



南京大学  
NANJING UNIVERSITY

工程管理学院  
SCHOOL OF MANAGEMENT & ENGINEERING

# 数据分割与强规则过滤

汇报人：杨超凡

小组成员：王一伦、钱创杰、温子祺、陈波

# 目录

## CONTENTS

1

工作回顾

2

强规则过滤

3

数据分割

4

优化效果与总结



以 15 号风机各变量方差的百分之二十五分位数作为阈值进行过滤，剔除掉方差小于阈值的变量，保留21 个显著变量。

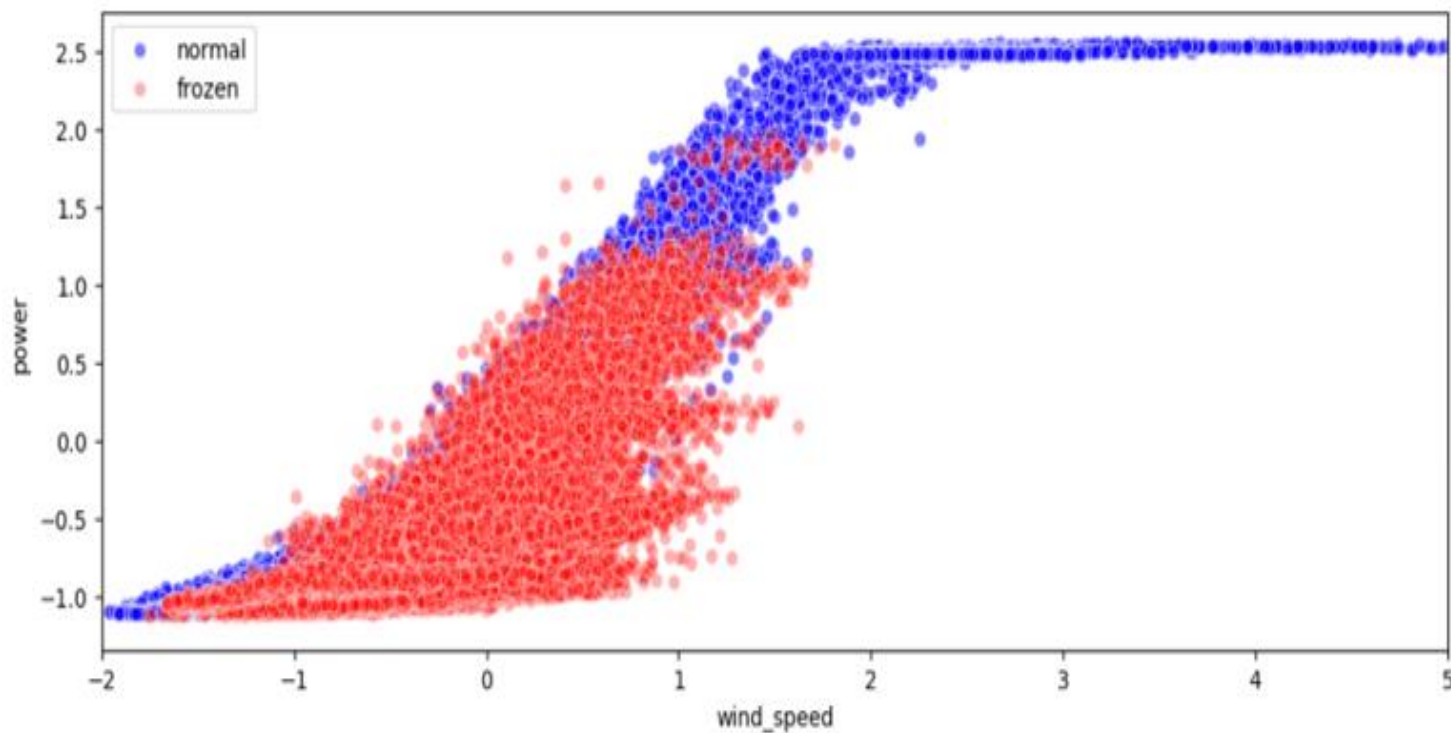
```
['wind_speed',  
 'generator_speed',  
 'power',  
 'wind_direction',  
 'wind_direction_mean',  
 'yaw_position',  
 'pitch1_angle',  
 'pitch2_angle',  
 'pitch3_angle',  
 'pitch1_moto_tmp',  
 'pitch2_moto_tmp',  
 'pitch3_moto_tmp',  
 'acc_x',  
 'acc_y',  
 'environment_tmp',  
 'int_tmp',  
 'pitch1_ng5_DC',  
 'pitch2_ng5_DC',  
 'pitch3_ng5_DC',  
 'group',  
 'timestamp']
```

