



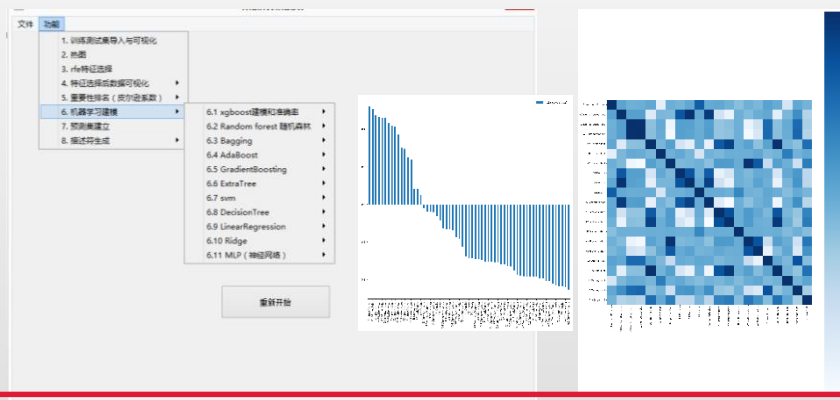
# NJmat软件

# NJmatML: 材料科学机器学习软件

## Python版



## windows版



## 网页版

ZL Lab Nanjing

Machine Learning Data Import Feature Engineering Machine Learning models Symbolic Regression Download Report

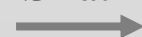
- pip install NJmatML
- <https://pypi.org/project/NJmatML/>
- publicly available
- 用户友好: 降低pandas, matplotlib, sklearn, matn
- 功能: 文件导入、可视化、机器学习建模、特征生成 (重要性排名)、特征选择、准确率计算、遗传算法、
- 说明书: <https://github.com/Zhang-NJ-Lab/NJmatML/blob/main/2022-11-21/NJ>

- NJmatML的windows版本
- 下载地址: <https://figshare.com/articles/software/NJmatML/24607893>
- doi: <https://doi.org/10.6084/m9.figshare.24607893.v1.exe>



NJmatML

实验数据 (csv)



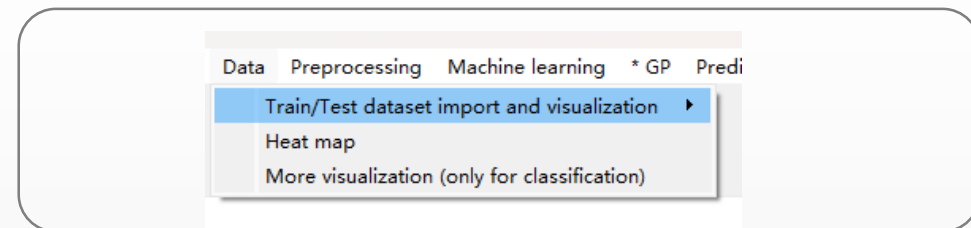
材料预测与实验验证

# NJmatML: 材料科学机器学习软件

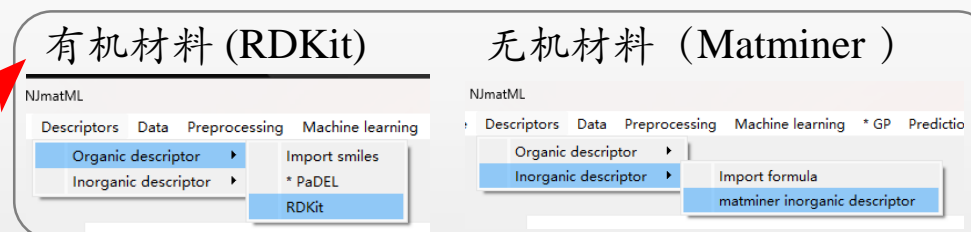
NJmatML主界面



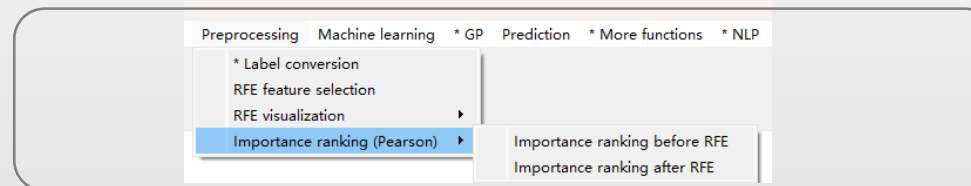
训练测试集导入



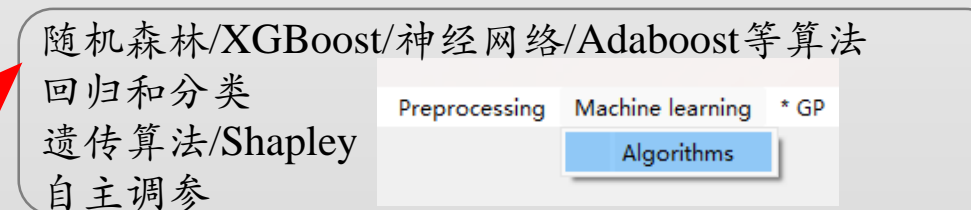
描述符辅助生成:



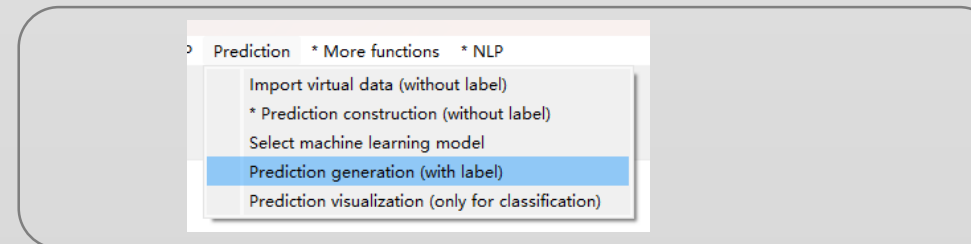
特征选择与可视化



机器学习建模



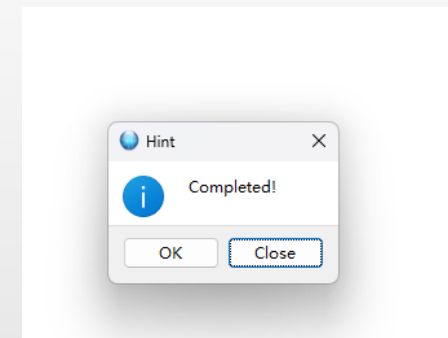
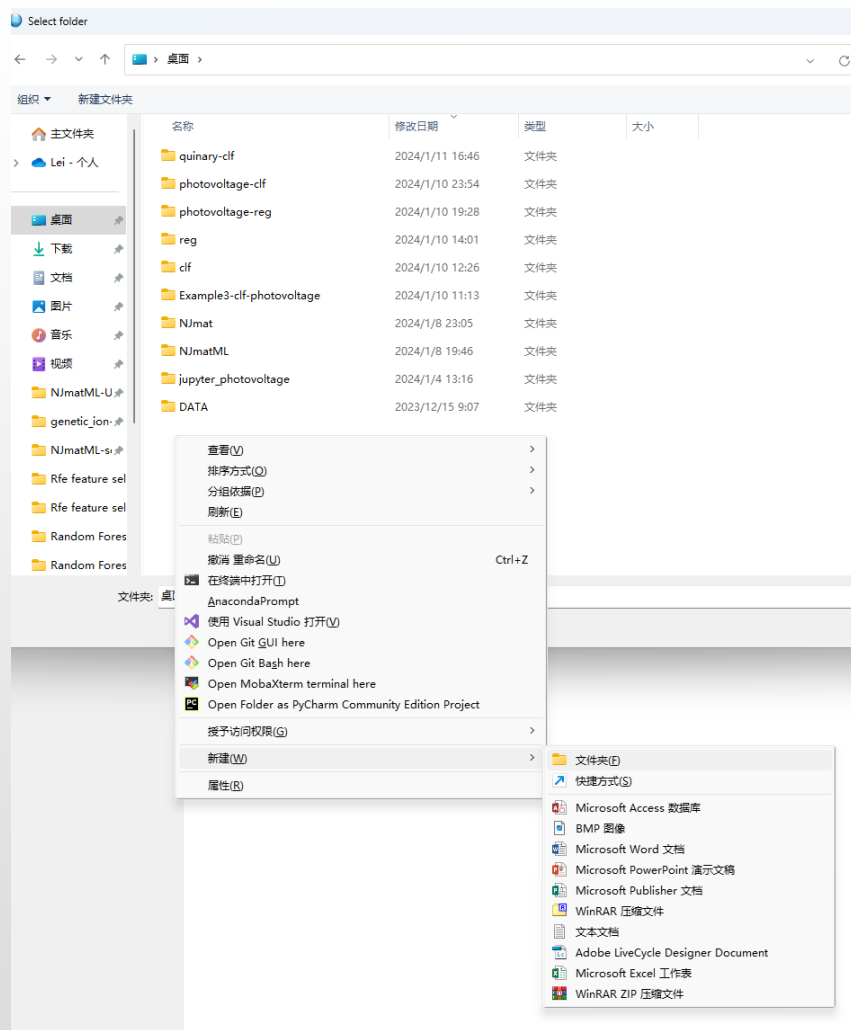
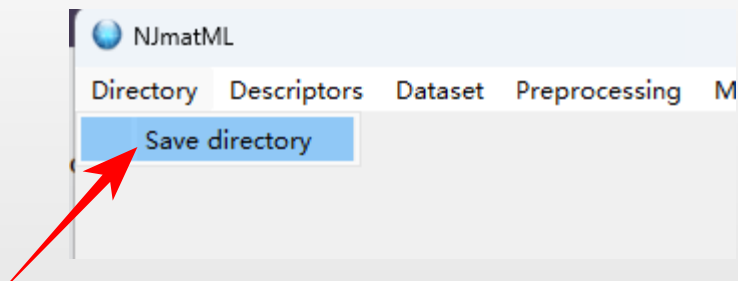
机器学习预测





# 1 选择保存路径

# 1 选择项目存储文件夹（自定义）



# 1 选择项目存储文件夹（自定义）

| 面 > reg >   |   |         |            |
|---|---|---------|------------|
|                            |  | ↑↓ 排序 ▾ | ≡ 查看 ▾ ... |
| 名称  | 修改日期  | 类型      | 大小         |
|  Data importing            | 2024/1/12 1:13  | 文件夹     |            |
|  Descriptor generation     | 2024/1/12 0:52  | 文件夹     |            |
|  GP                        | 2024/1/10 14:01   | 文件夹     |            |
|  Machine Learning Modeling | 2024/1/10 13:37   | 文件夹     |            |
|  Prediction                | 2024/1/10 13:37   | 文件夹     |            |
|  Preprocessing             | 2024/1/10 13:37   | 文件夹     |            |
|  Shapley                  | 2024/1/10 13:37   | 文件夹     |            |

数据集导入  
描述符生成  
遗传算法建模  
传统机器学习建模（分类和回归）  
预测虚拟空间  
预处理（特征选择等）  
Shapley描述符分析

后续在自定义文件夹中生成新的数据文件夹



## 2 材料特征化 Featurizer

# 2.1 Featurizer 材料特征化按钮（无机和有机）

## 2in1:有机无机2合1

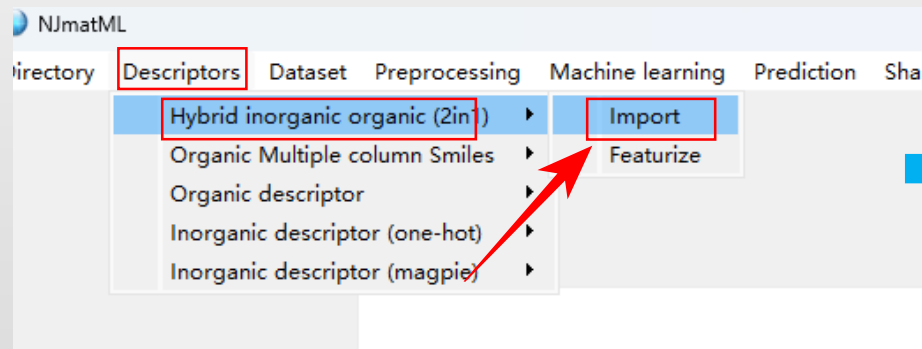
Smiles和化学式都有

特征化按钮（训练测试集辅助生成）

有机RDKit描述符，无机为Matminer magpie 描述符（元素、原子质量等等）

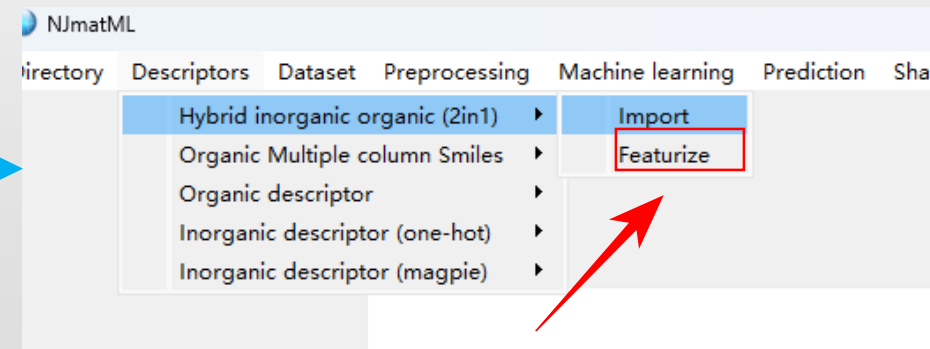
Descriptors: Hybrid inorganic organic (2in1) → Import → Featurize

