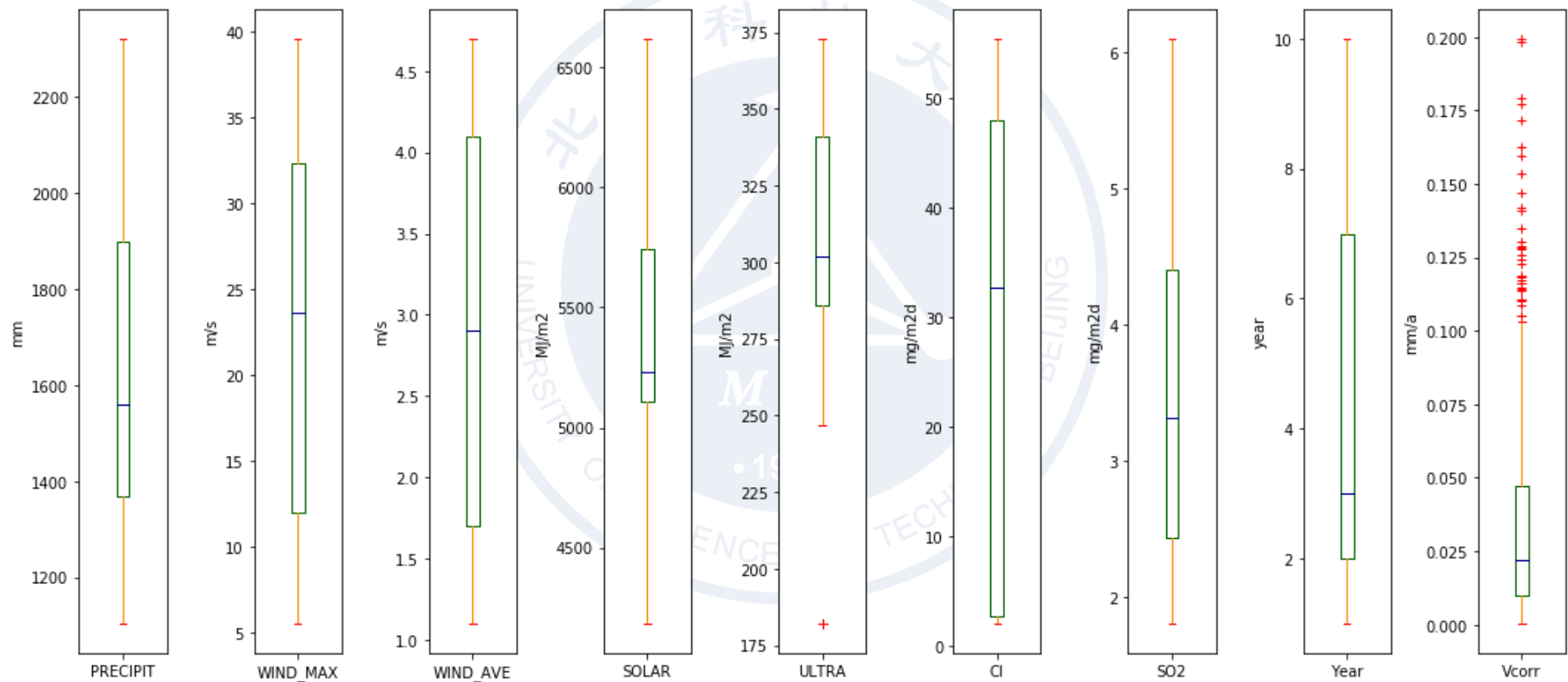




3. 机器学习算法



Different boxplots





3. 机器学习算法



◆ 上图的代码

```
: fig, axes = plt.subplots(1,9,figsize=(18,8)) #这个函数是设置子图参数的
color = dict(boxes='DarkGreen', whiskers='DarkOrange', medians='DarkBlue', caps='Red')
# boxes表示箱体, whisker表示触须线
# medians表示中位数, caps表示最大与最小值界限
df2.plot(kind='box',ax=axes, subplots=True, title='Different boxplots', color=color, sym='r+')
# sym参数表示异常值标记的方式
axes[0].set_ylabel('mm')
axes[1].set_ylabel('m/s')
axes[2].set_ylabel('m/s')
axes[3].set_ylabel('MJ/m2')
axes[4].set_ylabel('MJ/m2')
axes[5].set_ylabel('mg/m2d')
axes[6].set_ylabel('mg/m2d')
axes[7].set_ylabel('year')
axes[8].set_ylabel('mm/a')
fig.subplots_adjust(wspace=1,hspace=1) # 调整子图之间的间距
plt.show()
```