如果属性变量太多容易导致过拟合，为提高模型泛化能力通过卡方检验进一步筛选，最后保留16个原始变量。

```
['wind_speed',  
 'generator_speed',  
 'power',  
 'wind_direction',  
 'wind_direction_mean',  
 'yaw_position',  
 'pitch1_angle',  
 'pitch1_moto_tmp',  
 'pitch2_moto_tmp',  
 'pitch3_moto_tmp',  
 'acc_x',  
 'acc_y',  
 'environment_tmp',  
 'int_tmp',  
 'pitch1_ng5_DC',  
 'group']
```

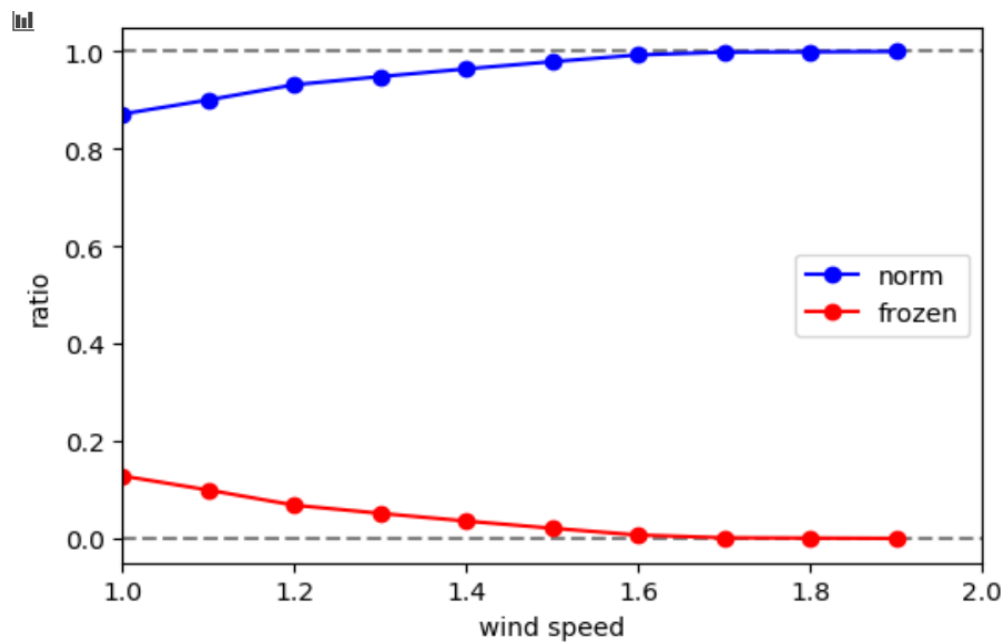
Model	RandomForest	DecisionTreeClassifier	GBDT	GBDT+LR	RF+LR
Cross_val_score	0.756	0.808	0.844	0.899	0.869