### Import raw



| 1              | 2              | 3                            | 4                 | 5 |
|----------------|----------------|------------------------------|-------------------|---|
| layer1_Formula | layer2_Formula | SMILES_Smiles                | conservation rate |   |
| 1              | OsPbBr2        | [NH4+][Cl-]                  | 0.941727          |   |
| 2              | OsPbBr2        | CCCC[N+](CCCC)(CCCC)CCCC[I-] | 0.326007          |   |
| 3              | OsPbBr2        | [Br-][Br-][Sn+2]             | 0.257407          |   |
| 4              | OsPbBr2        | [Br-][Br-][I-][I-]           | 0.721138          |   |
| 5              | OsPbBr2        | [Sn+2][I-][I-]               | 0.143376          |   |
| 6              | OsPbBr2        | [Cl-][Ag+]                   | 0.132046          |   |
| 7              | OsPbBr2        | CH3NH3PbBr2                  | 0.173264          |   |
| 8              | OsPbCl2Br      | CH3NH3SnBr2                  | 0.375907          |   |
| 9              | OsPbCl2Br      | CH3NH3SnBr2                  | 0.188991          |   |
| 10             | OsPbCl2        | [Cl-][K+]                    | 0.185522          |   |
| 11             | OsPbCl2        | [NH4+][Cl-]                  | 0.137411          |   |
| 12             | OsPbCl2        | [I-][I-][Pb+2]               | 0.528662          |   |
| 13             | OsPbCl2        | CCCC[N+](CCCC)(CCCC)CCCC[I-] | 0.359603          |   |
| 14             | OsPbCl2        | [Br-][Br-][Pb+2]             | 0.287676          |   |
| 15             | OsPbCl2        | [Br-][Br-][Sn+2]             | 0.127879          |   |
| 16             | OsPbCl2        | [I-][I-][Cl-]                | 0.528662          |   |
| 17             | OsPbCl2        | CH3NH3PbCl3                  | 0.265946          |   |
| 18             | OsSnBr2        | [NH4+][Cl-]                  | 0.9163            |   |
| 19             | OsSnBr2        | (NH2)2CHSnBr2                | 1.306513          |   |
| 20             | OsSnBr2        | (NH2)2CHSnBr2                | 0.465672          |   |
| 21             | OsSnBr2        | (NH2)2CHSnBr2                | 0.333333          |   |
| 22             | OsSnBr2        | CH3NH3SnCl2                  | 0.985401          |   |
| 23             | OsSnBr2        | CH3NH3SnCl2                  | 0.634021          |   |
| 24             | OsSnCl3        | (NH2)2CHSnCl2                | 0.879226          |   |
| 25             | OsSnCl3        | (NH2)2CHSnCl2                | 0.536913          |   |
| 26             | OsSnCl2        | (NH2)2CHSnCl2                | 0.628713          |   |
| 27             | OsSnCl2        | (NH2)2CHSnCl2                | 0.287739          |   |
| 28             | OsSnCl2        | (NH2)2CHSnCl3                | 0.232203          |   |
| 29             | OsSnCl2        | (NH2)2CHSnCl3                | 0.334286          |   |
| 30             | OsSnCl2        | CH3NH3PbCl2                  | 0.298056          |   |
| 31             | OsSnCl2        | CH3NH3PbCl2                  | 0.723404          |   |
| 32             | (NH2)2CHPbBr2  | CCCC[N+](CCCC)(CCCC)CCCC[I-] | 0.811321          |   |
| 33             | (NH2)2CHPbBr2  | [Cl-][Cl-][Pb+2]             | 0.573566          |   |
| 34             | (NH2)2CHPbBr2  | [Br-][Br-][I-][I-]           | 1.015625          |   |
| 35             | (NH2)2CHPbBr2  | OsSnBr2                      | 0.605159          |   |
| 36             | (NH2)2CHPbCl2  | OsSnBr2                      | 0.272121          |   |
| 37             | (NH2)2CHPbCl2  | OsPbCl2                      | 0.344214          |   |
| 38             | (NH2)2CHPbCl3  | CH3NH3PbBr2                  | 0.219089          |   |
| 39             | (NH2)2CHPbCl3  | CH3NH3PbBr2                  | 0.291083          |   |
| 40             | (NH2)2CHPbCl3  | CH3NH3PbCl3                  | 0.91573           |   |
| 41             | (NH2)2CHPbCl3  | CH3NH3PbCl3                  | 0.412811          |   |
| 42             | (NH2)2CHPbCl3  | CH3NH3SnCl2                  | 0.362069          |   |
| 43             | (NH2)2CHPbCl3  | CH3NH3SnCl2                  | 0.794304          |   |
| 44             | (NH2)2CHSnCl3  | CH3NH3SnCl2                  | 0.363465          |   |

### Featurize



选择起始原始raw数据集（只接受csv格式）



用户需要输入的csv格式:

无机: 列名以Formula结尾

有机: 列名以Smiles结尾

更准确:

化学式: 列名以Formula结尾

Smiles码: 列名以Smiles结尾

最右一列为目标数据

其它 (左侧) 为输入数据

拟合和预测交给软件!

输入变量

输出目标

|    | layer1_Formula | layer2_Formula | SMILES_Smiles                 | conservation rate |
|----|----------------|----------------|-------------------------------|-------------------|
| 1  | CsPbBr2I       | CsPbBr3        | [NH4+].[Cl-]                  | 0.041727          |
| 2  | CsPbBr2I       | CsPbBr3        | CCCC[N+](CCCC)(CCCC)CCCC.[I-] | 0.326007          |
| 3  | CsPbBr2I       | CsPbClI2       | [Br-].[Br-].[Sn+2]            | 0.257407          |
| 4  | CsPbBr2I       | CsPbClI2       | [Sn+2].[I-].[I-]              | 0.721138          |
| 5  | CsPbBr3        | CH3NH3PbBr2I   | [Cl-].[Ag+]                   | 0.143376          |
| 6  | CsPbBr3        | CH3NH3PbBr2I   | C[NH3+].[I-]                  | 0.132046          |
| 7  | CsPbCl2Br      | CH3NH3SnBrI2   | [Br-].[Ag+]                   | 0.173264          |
| 8  | CsPbCl2Br      | CH3NH3SnBrI2   | C(=N)N.[I-]                   | 0.375907          |
| 9  | CsPbCl2I       | CsPbCl2Br      | [Cl-].[K+]                    | 0.188991          |
| 10 | CsPbCl2I       | CsPbCl2Br      | [NH4+].[Cl-]                  | 0.185522          |
| 11 | CsPbI3         | CsPbBr2I       | [I-].[I-].[Pb+2]              | 0.137411          |
| 12 | CsPbI3         | CsPbBr2I       | CCCC[N+](CCCC)(CCCC)CCCC.[I-] | 0.506803          |
| 13 | CsPbI3         | CsPbI3         | [Br-].[Br-].[Pb+2]            | 0.287676          |
| 14 | CsPbI3         | CsPbI3         | [Br-].[Br-].[Sn+2]            | 0.127879          |
| 15 | CsPbI3         | CH3NH3PbI3     | [Li+].[Cl-]                   | 0.528662          |
| 16 | CsPbI3         | CH3NH3PbI3     | C[NH3+].[I-]                  | 0.285946          |
| 17 | CsSnBrI2       | (NH2)2CHSnBr2I | C(=N)N.Cl                     | 0.9163            |
| 18 | CsSnBrI2       | (NH2)2CHSnBr2I | CN.Cl                         | 1.306513          |
| 19 | CsSnBrI2       | (NH2)2CHSnBr2I | [Br-].[Br-].[Pb+2]            | 0.465672          |
| 20 | CsSnBrI2       | (NH2)2CHSnBr2I | [Cl-].[Cl-].[Pb+2]            | 0.333333          |
| 21 | CsSnBrI2       | CH3NH3SnClI2   | [Na+].[Cl-]                   | 0.985401          |
| 22 | CsSnBrI2       | CH3NH3SnClI2   | [Cl-].[Cl-].[Pb+2]            | 0.634021          |
| 23 | CsSnCl3        | (NH2)2CHSnClF2 | [Cl-].[K+]                    | 0.873926          |
| 24 | CsSnCl3        | (NH2)2CHSnClF2 | [Li+].[Cl-]                   | 0.536913          |
| 25 | CsSnClBr2      | (NH2)2CHSnCl2I | [F-].[F-].[F-].[Al+3]         | 0.628713          |
| 26 | CsSnClBr2      | (NH2)2CHSnCl2I | [Cl-].[Cl-].[Cu+2]            | 0.287739          |
| 27 | CsSnClBr2      | (NH2)2CHSnCl3  | [Br-].[Br-].[Pb+2]            | 0.232203          |
| 28 | CsSnClBr2      | (NH2)2CHSnCl3  | [I-].[I-].[Pb+2]              | 0.334286          |
| 29 | CsSnClF2       | CH3NH3PbCl2I   | [NH4+].[Cl-]                  | 0.298056          |
| 30 | CsSnClF2       | CH3NH3PbCl2I   | CCCC[N+](CCCC)(CCCC)CCCC.[I-] | 0.723404          |
| 31 | (NH2)2CHPbBr2I | CsPbI3         | [Cl-].[Cl-].[Pb+2]            | 0.811321          |
| 32 | (NH2)2CHPbBr2I | CsPbI3         | [Sn+2].[I-].[I-]              | 0.573566          |
| 33 | (NH2)2CHPbBr2I | CsSnBrI2       | C(=N)N.Cl                     | 1.015625          |
| 34 | (NH2)2CHPbBr2I | CsSnBrI2       | C(=N)N.[I-]                   | 0.605159          |
| 35 | (NH2)2CHPbCl2I | CsPbCl2I       | CN.Br                         | 0.272121          |
| 36 | (NH2)2CHPbCl2I | CsPbCl2I       | [Sn+2].[I-].[I-]              | 0.344214          |
| 37 | (NH2)2CHPbI3   | CH3NH3PbBr2I   | [Br-].[Br-].[Sn+2]            | 0.219089          |
| 38 | (NH2)2CHPbI3   | CH3NH3PbBr2I   | [Sn+2].[I-].[I-]              | 0.291083          |
| 39 | (NH2)2CHPbI3   | CH3NH3PbI3     | CN.Br                         | 0.91573           |
| 40 | (NH2)2CHPbI3   | CH3NH3PbI3     | CCCC[N+](CCCC)(CCCC)CCCC.[I-] | 0.412811          |
| 41 | (NH2)2CHPbI3   | CH3NH3SnCl2Br  | C(=N)N.Cl                     | 0.362069          |
| 42 | (NH2)2CHPbI3   | CH3NH3SnCl2Br  | C(=N)N.[I-]                   | 0.794304          |
| 43 | (NH2)2CHSnCl3  | CH3NH3SnClF2   | [F-].[F-].[F-].[Al+3]         | 0.363465          |

# Formula Formula Smiles T, t, C Output

需要特征化的变量，例如：

第一列：化学式

第二列：化学式

第三列：Smiles码

....

需要注意：

化学式列需要以Formula结尾

Smiles码列的列名需要以Smiles结尾

第四/五等等列：时间、温度、浓度等其它不需要再次特征化的实验变量

最右一列：输出目标（稳定性、效率、带隙等）

| 1              | 2              | 3                             | 4 | 5                 |
|----------------|----------------|-------------------------------|---|-------------------|
| layer1_Formula | layer2_Formula | SMILES_Smiles                 |   | conservation rate |
| CsPbBr2I       | CsPbBr3        | NH4+].[Cl-]                   |   | 0.041727          |
| CsPbBr2I       | CsPbBr3        | CCCC[N+](CCCC)(CCCC)CCCC.[I-] |   | 0.326007          |
| CsPbBr2I       | CsPbClI2       | Br-].[Br-].[Sn+2]             |   | 0.257407          |
| CsPbBr2I       | CsPbClI2       | Sn+2].[I-].[I-]               |   | 0.721138          |
| CsPbBr3        | CH3NH3PbBr2I   | Cl-].[Ag+]                    |   | 0.143376          |
| CsPbBr3        | CH3NH3PbBr2I   | C[NH3+].[I-]                  |   | 0.132046          |
| CsPbCl2Br      | CH3NH3SnBrI2   | Br-].[Ag+]                    |   | 0.173264          |
| CsPbCl2Br      | CH3NH3SnBrI2   | C(=N)N.[I-]                   |   | 0.375907          |
| CsPbCl2I       | CsPbCl2Br      | Cl-].[K+]                     |   | 0.188991          |
| CsPbCl2I       | CsPbCl2Br      | NH4+].[Cl-]                   |   | 0.185522          |
| CsPbI3         | CsPbBr2I       | I-].[I-].[Pb+2]               |   | 0.137411          |
| CsPbI3         | CsPbBr2I       | CCCC[N+](CCCC)(CCCC)CCCC.[I-] |   | 0.506803          |
| CsPbI3         | CsPbI3         | Br-].[Br-].[Pb+2]             |   | 0.287676          |
| CsPbI3         | CsPbI3         | Br-].[Br-].[Sn+2]             |   | 0.127879          |
| CsPbI3         | CH3NH3PbI3     | Li+].[Cl-]                    |   | 0.528662          |
| CsPbI3         | CH3NH3PbI3     | C[NH3+].[I-]                  |   | 0.285946          |
| CsSnBrI2       | (NH2)2CHSnBr2I | C(=N)N.Cl                     |   | 0.9163            |
| CsSnBrI2       | (NH2)2CHSnBr2I | CN.Cl                         |   | 1.306513          |
| CsSnBrI2       | (NH2)2CHSnBr2I | Br-].[Br-].[Pb+2]             |   | 0.465672          |
| CsSnBrI2       | (NH2)2CHSnBr2I | Cl-].[Cl-].[Pb+2]             |   | 0.333333          |
| CsSnBrI2       | (NH2)2CHSnBrI2 | NH4+].[Cl-]                   |   | 0.934002          |
| CsSnBrI2       | CH3NH3SnClI2   | Cl-].[Cl-].[Pb+2]             |   | 0.634021          |
| CsSnBrI2       | CH3NH3SnClI2   | Li+].[Cl-]                    |   | 0.873021          |
| CsSnCl3        | (NH2)2CHSnClF2 | C(=N)N.Cl                     |   | 0.536913          |
| CsSnCl3        | (NH2)2CHSnClI2 | C(=N)N.[I-]                   |   | 0.634002          |
| CsSnClBr2      | (NH2)2CHSnClI2 | Cl-].[Cl-].[Cu+2]             |   | 0.287739          |
| CsSnClBr2      | (NH2)2CHSnClI3 | Br-].[Br-].[Pb+2]             |   | 0.234039          |
| CsSnClBr2      | (NH2)2CHSnClI3 | I-].[I-].[Pb+2]               |   | 0.334286          |
| CsSnClF2       | CH3NH3PbCl2I   | NH4+].[Cl-]                   |   | 0.234039          |
| CsSnClF2       | CH3NH3PbCl2I   | CCCC[N+](CCCC)(CCCC)CCCC.[I-] |   | 0.733404          |
| (NH2)2CHPbBr2  | CsPbI3         | Cl-].[Cl-].[Pb+2]             |   | 0.812222          |
| (NH2)2CHPbBr2  | CsPbI3         | Sn+2].[I-].[I-]               |   | 0.573566          |
| (NH2)2CHPbBr2  | CsSnBrI2       | C(=N)N.Cl                     |   | 1.013255          |
| (NH2)2CHPbBr2  | CsSnBrI2       | C(=N)N.[I-]                   |   | 0.605159          |
| (NH2)2CHPbCl2  | CsPbCl2I       | CN.Br                         |   | 0.272121          |
| (NH2)2CHPbCl2  | CsPbCl2I       | Sn+2].[I-].[I-]               |   | 0.344214          |
| (NH2)2CHPbI3   | CH3NH3PbBr2I   | Br-].[Br-].[Sn+2]             |   | 0.219089          |
| (NH2)2CHPbI3   | CH3NH3PbBr2I   | Sn+2).[I-].[I-]               |   | 0.291083          |
| (NH2)2CHPbI3   | CH3NH3PbI3     | CN.Br                         |   | 0.91573           |
| (NH2)2CHPbI3   | CH3NH3PbI3     | CCCC[N+](CCCC)(CCCC)CCCC.[I-] |   | 0.412811          |
| (NH2)2CHPbI3   | CH3NH3SnCl2Br  | C(=N)N.Cl                     |   | 0.362069          |
| (NH2)2CHPbI3   | CH3NH3SnCl2Br  | C(=N)N.[I-]                   |   | 0.794304          |
| (NH2)2CHSnCl3  | CH3NH3SnClF2   | F-].[F-].[Al+3]               |   | 0.363465          |