3. 机器学习算法

◆ 数据分布的展示分析

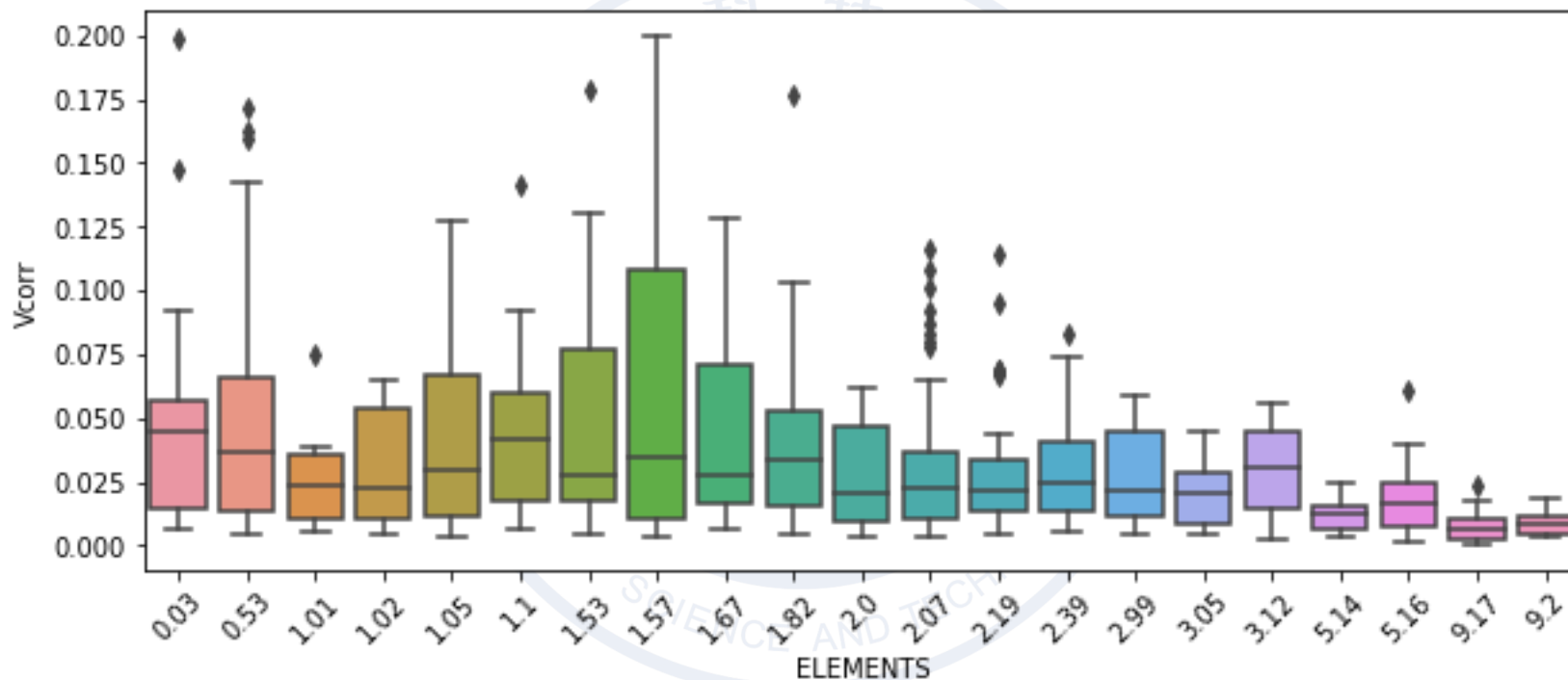
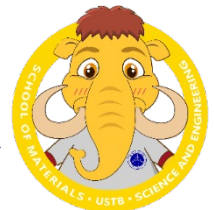
可视化分析

```
[8]: #1) 合金化程度与腐蚀速率
plt.figure(figsize=(10, 4))
plt.xticks(rotation=45)
sns.boxplot(x='ELEMENTS', y='Vcorr', data=df)
plt.show()
```