## 基于特征属性

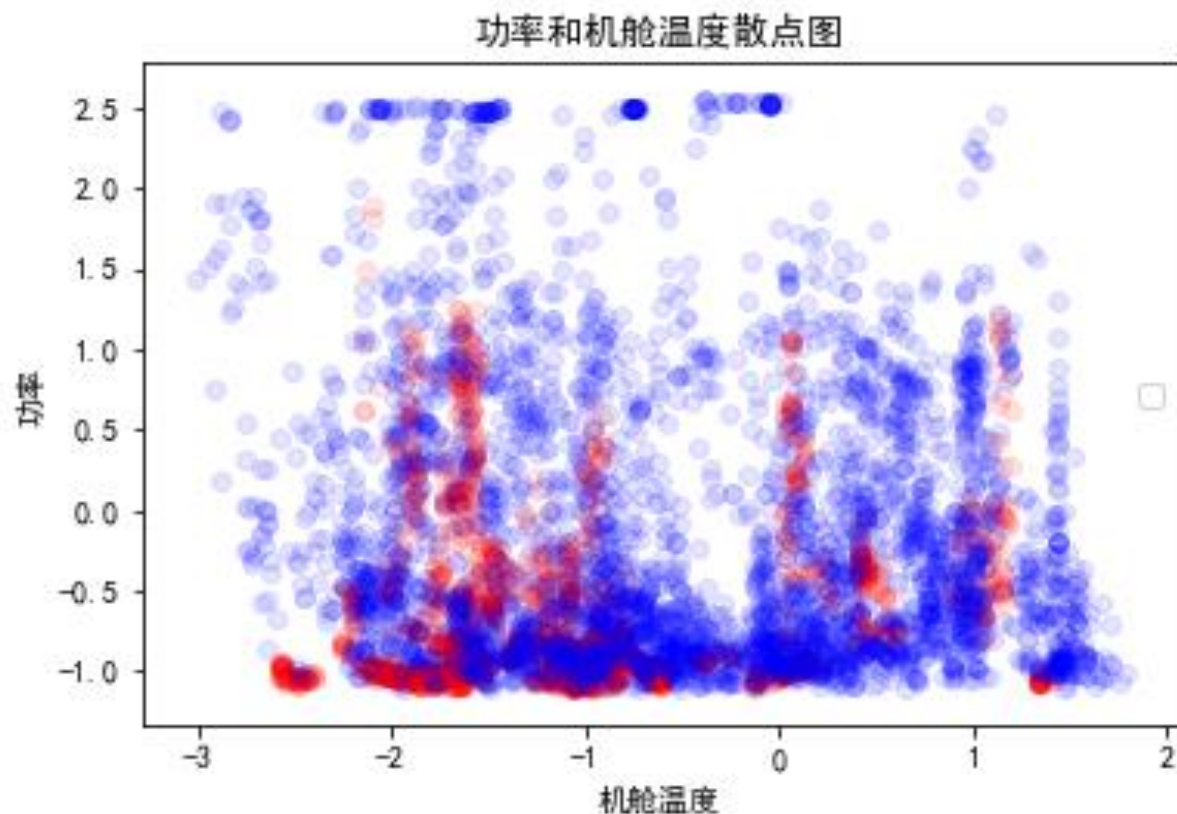


可以看出在风速和功率 $<2$ 的范围内，模型对两类样本的区分度不大。



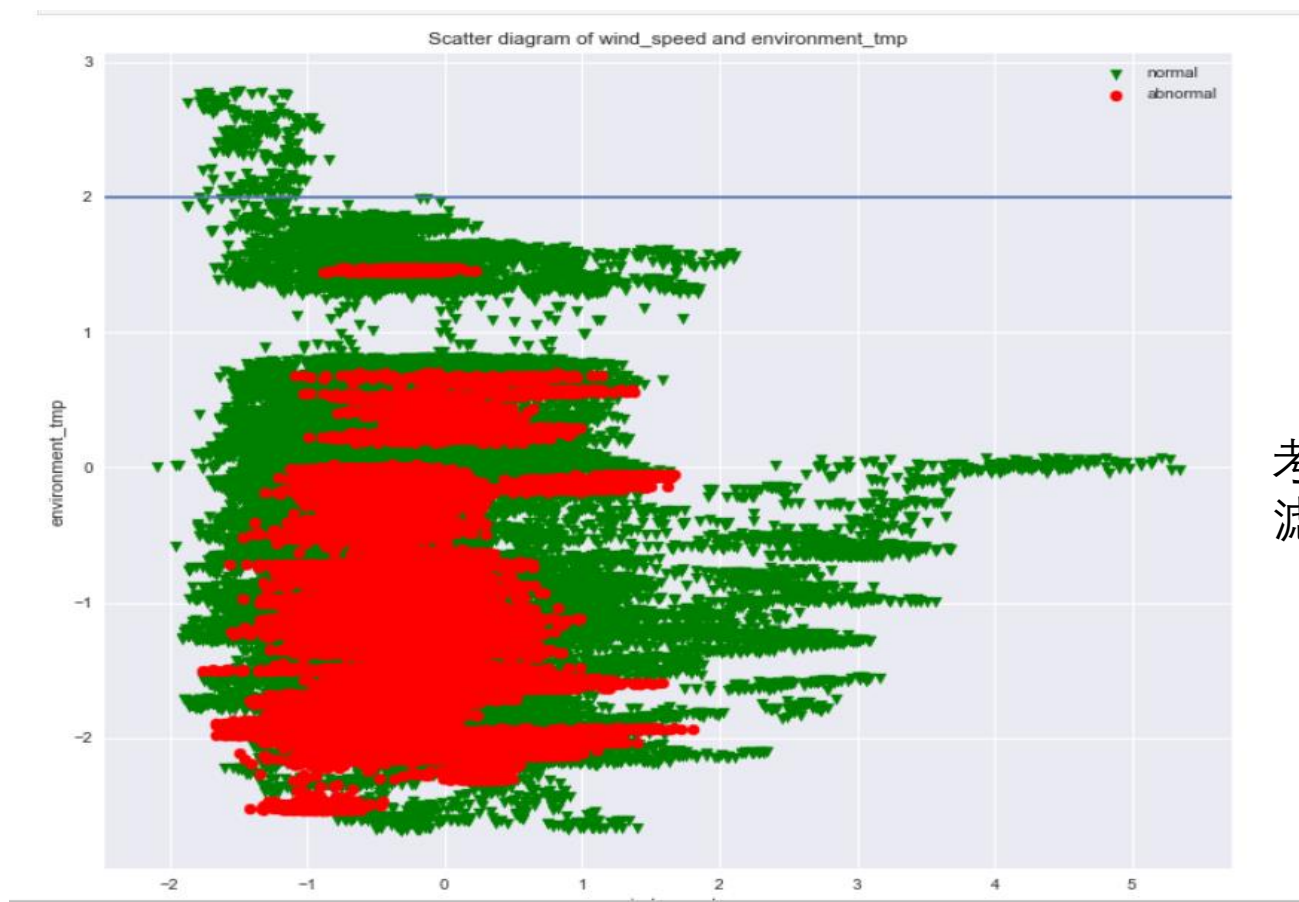