变量1  
需要  
特征化

变量2  
需要  
特征化

变量3  
需要  
特征化

其它  
变量：  
时间  
温度  
浓度  
动作  
等等

输出  
目标：  
稳定性  
效率  
带隙  
性能  
等等

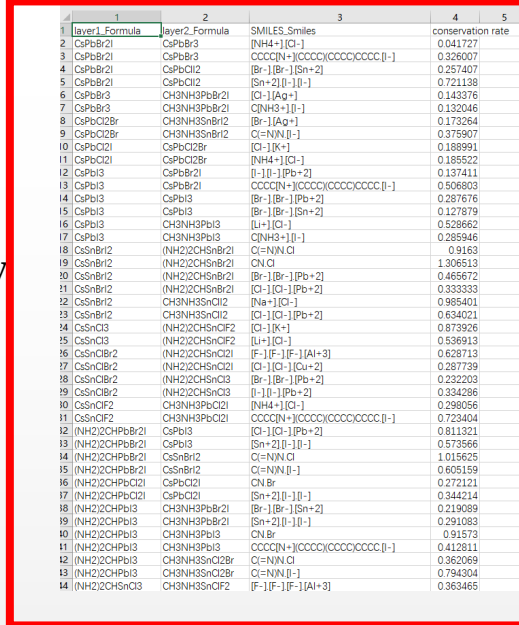
时间、温度、浓度

Layer 3 Smiles

Layer 2 Formula

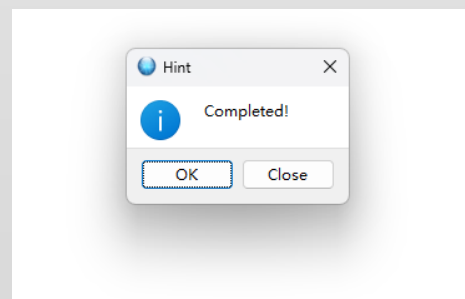
Layer 1 Formula

自定义文件夹→Descriptor generation →Hybrid inorganic organic (2in1) →train\_test\_dataset.csv



## Smiles码： RDKit描述符

## 自动特征化



## 训练测试集

但是需要后面进行特征选择!!!!

[illegible]

特征化按钮 (训练测试集辅助生成) 多种有机分子:例如双分子SMILES码 → RDKit描述符

## Descriptors: Organic Multiple column Smiles → Import → Featurize

12



可以是已经特征化好的数字csv

也可以是最原始的Smiles码或者化学式

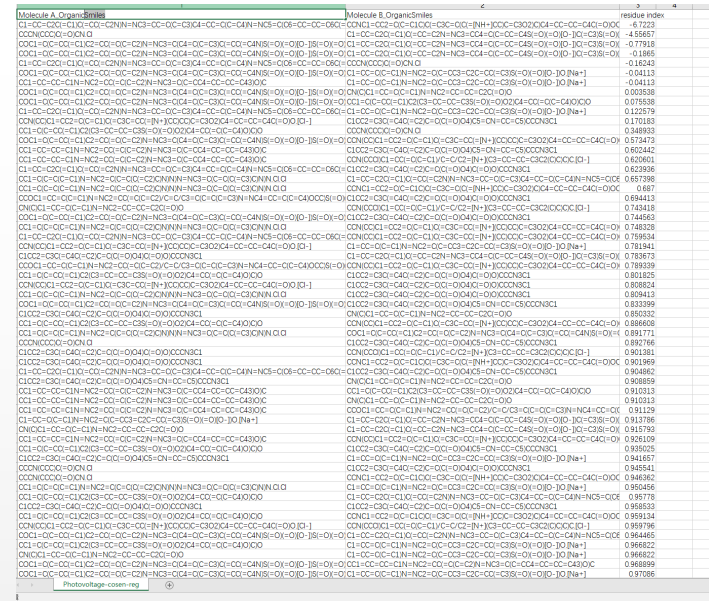
**如果Smiles码，  
请将列名后缀标为Smiles**

输入变量1 (分子A)

输入变量2(分子B)

输出列:  
稳定性  
效率  
带隙  
性能  
等等

## The logo of Nanjing University of Information Engineering (NUIE) is a circular emblem. It features a stylized blue wave or 'S' shape in the center. The text '南京信息工程大学' (Nanjing University of Information Engineering) is written in Chinese characters along the top arc, and '1960' is at the bottom. The English name 'NANJING UNIVERSITY OF INFORMATION SCIENCE &amp; TECHNOLOGY' is written along the bottom arc.



但是需要后面进行特征选择!!!

[illegible]

# Smiles: 可以从Pubchem等网站中查询





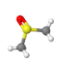

Pubchem

NIH National Library of Medicine  
National Center for Biotechnology Information

PubChem About Docs Submit Contact

COMPOUND SUMMARY

**Dimethyl Sulfoxide**

|                   |  |
|-------------------|--|
| PubChem CID       | 679  |
| Structure         | <div><br/>2D</div> <div><br/>3D</div> <div><br/>Crystal</div> |
| Chemical Safety   | <br>Irritant<br><a href="#">Laboratory Chemical Safety Summary (LCSS) Datasheet</a>  |
| Molecular Formula | C <sub>2</sub> H <sub>6</sub> OS<br>(CH <sub>3</sub> ) <sub>2</sub> SO   |
| Synonyms          | dimethyl sulfoxide<br>DMSO   |

► PubChem

## 2.1.4 Canonical SMILES

**CS(=O)C**

*Computed by OEChem 2.3.0 (PubChem release 2021.10.14)*

► PubChem

## 2.3 其他特征化按钮

NJmatML

| Directory | Descriptors                     | Dataset | Preprocessing | Machine learning        | Prediction | Share |
|-----------|---------------------------------|---------|---------------|-------------------------|------------|-------|
| 推荐        | Hybrid inorganic organic (2in1) |         |               | Import                  |            |       |
|           | Organic Multiple column Smiles  |         |               | Featurize               |            |       |
|           | Organic descriptor              |         |               | 单列Smiles码 → 生成RDKit描述符  |            |       |
|           | Inorganic descriptor (one-hot)  |         |               | 单列化学式 → 生成元素周期表独热编码     |            |       |
|           | Inorganic descriptor (magpie)   |         |               | 单列化学式 → 生成magpie元素基础描述符 |            |       |



## 2.3 其他化学式特征化按钮：独热编码

|   | formula      |
|---|--------------|
| 0 | (Fe2AgCu2)O3 |
| 1 | Fe2O3        |
| 2 | CsPbI3       |
| 3 | MoS2         |
| 4 | CuInGaSe     |
| 5 | Si           |
| 6 | TiO2         |

**TiO<sub>2</sub>:**  
**O: 0.666666为O**  
**Ti: 0.333333为Ti**



| H | He | Li | Be | B | C | N | O        | F | Ne | Na | Mg | Al | Si | P | S        | Cl | Ar | K | Ca |
|---|----|----|----|---|---|---|----------|---|----|----|----|----|----|---|----------|----|----|---|----|
| 0 | 0  | 0  | 0  | 0 | 0 | 0 | 0.375    | 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0        | 0  | 0  | 0 | 0  |
| 0 | 0  | 0  | 0  | 0 | 0 | 0 | 0.6      | 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0        | 0  | 0  | 0 | 0  |
| 0 | 0  | 0  | 0  | 0 | 0 | 0 | 0        | 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0        | 0  | 0  | 0 | 0  |
| 0 | 0  | 0  | 0  | 0 | 0 | 0 | 0        | 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0.666667 | 0  | 0  | 0 | 0  |
| 0 | 0  | 0  | 0  | 0 | 0 | 0 | 0        | 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0        | 0  | 0  | 0 | 0  |
| 0 | 0  | 0  | 0  | 0 | 0 | 0 | 0        | 0 | 0  | 0  | 0  | 0  | 0  | 1 | 0        | 0  | 0  | 0 | 0  |
| 0 | 0  | 0  | 0  | 0 | 0 | 0 | 0.666667 | 0 | 0  | 0  | 0  | 0  | 0  | 0 | 0        | 0  | 0  | 0 | 0  |

| P | S | Cl       | Ar | K | Ca | Sc | Ti       | V | Cr | Mn | Fe   | Co | Ni | Cu   | Zn | Ga   | Ge | As | Se |
|---|---|----------|----|---|----|----|----------|---|----|----|------|----|----|------|----|------|----|----|----|
| 0 | 0 | 0        | 0  | 0 | 0  | 0  | 0        | 0 | 0  | 0  | 0.25 | 0  | 0  | 0.25 | 0  | 0    | 0  | 0  | 0  |
| 0 | 0 | 0        | 0  | 0 | 0  | 0  | 0        | 0 | 0  | 0  | 0.4  | 0  | 0  | 0    | 0  | 0    | 0  | 0  | 0  |
| 0 | 0 | 0        | 0  | 0 | 0  | 0  | 0        | 0 | 0  | 0  | 0    | 0  | 0  | 0    | 0  | 0    | 0  | 0  | 0  |
| 0 | 0 | 0.666667 | 0  | 0 | 0  | 0  | 0        | 0 | 0  | 0  | 0    | 0  | 0  | 0    | 0  | 0    | 0  | 0  | 0  |
| 0 | 0 | 0        | 0  | 0 | 0  | 0  | 0        | 0 | 0  | 0  | 0    | 0  | 0  | 0.25 | 0  | 0.25 | 0  | 0  | 0  |
| 1 | 0 | 0        | 0  | 0 | 0  | 0  | 0        | 0 | 0  | 0  | 0    | 0  | 0  | 0    | 0  | 0    | 0  | 0  | 0  |
| 0 | 0 | 0        | 0  | 0 | 0  | 0  | 0.333333 | 0 | 0  | 0  | 0    | 0  | 0  | 0    | 0  | 0    | 0  | 0  | 0  |



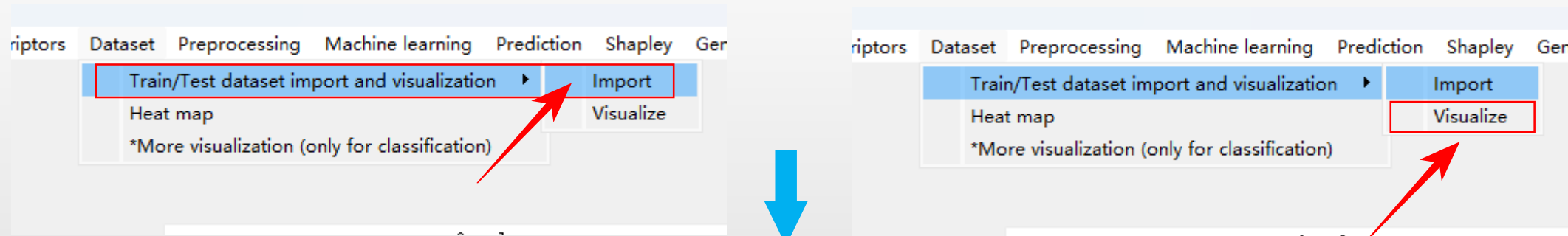
### 3 导入数据集 Dataset import

# 3.1 导入数据集

数据集为之前软件自动特征化的train\_test\_dataset.csv

例如：自定义文件夹→Descriptor generation →Hybrid inorganic organic (2in1) → train\_test\_dataset.csv

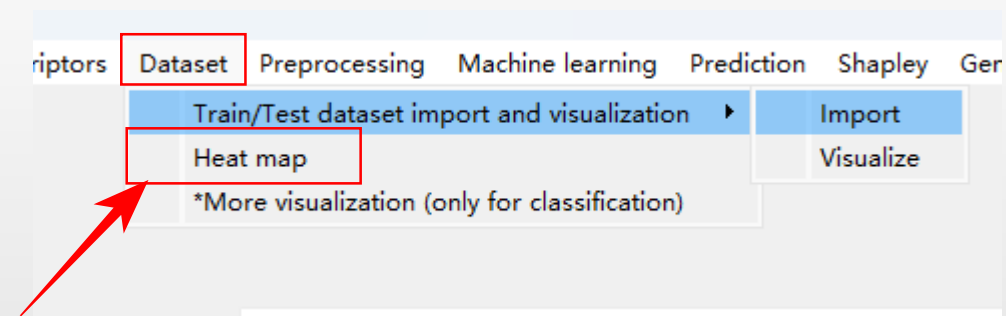
更好的方式：领域知识描述符 (customized)