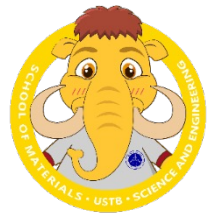


3. 机器学习算法



学厚质朴 百炼成材





◆绘制出变量 ‘T_MAX’, ‘T_MIN’, ‘T_AVE’, ‘RH_MIN’,
‘RH_AVE’, ‘SUN’, ‘TOW’, ‘PRECIPIT’,
‘WIND_MAX’, ‘WIND_AVE’, ‘SOLAR’, ‘ULTRA’,
‘Cl’, ‘SO2’, ‘Year’, ‘Vcorr’ 的数值分布图



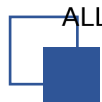


3. 机器学习算法

◆ 特征工程

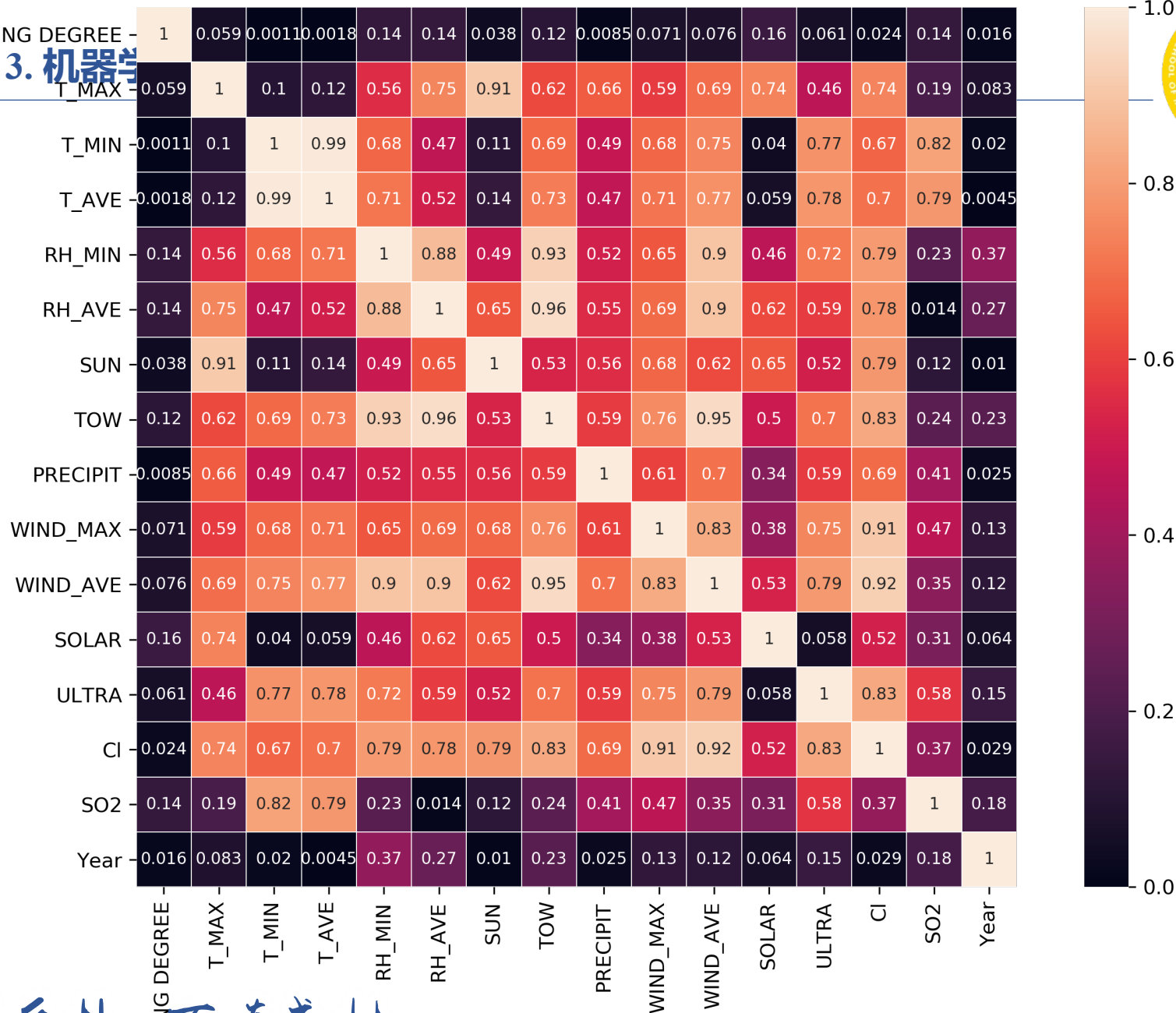
```
#归一化
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
Data_Nor = sc.fit_transform(df)
mater_Data2 = Data_Nor[:, 0:16].copy() #除去了Vcorr那一列
mater_Data2 = pd.DataFrame(mater_Data2) #归一化后矩阵变成np.array格式, 不能做.corr(), 所以要转换成dataframe
mater_Data2.columns = ['ALLOYING DEGREE', 'T_MAX', 'T_MIN', 'T_AVE', 'RH_MIN', 'RH_AVE', 'SUN',
                      'TOW', 'PRECIPIT', 'WIND_MAX', 'WIND_AVE', 'SOLAR', 'ULTRA', 'Cl',
                      'S02', 'Year'] #转换后的dataframe列名称都丢了, 重新赋值, 否则图里的特征名称都是0.1.2这样
correlation2 = mater_Data2.corr(method='pearson')
correlation2_abs = correlation2.abs()
plt.figure(figsize=(12, 8), dpi=300)
#绘制热力图
sns.heatmap(correlation2_abs, linewidths=0.2, vmax=1, vmin=0, linecolor='w',
            annot=True, annot_kws={'size': 8}, square=True)
```