当风速较大时，风机为正常的比重较高，当风速为1.以上的时候正常风机的比例为100%，认为可剔除风速>1.9的数据。



当机舱温度大于1.5或小于-2.5时，数据基本为不结冰数据，可以考虑将“`int_tmp>1.5 & int_tmp<-2.5`”作为强规则对数据进行过滤。





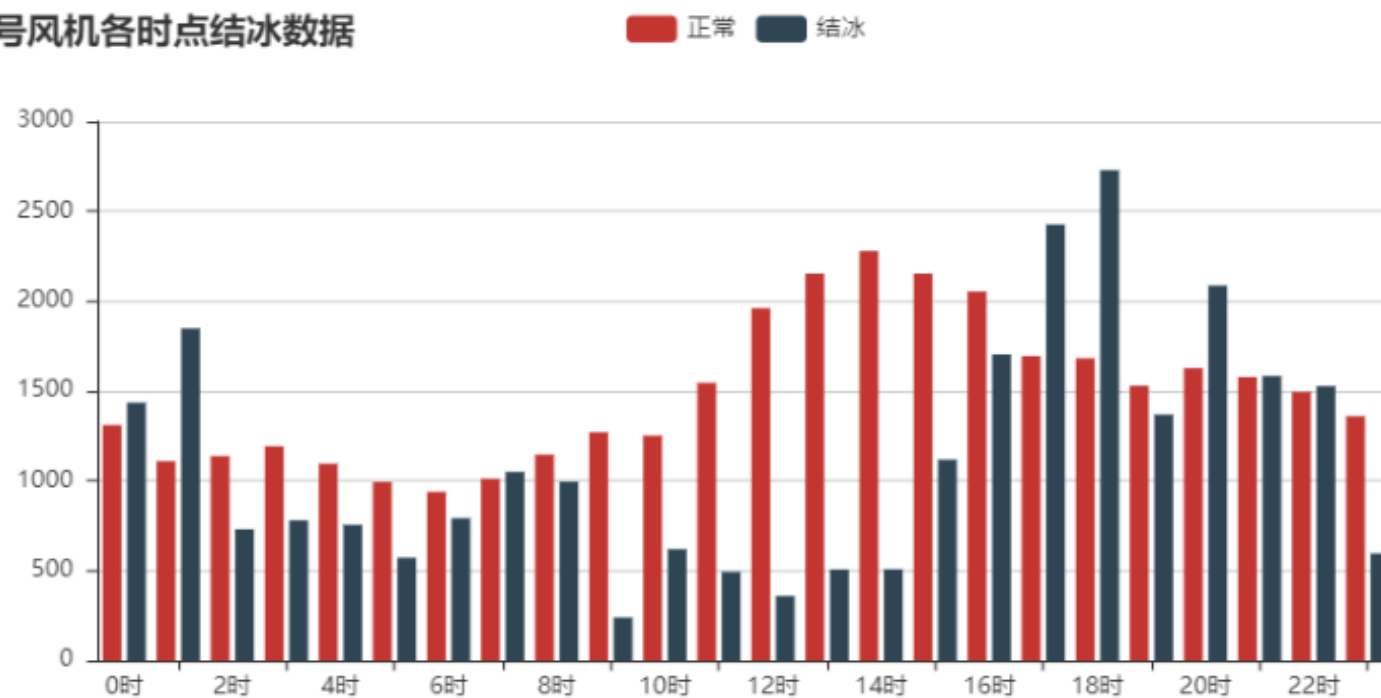
当环境温度大于2时，数据基本为不结冰数据，考虑将“environment\_tmp>2”作为强规则对数据过滤。