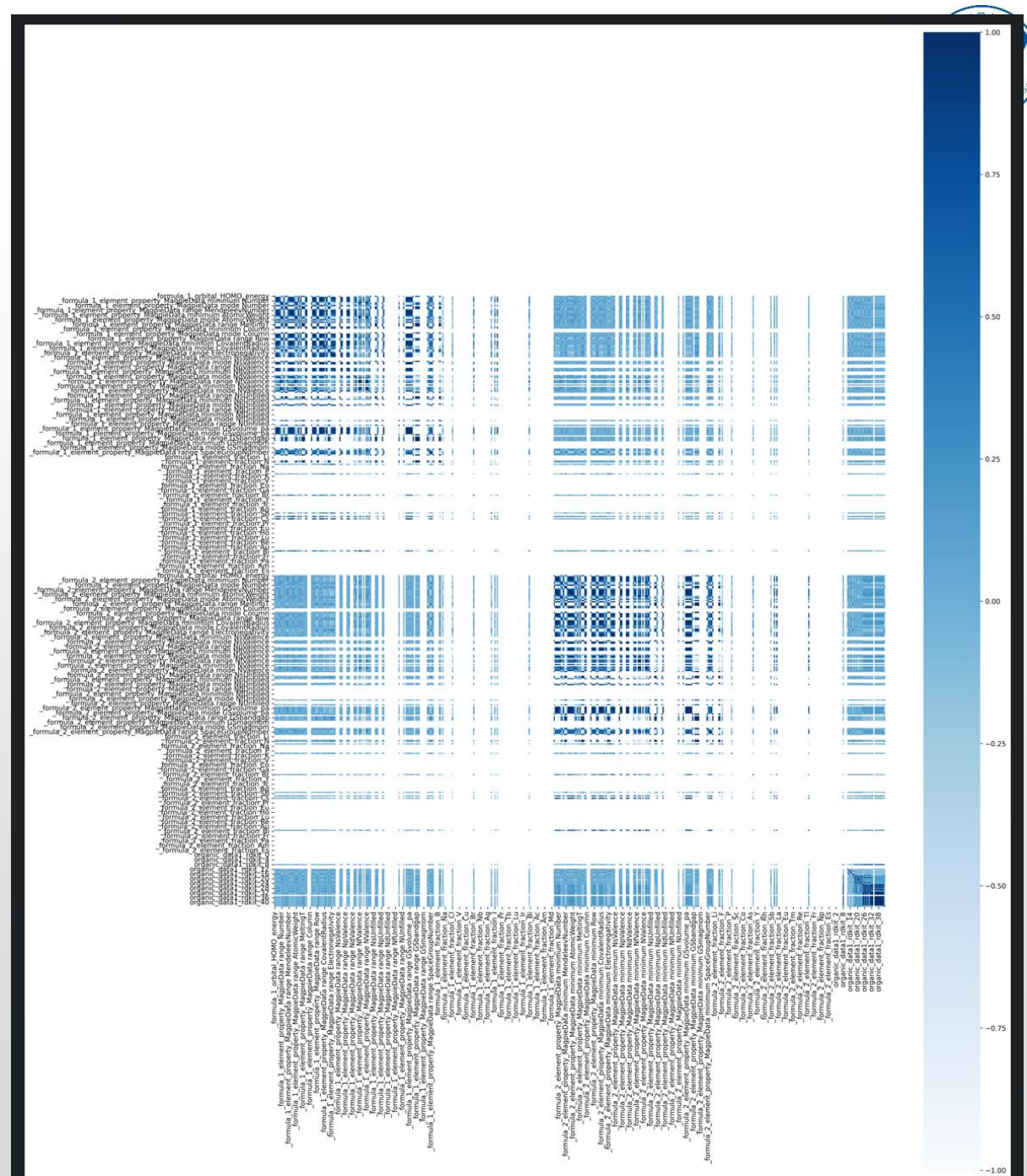
特征选择前的数据可视化



### 3.2 导入数据集后的热图



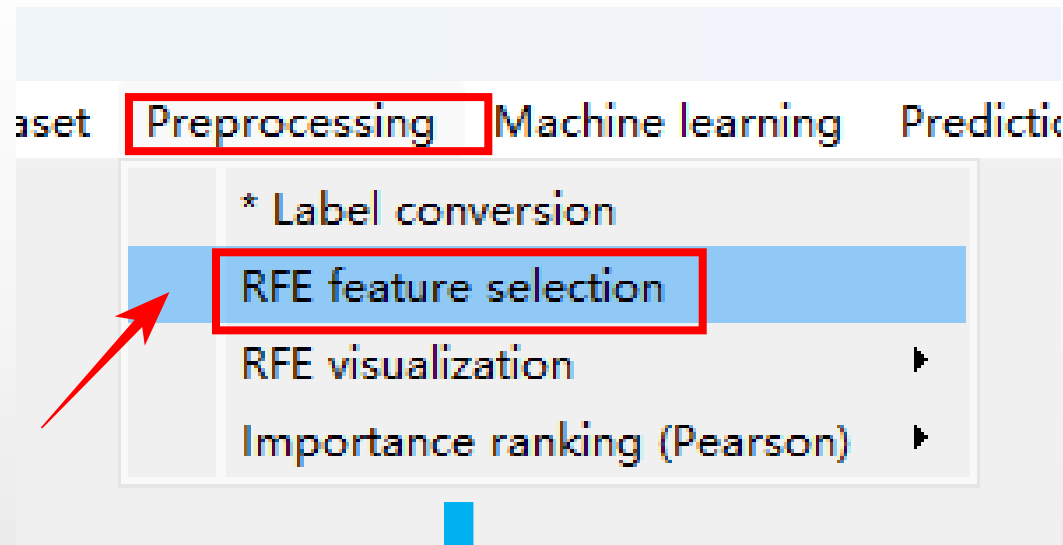
# Heat Map



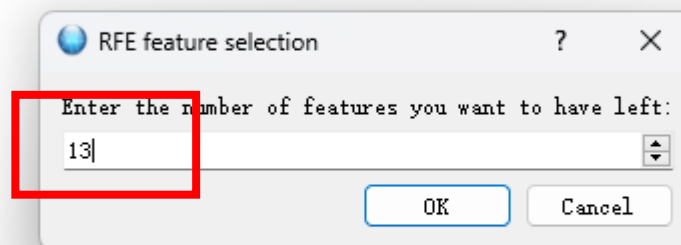
# 4 Preprocessing

# 4.1 Preprocessing: RFE 特征选择

特征选择: RFE  
Recursive feature elimination



输入想要保留的特征数目  
例如 13





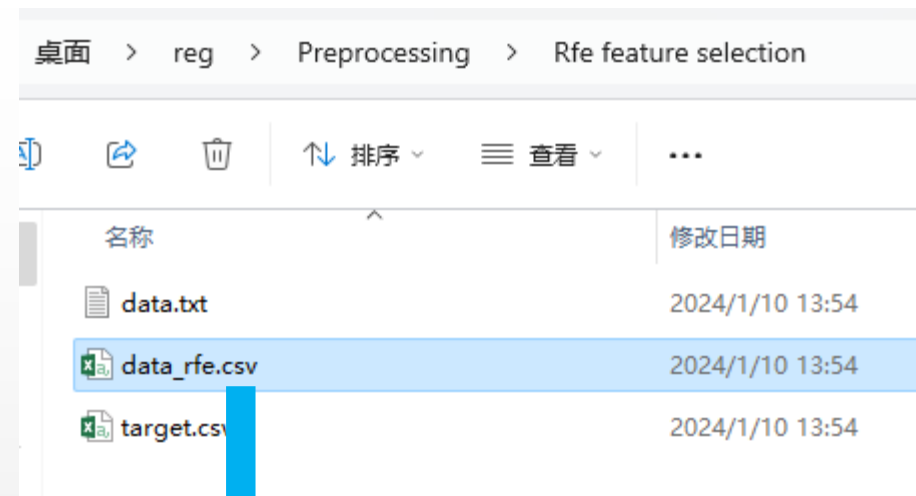
# 4.1 Preprocessing

特征选择后的数据集保存在data\_rfe.csv:

自定义folder → Proprocessing →  
Rfe feature selection → data\_rfe.csv

该数据集实质用于机器学习建模

如果对话框填写13, 自动生成的data\_rfe.csv中的特征为13列, 左右一列为输出数据



输入特征

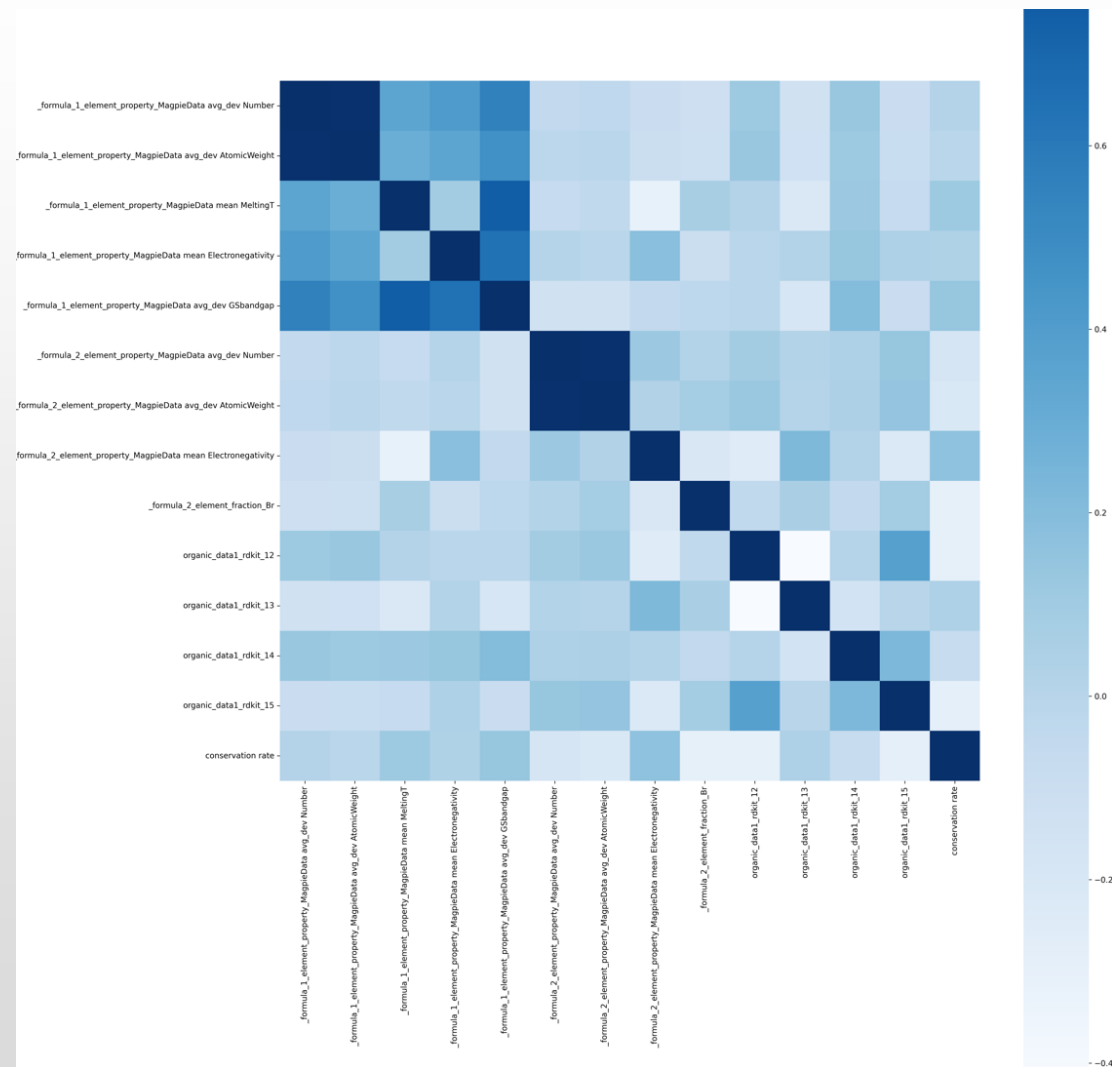
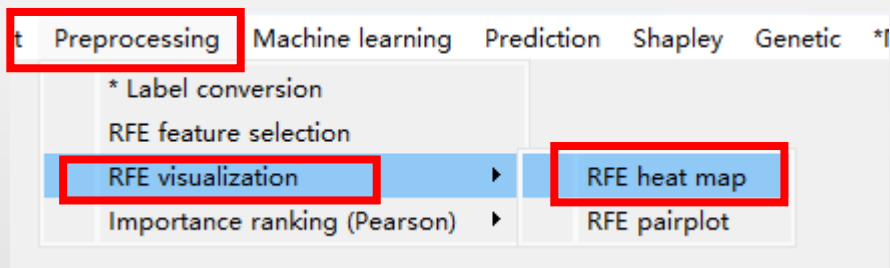
输出目标

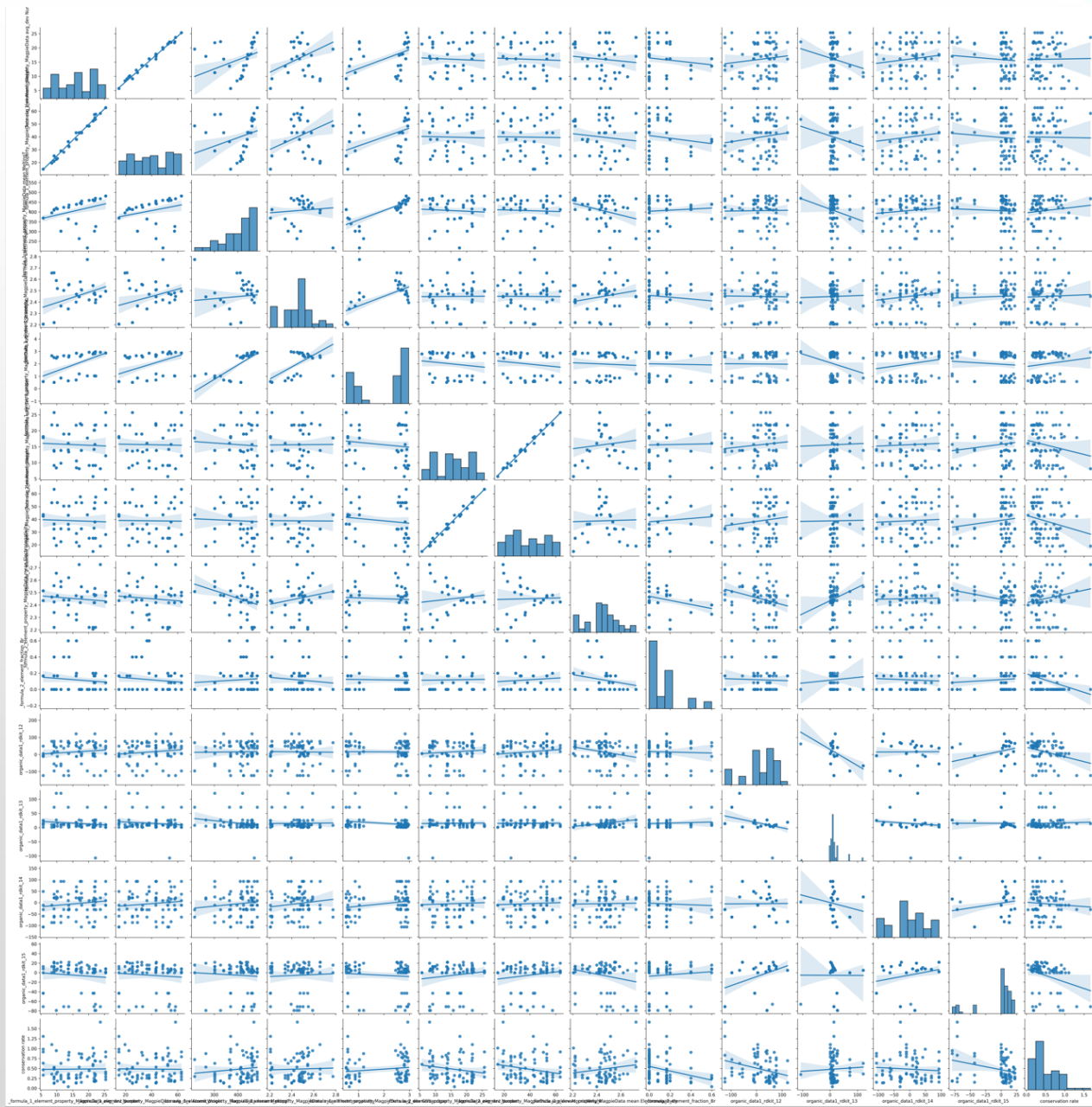
| formula_1 | formula_1 | formula_1 | formula_1 | formula_1 | formula_2 | formula_2 | formula_2 | formula_2 | organic_de | organic_de | organic_de | organic_de | conservation rate |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------|------------|------------|------------|-------------------|
| 13.6      | 36.36767  | 364.13    | 0.7952    | 0.63616   | 16.08     | 2.4       | 6         | 0.6       | 69         | 8          | -31        | 16         | 0.041727          |
| 13.6      | 36.36767  | 364.13    | 0.7952    | 0.63616   | 16.08     | 2.4       | 6         | 0.6       | -96        | 72         | -63        | 0          | 0.326007          |
| 13.6      | 36.36767  | 364.13    | 0.7952    | 0.63616   | 14        | 2.32      | 8         | 0         | 49         | 10         | 37         | 8          | 0.257407          |
| 13.6      | 36.36767  | 364.13    | 0.7952    | 0.63616   | 14        | 2.32      | 8         | 0         | 77         | 10         | 9          | 8          | 0.721138          |
| 16.08     | 43.27139  | 339.92    | 0.8742    | 0.69936   | 21.72222  | 2.475     | 7.944444  | 0.166667  | 121        | 19         | -83        | 5          | 0.143376          |
| 16.08     | 43.27139  | 339.92    | 0.8742    | 0.69936   | 21.72222  | 2.475     | 7.944444  | 0.166667  | 53         | 4          | 21         | 4          | 0.132046          |
| 21.84     | 57.49571  | 302.24    | 1.2886    | 1.03088   | 20.16667  | 2.419167  | 6.388889  | 0.083333  | 49         | 10         | 69         | 7          | 0.173264          |
| 21.84     | 57.49571  | 302.24    | 1.2886    | 1.03088   | 20.16667  | 2.419167  | 6.388889  | 0.083333  | 5          | 1          | -11        | 0          | 0.375907          |
| 22.24     | 57.70415  | 326.45    | 1.2096    | 1.02672   | 21.84     | 2.48      | 8.4       | 0.2       | 41         | 27         | 16         | -79        | 0.188991          |
| 22.24     | 57.70415  | 326.45    | 1.2096    | 1.02672   | 21.84     | 2.48      | 8.4       | 0.2       | 69         | 8          | -31        | 16         | 0.185522          |
| 9.12      | 25.21449  | 412.55    | 0.6372    | 0.50976   | 13.6      | 2.34      | 6         | 0.4       | 77         | 10         | -87        | 13         | 0.137411          |
| 9.12      | 25.21449  | 412.55    | 0.6372    | 0.50976   | 13.6      | 2.34      | 6         | 0.4       | -96        | 72         | -63        | 0          | 0.506803          |
| 9.12      | 25.21449  | 412.55    | 0.6372    | 0.50976   | 9.12      | 2.22      | 6         | 0         | 49         | 10         | -59        | 13         | 0.287676          |
| 9.12      | 25.21449  | 412.55    | 0.6372    | 0.50976   | 9.12      | 2.22      | 6         | 0         | 49         | 10         | 37         | 8          | 0.127879          |
| 9.12      | 25.21449  | 412.55    | 0.6372    | 0.50976   | 25.72222  | 2.425     | 7.944444  | 0         | 17         | 5          | 21         | 20         | 0.528662          |
| 9.12      | 25.21449  | 412.55    | 0.6372    | 0.50976   | 25.72222  | 2.425     | 7.944444  | 0         | 53         | 4          | 21         | 4          | 0.285946          |
| 5.68      | 14.86467  | 369.234   | 0.7162    | 0.57296   | 17.83333  | 2.514167  | 6.166667  | 0.166667  | 5          | 1          | -11        | 0          | 0.9163            |
| 5.68      | 14.86467  | 369.234   | 0.7162    | 0.57296   | 17.83333  | 2.514167  | 6.166667  | 0.166667  | -123       | 14         | -30        | -71        | 1.306513          |
| 5.68      | 14.86467  | 369.234   | 0.7162    | 0.57296   | 17.83333  | 2.514167  | 6.166667  | 0.166667  | 49         | 10         | -59        | 13         | 0.465672          |