ALLOYING DEGREE

3. 机器学习

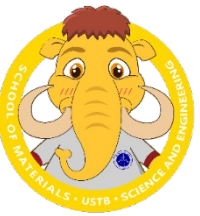


学厚质朴 百炼成材





3. 机器学习算法

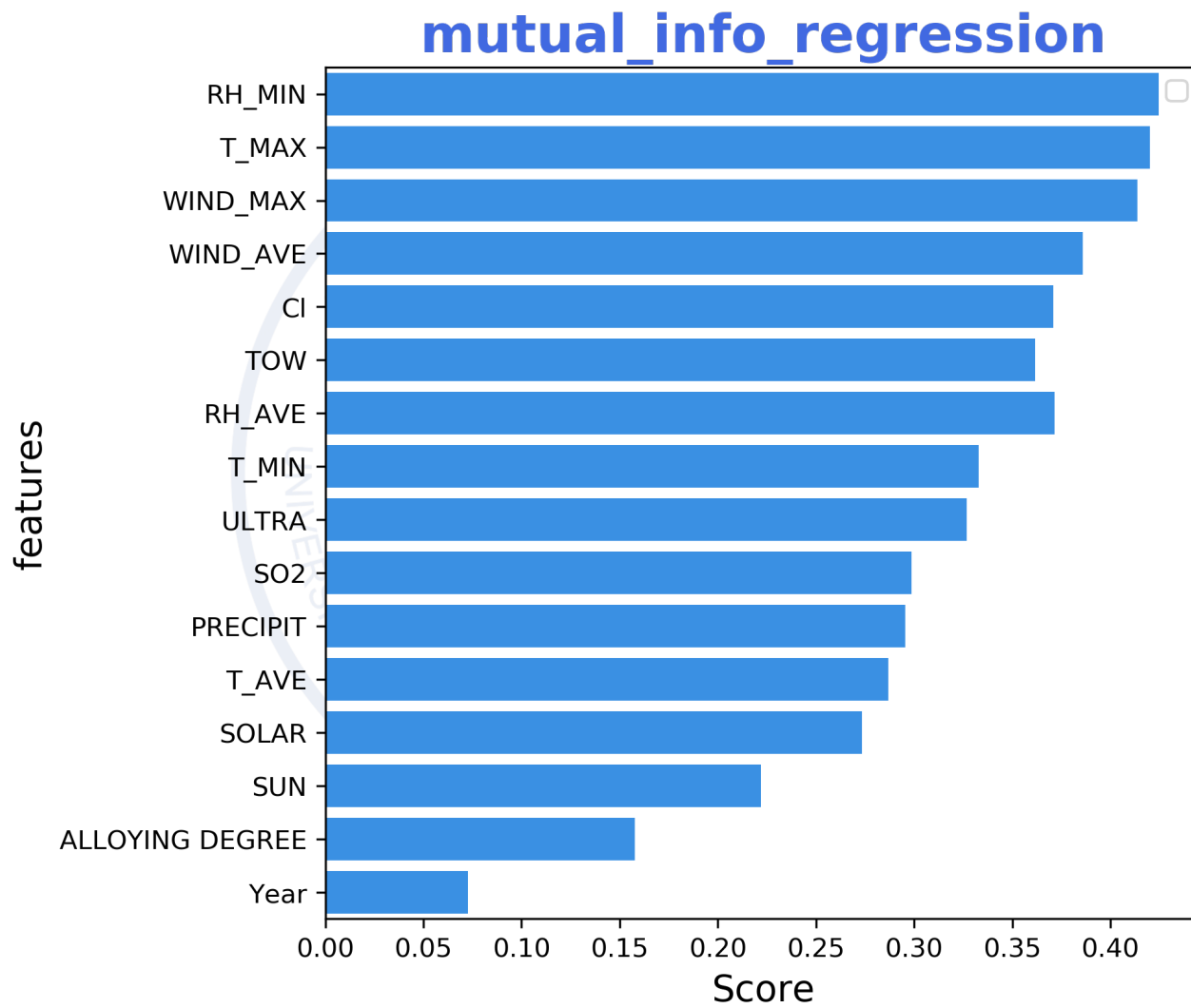
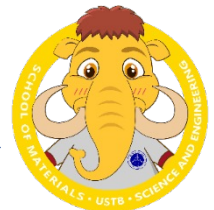


```
from sklearn.feature_selection import mutual_info_regression
features = mater_Data2
target = Data_Nor[:, 16].copy()
score = mutual_info_regression(features, target)
X = ['ALLOYING DEGREE', 'T_MAX', 'T_MIN', 'T_AVE', 'RH_MIN', 'RH_AVE', 'SUN',
      'TOW', 'PRECIPIT', 'WIND_MAX', 'WIND_AVE', 'SOLAR', 'ULTRA', 'Cl',
      'S02', 'Year']
y = score
plt.figure(figsize=(6, 6), dpi=300)
sns.barplot(y, X, color="dodgerblue", order=[ 'RH_MIN', 'T_MAX', 'WIND_MAX', 'WIND_AVE',
      'S02', 'PRECIPIT', 'T_AVE', 'SOLAR', 'SUN', 'ALLOYING DEGREE', 'Year'])
plt.xlabel('Score', fontsize=14)
plt.ylabel('features', fontsize=14)
plt.legend()
plt.title('mutual_info_regression', fontsize=20, fontweight='bold', color='royalblue')
plt.show()
```