### 特征选择后的数据集保存在data\_rfe.csv:

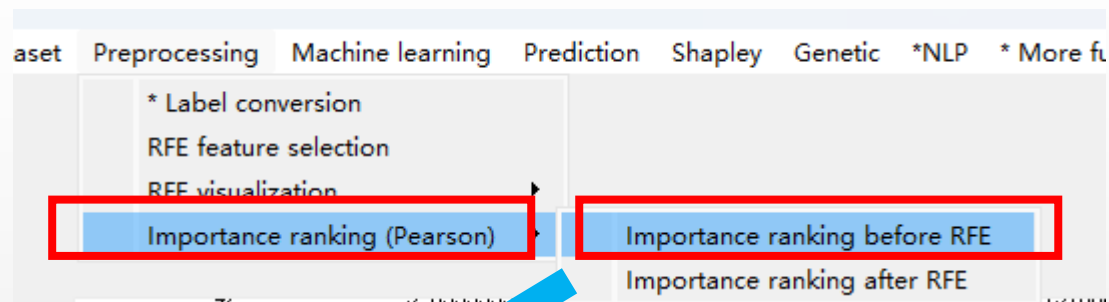


## 4.2 Preprocessing → 特征选择后热图

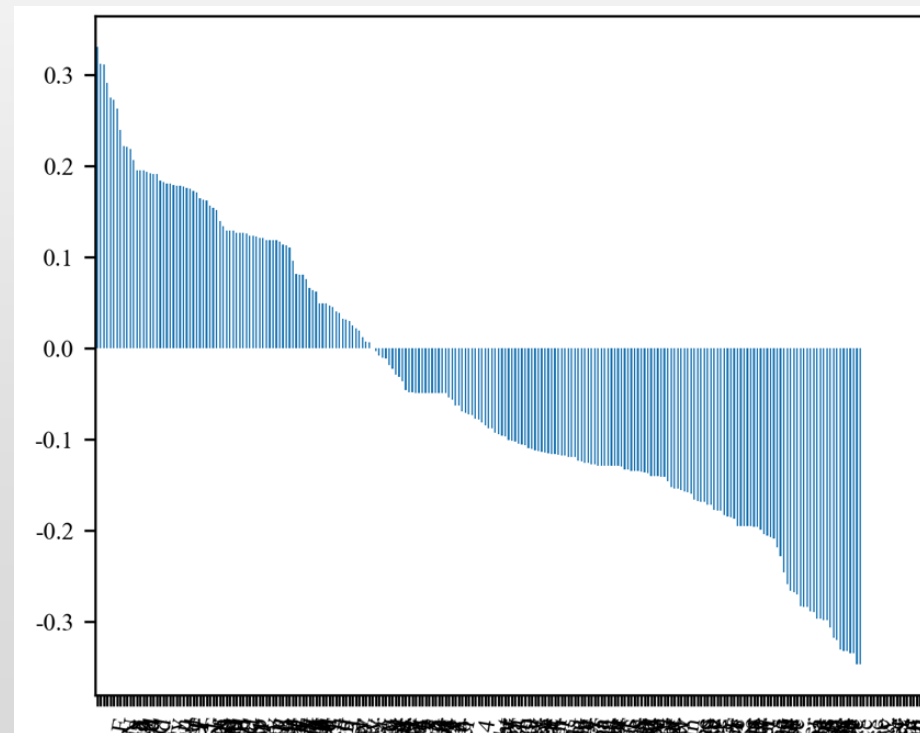
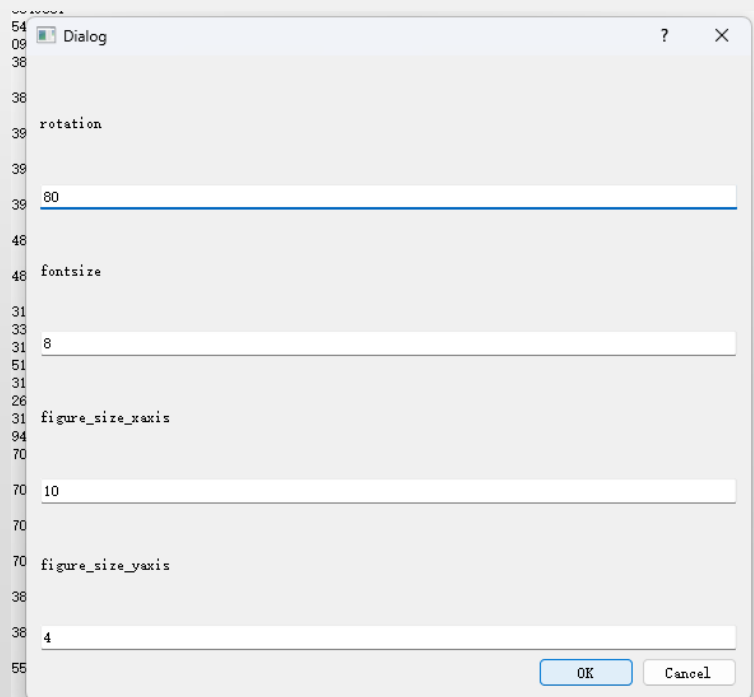




## 4.3.1 Preprocessing→RFE前的特征重要性排名

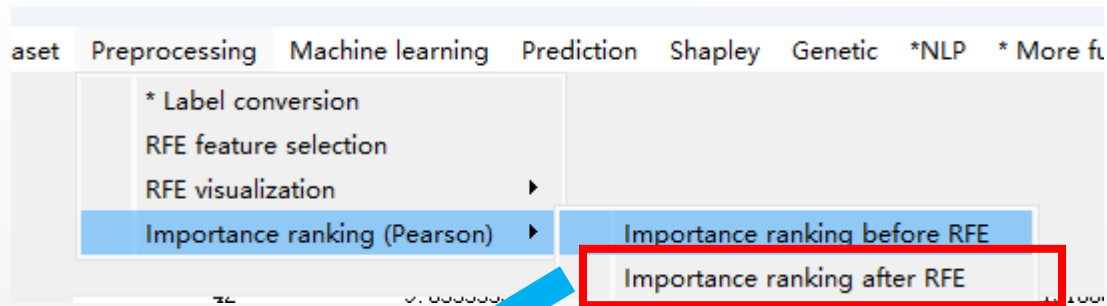


特征重要性排名 (Pearson)

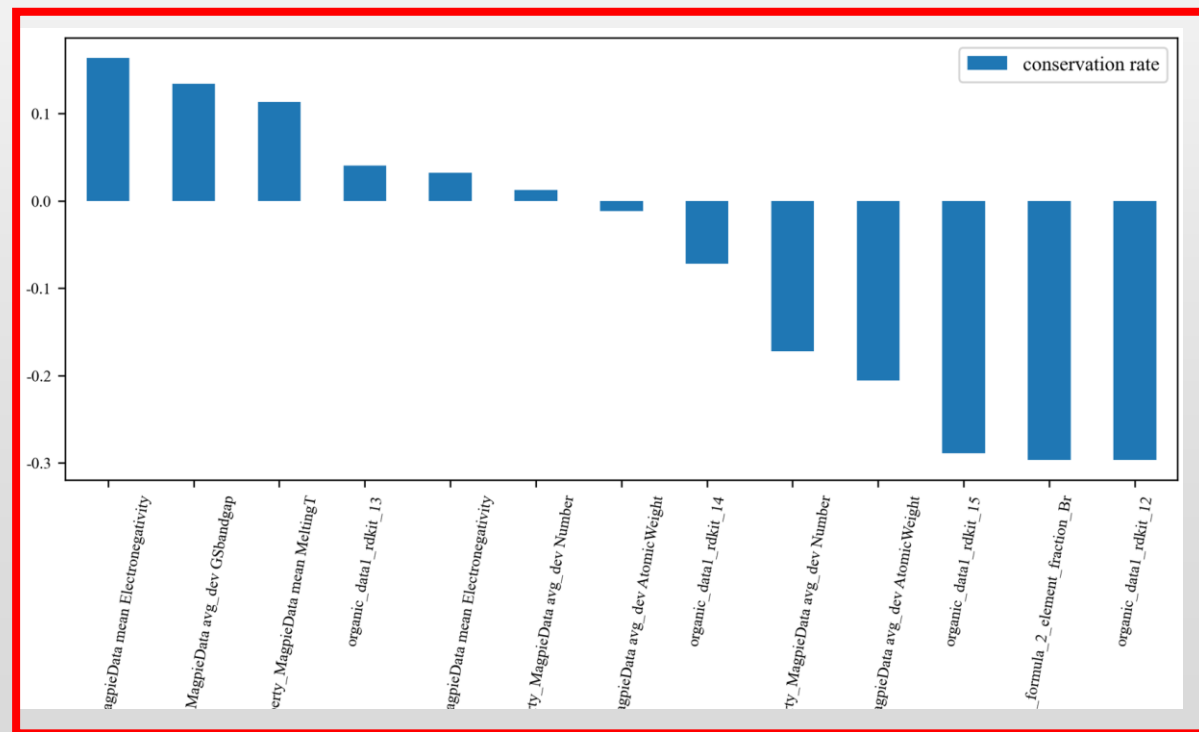
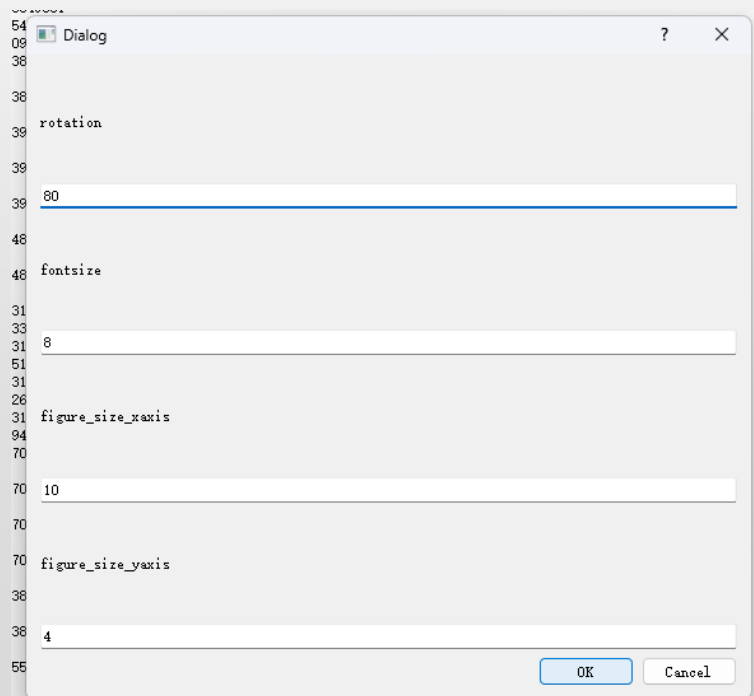


Shap排名更被接受

## 4.3.2 Preprocessing → RFE后的特征重要性排名



RFE特征选择后特征重要性排名 (Pearson)

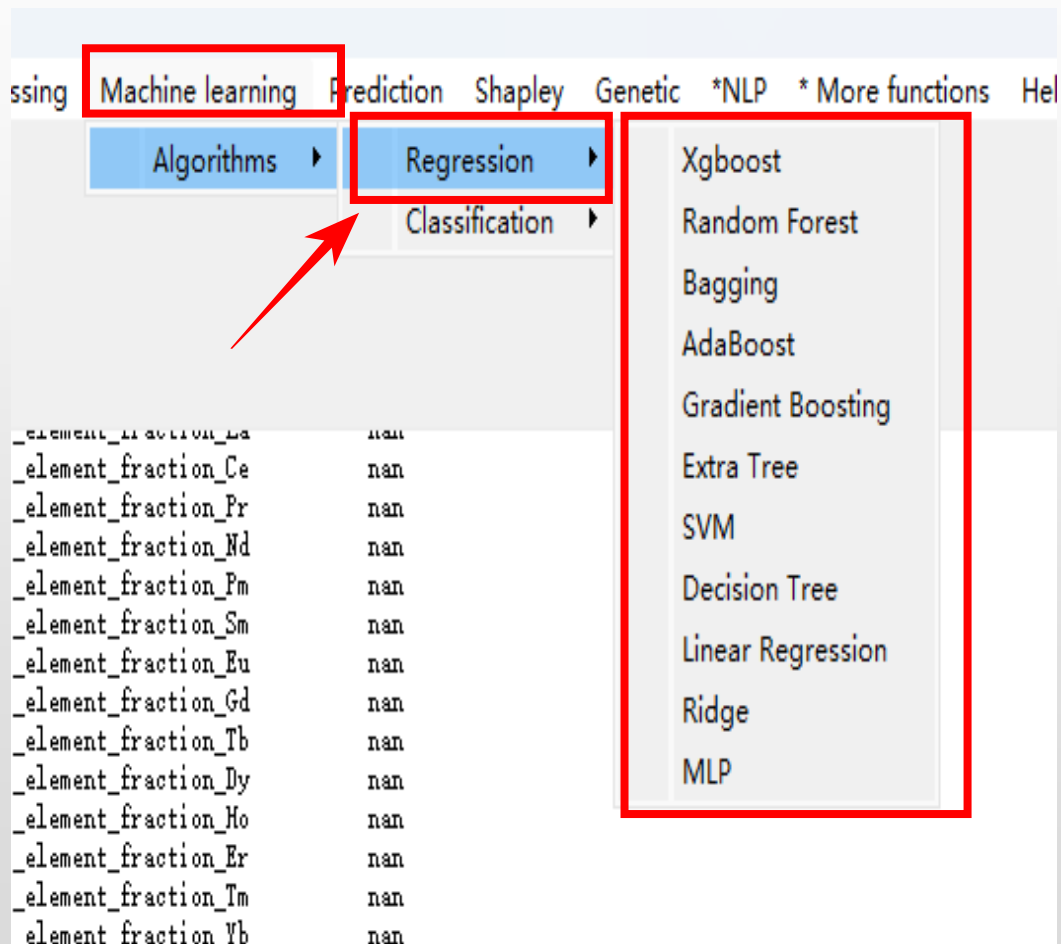


Shap排名更被接受

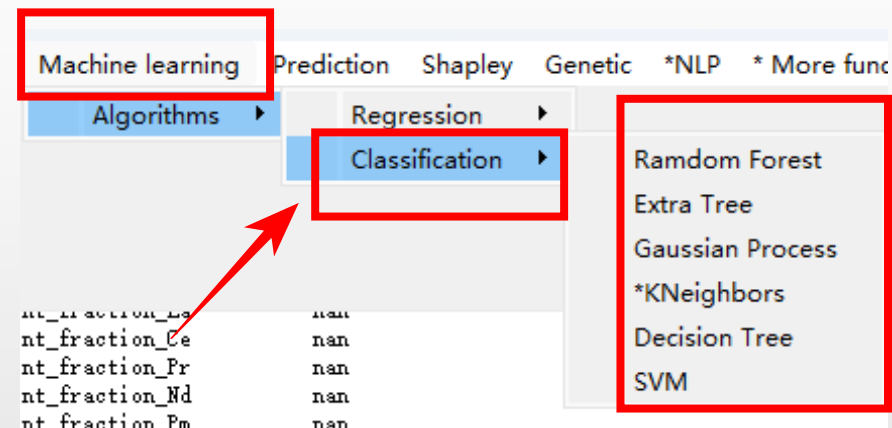
## 5 机器学习建模

# 5.1 机器学习建模

## 机器学习回归任务



## 机器学习分类任务



## 5.2 机器学习:超参数设置 (单独对话框)

例: XGBoost 超参数默认, 可调参

A dialog box titled "Dialog" with a question mark icon and a close button. It contains several input fields for XGBoost parameters. The parameters and their values are: n\_estimators (1000), max\_depth (200), eta (0.2), gamma (0), subsample (0.9), colsample\_bytree (0.8), and learning\_rate (0.2). At the bottom right, there are "OK" and "Cancel" buttons.

### Machine learning

按钮下所有回归和分类机器学习算法皆可自定义超参数

例: 随机森林超参数默认, 可调参

A dialog box titled "Dialog" with a question mark icon and a close button. It contains several input fields for Random Forest parameters. The parameters and their values are: max\_depth (7), random\_state (0), min\_samples\_leaf (1), max\_features (1), min\_samples\_split (2), and n\_estimators (100). At the bottom right, there are "OK" and "Cancel" buttons.

# 5.3 机器学习:数据拟合、准确率与超参数重新设置

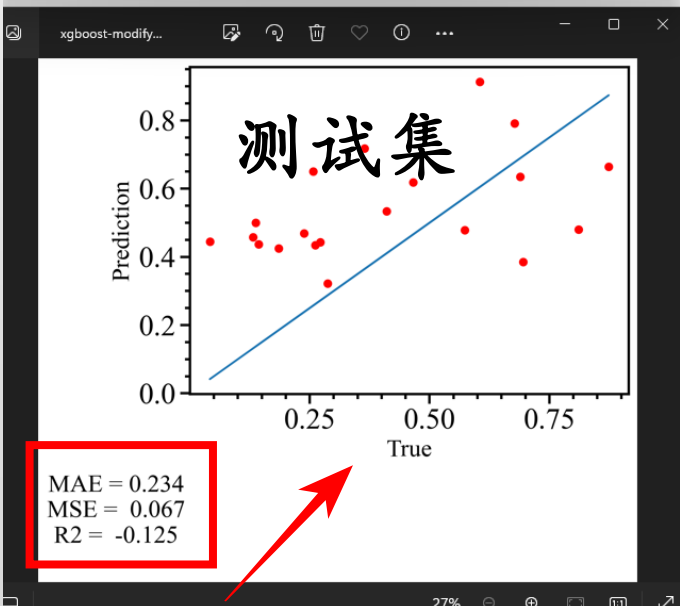
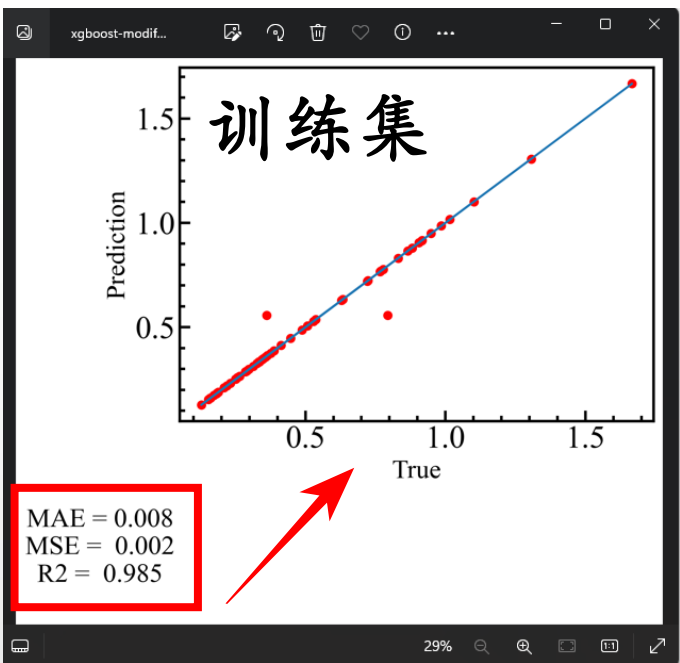
默认超算数

重新调节超算数

数据拟合

准确率提升方法:

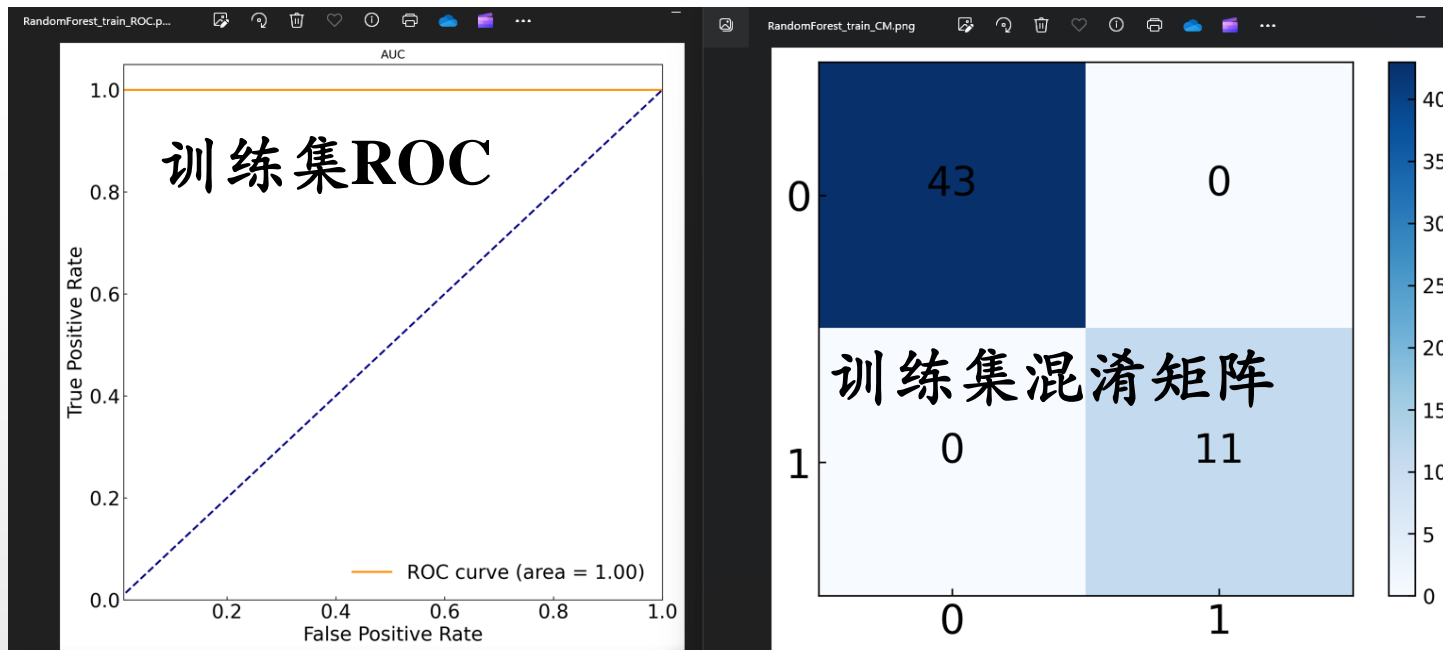
1. Customized 描述符!!!
2. 机器学习 (算法与超参数)



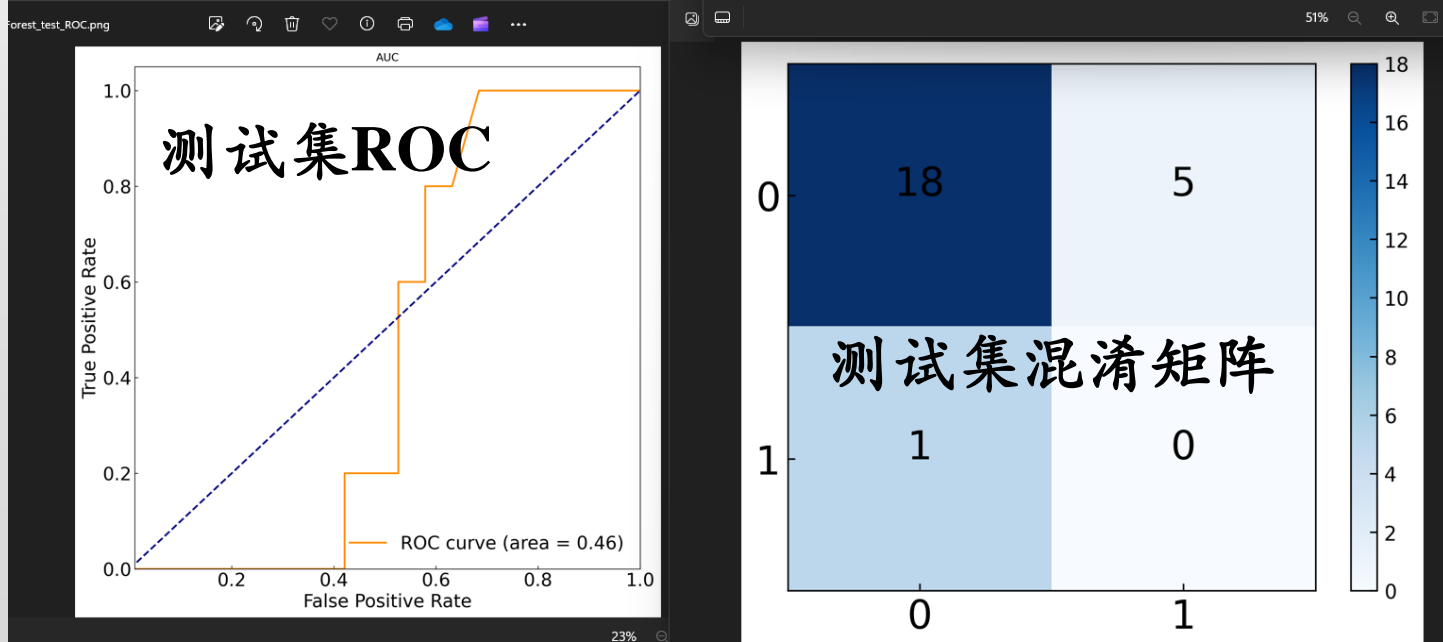
反复调参, 直至测试集准确



## 5.4 机器学习分类任务 (e.g., 稳定与否)



反复调参，直至测试集准确



Again, 准确率提升方法:

1. Customized 描述符! (重点)
2. 机器学习 (算法与超参数)

## 6 虚拟空间预测

# 预测虚拟空间

如果需要辅助特征化

从import and generate做起

按照下述步骤：

Prediction → Import and generate (保持与前一一致) →  
Import → Generate → Import virtual data (without  
label) →  
Select machine learning model → Prediction generation  
(with label)