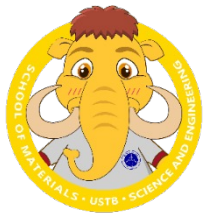


3. 机器学习算法





3. 机器学习算法

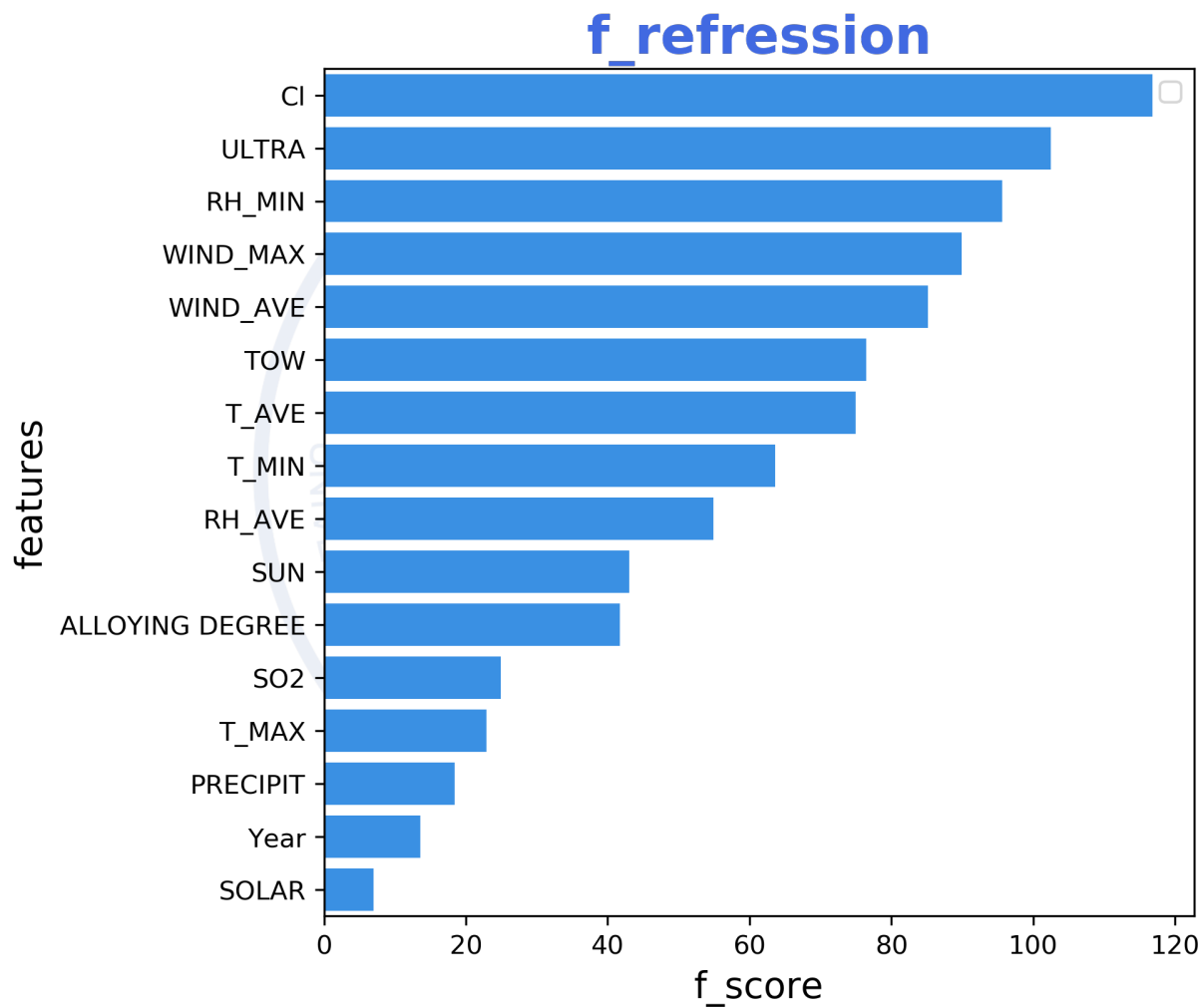
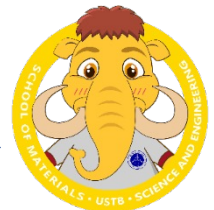


```
from sklearn.feature_selection import f_regression
score2 = f_regression(features, target)[0]
#不加最后的【0】，则会返回包含连个阵列的数值，前一个是各特征的F值，后一个是p值。
X = ['ALLOYING DEGREE', 'T_MAX', 'T_MIN', 'T_AVE', 'RH_MIN', 'RH_AVE', 'SUN',
      'TOW', 'PRECIPIT', 'WIND_MAX', 'WIND_AVE', 'SOLAR', 'ULTRA', 'Cl',
      'SO2', 'Year']
y = score2
plt.figure(figsize=(6, 6), dpi=300)
sns.barplot(y, X, color="dodgerblue", order=[ 'Cl', 'ULTRA', 'RH_MIN', 'WIND_MAX',
      'SUN', 'ALLOYING DEGREE', 'SO2', 'T_MAX', 'PRECIPIT', 'Year', 'SOLAR'])
plt.xlabel('f_score', fontsize=14)
plt.ylabel('features', fontsize=14)
plt.legend()
plt.title('f_refression', fontsize=20, fontweight='bold', color='royalblue')
plt.show()
```





3. 机器学习算法





3. 机器学习算法



```