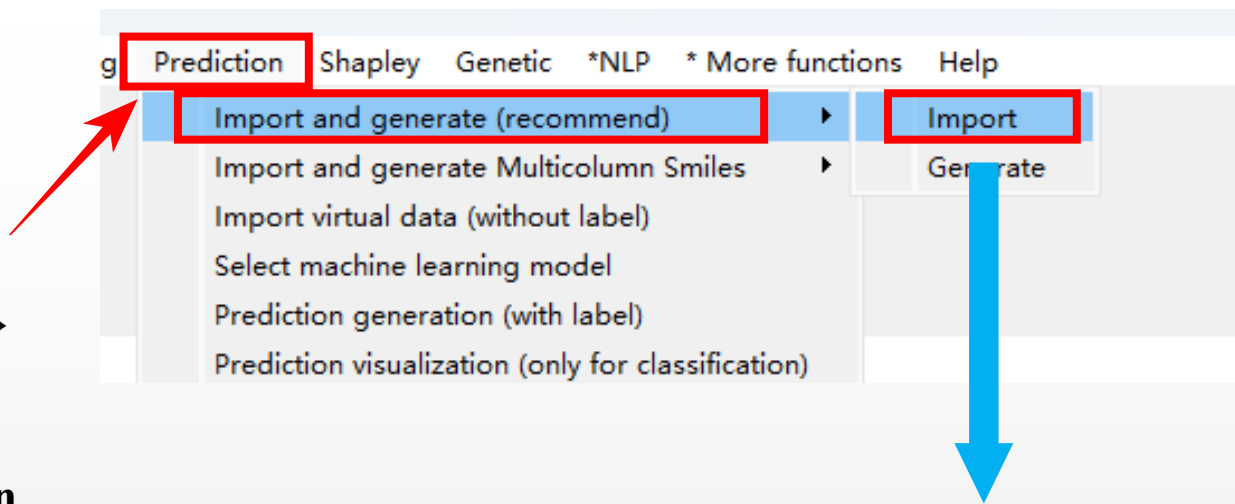
如果不需要辅助特征化、已经有特征化好的虚拟空间，  
直接跳过本步骤：

直接选择Prediction → Import virtual data (without  
label) → Select machine learning model → Prediction  
generation (with label)

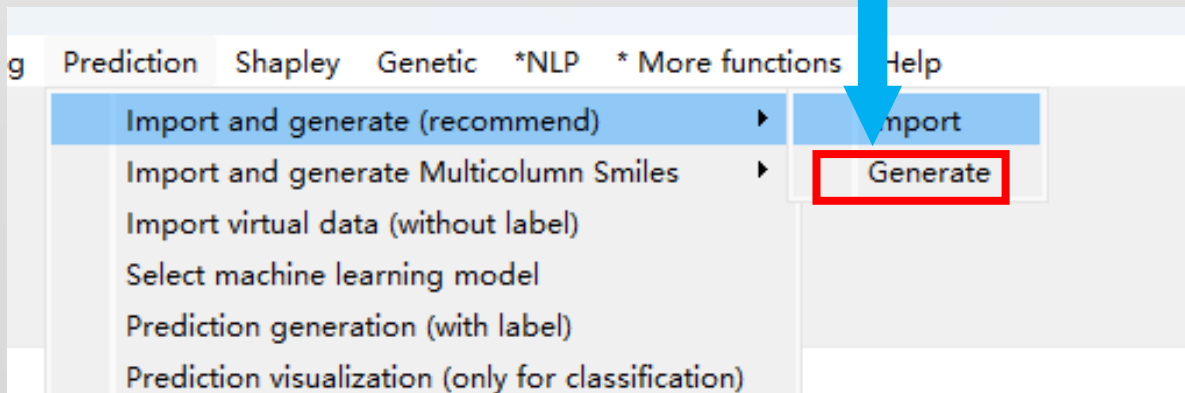
虚拟空间csv保持和原始csv格式一致

但最右一列为空（需要预测）

最右一列需有列名



|   | 1              | 2              | 3                            | 4      |
|---|----------------|----------------|------------------------------|--------|
| 1 | layer1_Formula | layer2_Formula | SMILES_Smiles                | target |
| 2 | CH3NH3SnF2Cl   | CH3NH3SnI3     | C1=CC(=CC(=C1)Br)/C=C/C(=O)O |        |
| 3 | CsSnClBr2      | CH3NH3SnI3     | [NH4+].[Cl-]                 |        |
| 4 | CH3NH3PbI3     | (NH2)2CHSnI3   | C(=N)N.[I-]                  |        |
| 5 | CsSnF2Cl       | (NH2)2CHSnI3   | [Li+].[Cl-]                  |        |
| 6 | (NH2)2CHSnI3   | CH3NH3SnI3     | [Na+].[Cl-]                  |        |
| 7 | CsSnBr2I       | CsPbI3         | [Mg+2].[Cl-].[Cl-]           |        |
| 8 | CsSnBr2I       | CsPbFCl2       | [Cl-].[K+]                   |        |



# 预测虚拟空间



**Import virtual data (without label)**  
**Virtual\_generate\_final.csv**

| 桌面 > reg > Prediction > Import and generate (recommend) |                |                     |
|---|----------------|---------------------|
| 名称  | 修改日期           | 类型                  |
| Formula_selected_columns.csv                            | 2024/1/12 1:58 | Microsoft Excel ... |
| merged_rdkit_result.csv                                 | 2024/1/12 1:58 | Microsoft Excel ... |
| merged_result.csv                                       | 2024/1/12 1:58 | Microsoft Excel ... |
| Smiles_selected_columns.csv                             | 2024/1/12 1:58 | Microsoft Excel ... |
| train_test_dataset.csv                                  | 2024/1/12 1:58 | Microsoft Excel ... |
| unselected_columns.csv                                  | 2024/1/12 1:58 | Microsoft Excel ... |
| virtual_generate_final.csv                              | 2024/1/12 1:58 | Microsoft Excel ... |

# 预测虚拟空间

Prediction → Import virtual data (without label) → Select machine learning model →

**Prediction generation (with label)**

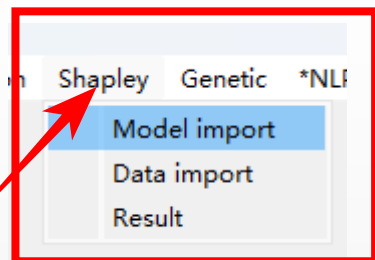
预测的数据（稳定性、效率等）  
最终数据

|   | 1         | 2         | 3         | 4         | 5         | 6         | 7         | 8         | 9         | 10         | 11         | 12         | 13         | 14       |
|---|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------|------------|------------|------------|----------|
| 1 | formula_1 | formula_1 | formula_1 | formula_1 | formula_1 | formula_2 | formula_2 | formula_2 | formula_2 | organic_da | organic_da | organic_da | organic_da | Output   |
| 2 | 8.388889  | 19.5697   | 396.1492  | 2.655833  | 2.656556  | 22.16667  | 53.7101   | 2.394167  | 0         | 62         | -107       | 3          | -66        | 0.938381 |
| 3 | 11.28     | 29.14595  | 301.974   | 2.366     | 0.86512   | 22.16667  | 53.7101   | 2.394167  | 0         | 69         | 8          | -31        | 16         | 0.435465 |
| 4 | 25.72222  | 63.54232  | 477.6058  | 2.425     | 2.972236  | 21.83333  | 52.98795  | 2.464167  | 0         | 5          | 1          | -11        | 0          | 0.519687 |
| 5 | 19.6      | 48.63574  | 217.054   | 2.774     | 1.02928   | 21.83333  | 52.98795  | 2.464167  | 0         | 17         | 5          | 21         | 20         | 0.441674 |
| 6 | 21.83333  | 52.98795  | 473.7317  | 2.464167  | 2.873903  | 22.16667  | 53.7101   | 2.394167  | 0         | -11        | 16         | 49         | 8          | 0.652173 |
| 7 | 8.48      | 22.20927  | 345.024   | 2.266     | 0.63616   | 9.12      | 25.21449  | 2.22      | 0         | 65         | 2          | -27        | 22         | 0.637941 |
| 8 | 8.48      | 22.20927  | 345.024   | 2.266     | 0.63616   | 26        | 67.2406   | 2.684     | 0         | 41         | 27         | 16         | -79        | 0.889752 |
| 9 |           |           |           |           |           |           |           |           |           |            |            |            |            |          |

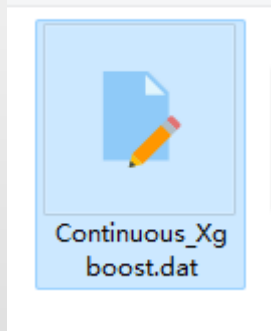
待实验验证！  
待模拟验证！

## 7 Shap 特征分析

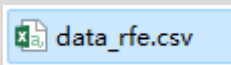
# Shap 特征分析



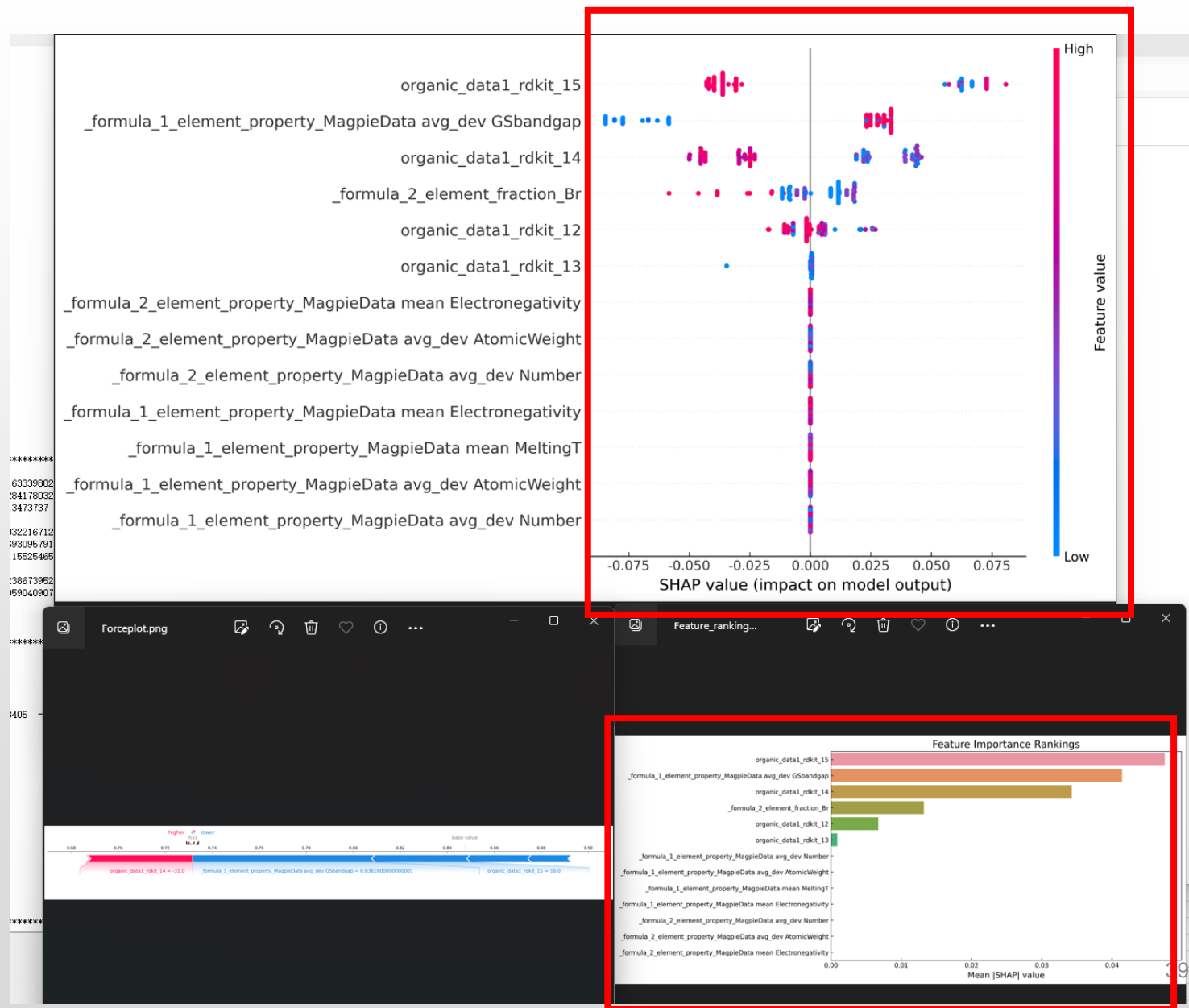
选择保存的机器学习模型.dat



Data import: 特征选择后的 Data\_rfe.csv



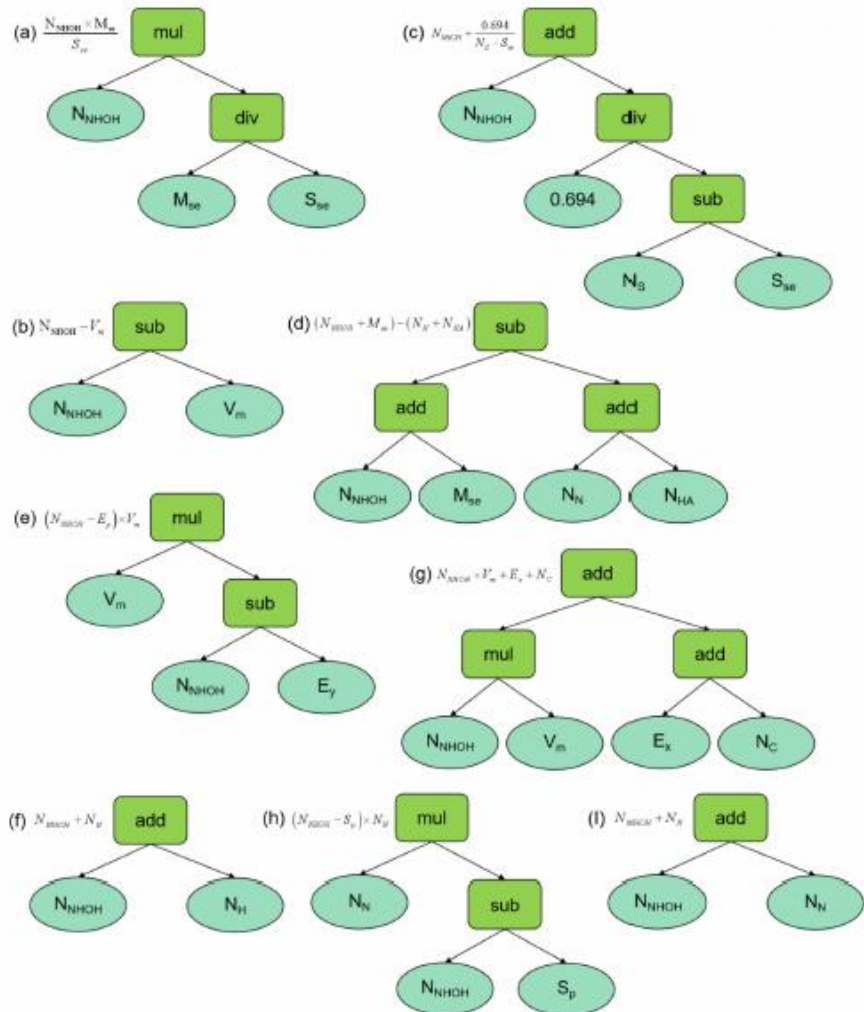
点击result: 自动生成Shapley图



## 8 遗传算法（符号回归/分类）

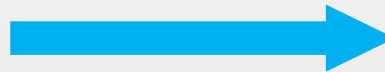


# 遗传算法:Mutation



Hoist Mutation

Point Mutation



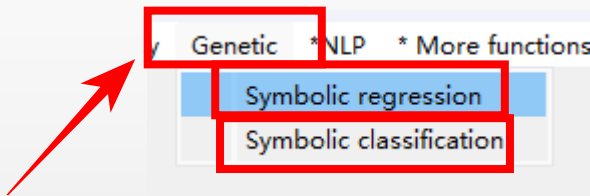
Crossover

Subtree mutation

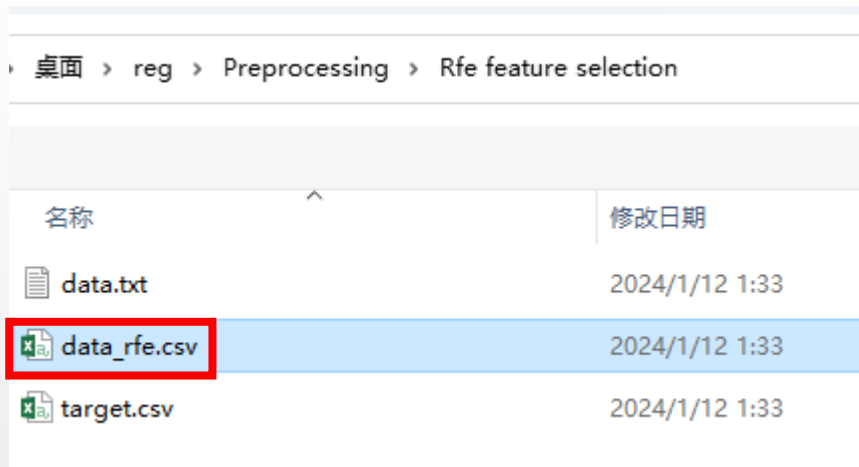
| symbol | mathematical formula  | correlation $r$ |
|--------|---|-----------------|
| F1     | $N_c(E_2^2 - M^3)$  | 0.50            |
| F2     | $N_c(E_2 - M)$  | 0.49            |
| F3     | $\frac{E_2^3}{M^2 + M_{ae}^2}$  | 0.48            |
| F4     | $(N_c^3 - \sqrt{E_2})(E_2 - M)$   | 0.51            |
| F5     | $E_2^2 - R$   | 0.46            |
| F6     | $E_2^2 - N_c M$   | 0.51            |
| F7     | $N_c^2 - N_c(M - E_9^2)$  | 0.51            |
| F8     | $\frac{E_2^2}{E_6^2 \left( M^2 + E_4^2 \frac{N_c^2}{N_{ar}^2} \right)}$ | 0.50            |
| F9     | $E_4^2 \left( E_2^2 + \frac{M^3}{M^2 + N_{ar}} \right)$                 | 0.51            |
| F10    | $2(E_2^3 - M^3 - M_{ae}^3) \ln(E_2)$                                    | 0.50            |
| F11    | $N_c(N_c^2 - M - R)$  | 0.52            |
| F12    | $\frac{B_1 E_2^2}{E_{10}(M^2 - E_4)}$                                   | 0.52            |

# 遗传算法建立模型：导入csv数据集

以特征选择后的Data\_rfe.csv为例



符号回归与符号分类按钮



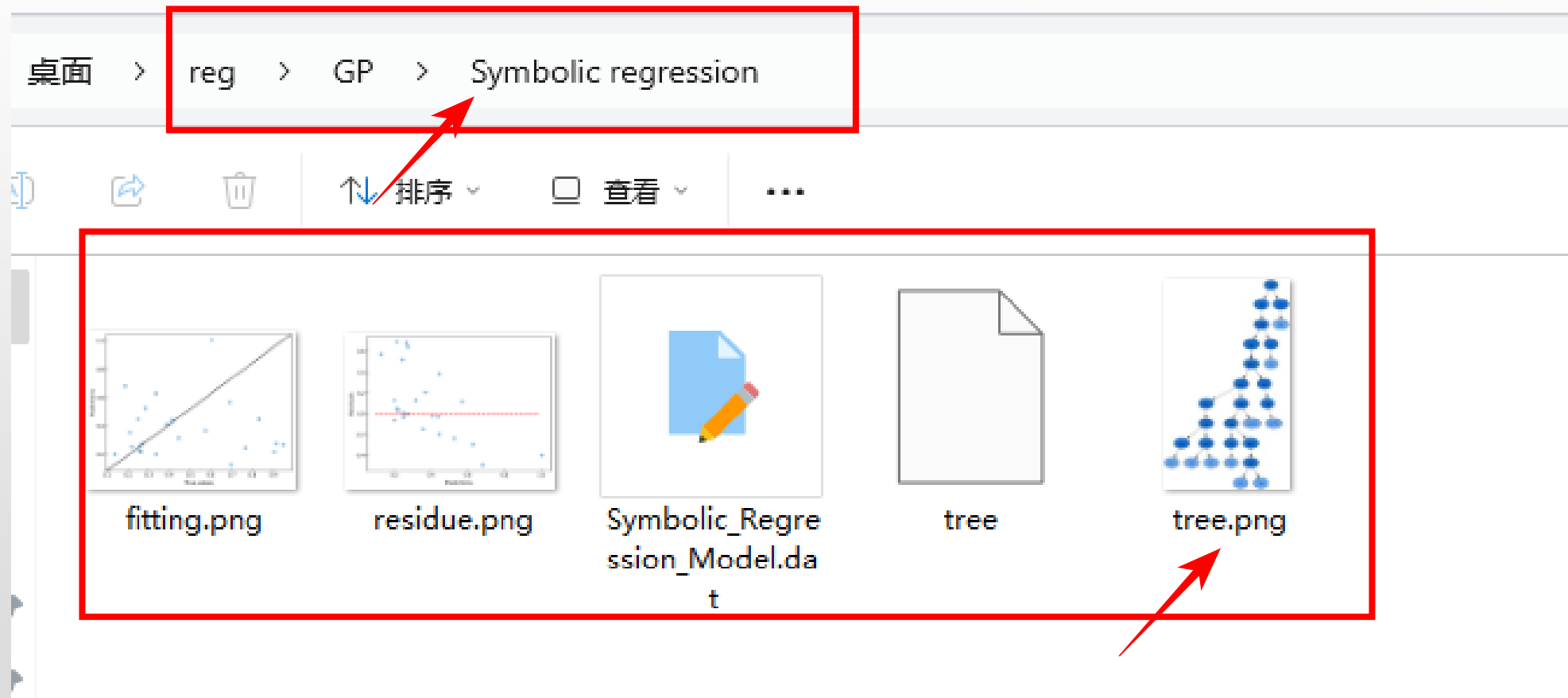
选择之前保存的特征选择好后的csv，例如data\_rfe.csv

| Population Average |        |          | Best Individual |          |             |           |
|--------------------|--------|----------|-----------------|----------|-------------|-----------|
| Gen                | Length | Fitness  | Length          | Fitness  | 00B Fitness | Time Left |
| 0                  | 9.22   | 1.78342  | 5               | 0.245292 | N/A         | 20.65s    |
| 1                  | 5.89   | 0.570097 | 14              | 0.220634 | N/A         | 20.99s    |
| 2                  | 6.23   | 0.480999 | 14              | 0.213933 | N/A         | 20.49s    |
| 3                  | 7.54   | 0.45332  | 9               | 0.209441 | N/A         | 21.24s    |
| 4                  | 9.61   | 0.450839 | 15              | 0.203901 | N/A         | 19.17s    |
| 5                  | 10.74  | 0.361537 | 18              | 0.198648 | N/A         | 18.18s    |
| 6                  | 10.06  | 0.268182 | 21              | 0.189408 | N/A         | 16.53s    |

遗传算法建模进行中！（目前默认超参数不可改，后期待续）

# 遗传算法建立模型：生成数据图

若未自动打开，查看图路径：自定义保存路径→ GP → Symbolic regression （分类图存在Symbolic classification）

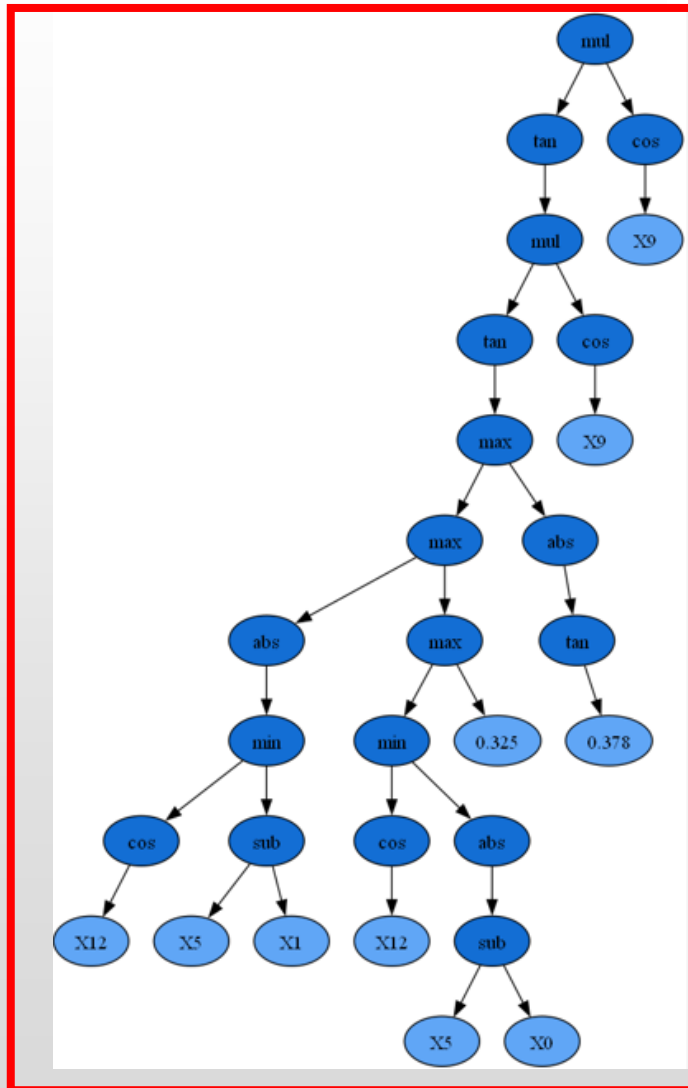
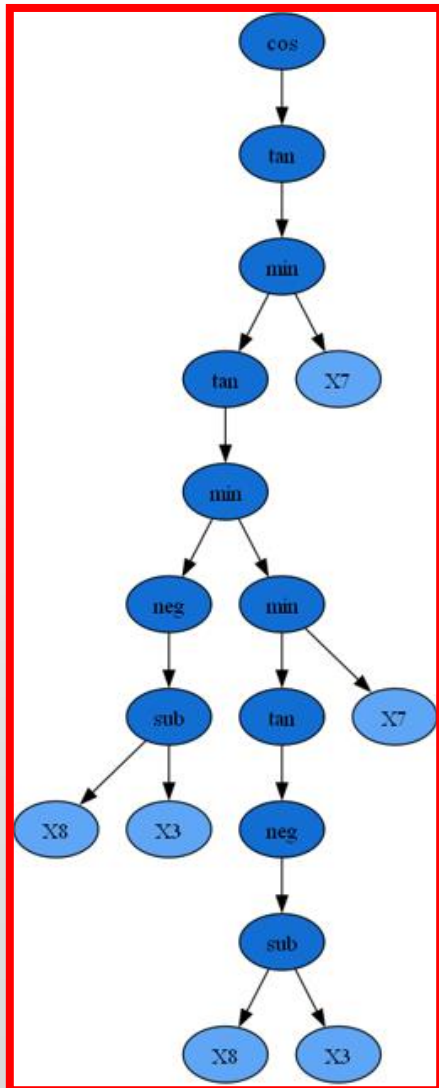


# 遗传算法建立模型：符号树图 Tree

可解释性！

科学！

白匣子！！



$$stability = f_1 + f_2 + f_3 \times f_4$$

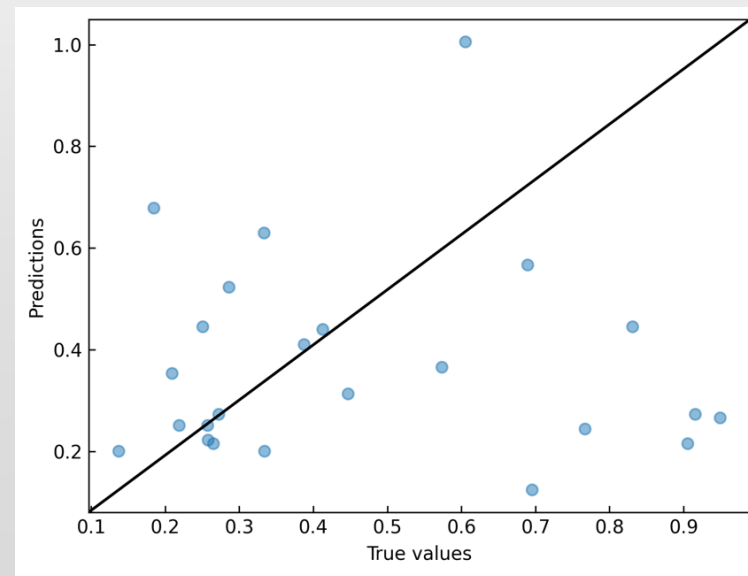
$$f_1 = R_A^2 + a + E_S + R_D$$

$$f_2 = \frac{D_{Ma}^2(E_B + A_S + b)}{h + b} + \frac{f}{h} + \frac{R_C}{D_{Fp}} + D_{Ma}(2d + g + 2C_{Me} + 2b)$$

$$f_3 = R_A^2 + a + \frac{f D_{Ma}}{b} + \frac{D_{Ma}(d + b)}{D_{Ma} + d + \frac{f}{d}}$$

$$f_4 = D_{Ma}(d + g + 2C_{Me}) + D_{Ma}^2(C_{Me} + b + A_S + 1) + \frac{f}{d} + d + C_S$$

$$a = \frac{R_A R_D}{R_A + E_S}, b = \frac{R_D}{D_{Ml}}, d = C_{Me} + R_C, f = D_{Fp} + D_S, g = \frac{D_{Fp}}{E_B + b}, h = \frac{D_{Fp} + D_S}{E_b + b}$$



未完待续（其它软件功能

**NLP/Featurizer/Algorithms/DFT/SLME/CSP**

**/Linux/CGCNN等）**

# Acknowledgement and Conclusions



- 材料设计中小数据问题仍然严重（无论模拟和实验：成本高，欠拟合，过拟合）
- 各类数据来源和算法各有利弊
- Natural language processing (自然语言处理) facilitates materials prediction
- 遗传算法Genetic programming帮助科学解释（公式/decouple解耦）
- 基于文本描述符的多模材料设计方法：准确率、成本与可解释性的较好均衡！
- NJmat软件开发（描述符与算法）



# Thanks !