X1 = df.iloc[:, 0:15].copy()
y1 = df.iloc[:, 15].copy()

from sklearn.ensemble import RandomForestRegressor
model = RandomForestRegressor(random_state=10)

def RFFeaImpor_(X_data,y_data):
    model.fit(X_data,y_data)
    result_ = {'var':X_data.columns.values
               , 'importances_':model.feature_importances_}
    feature_importances_ = pd.DataFrame(result_, columns=['var','importances_'], index=X_data.columns.values)
    return feature_importances_

feature_importances_ = RFFeaImpor_(X1,y1)
print (feature_importances_)
```





3. 机器学习算法



	var	importances_
ELEMENTS	ELEMENTS	0.419420
T_MAX	T_MAX	0.025188
T_MIN	T_MIN	0.014046
T_AVE	T_AVE	0.043115
RH_MIN	RH_MIN	0.063493
RH_AVE	RH_AVE	0.000468
SUN	SUN	0.062548
TOW	TOW	0.209315
PRECIPIT	PRECIPIT	0.000017
WIND_MAX	WIND_MAX	0.025196
WIND_AVE	WIND_AVE	0.011529
SOLAR	SOLAR	0.000352
ULTRA	ULTRA	0.099836
Cl	Cl	0.012945
S02	S02	0.012533