综上所述，构建出四条强规则对数据进行过滤，保留满足以下四个条件的数据

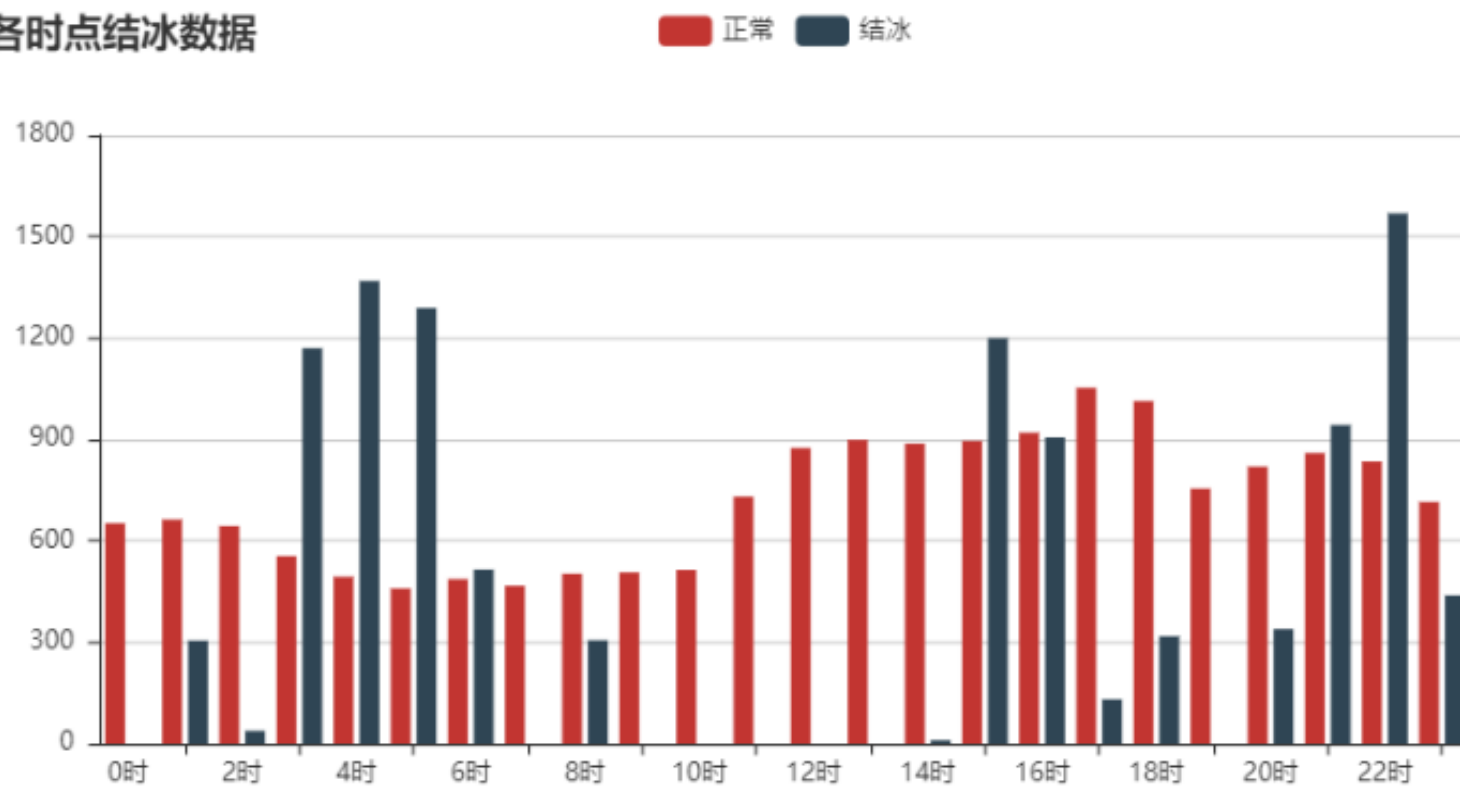
1.  $\text{power} < 2$
2.  $\text{wind\_speed} < 1.9$
3.  $\text{environment\_tmp} < 2$
4.  $-2.5 < \text{int\_tmp} < 1.5$

## 基于时间节点

15号风机各时点结冰数据



21各时点结冰数据



综上所述，我们尝试按数据时间点是否属于9-14时段对数据分割  
两种策略：把period作为标签属性作为特征学习或分割后各自训练

```
data15['period']=1
data21['period']=1
for i in range(len(data15)):
    if data15.loc[i,'time']<9 or data15.loc[i,'time']>14:
        data15.loc[i,'period']=0
for i in range(len(data21)):
    if data21.loc[i,'time']<9 or data21.loc[i,'time']>14:
        data21.loc[i,'period']=0
```

## 强规则过滤优化效果

```
data15=data15[(data15['power']<2) & (data15['wind_speed']<1.9) & (data15['environment_tmp']<2)&(data15['int_tmp']<1.5)&(data15['int_tmp']>-2.5)]  
data21=data21[(data21['power']<2) & (data21['wind_speed']<1.9) & (data21['environment_tmp']<2)&(data21['int_tmp']<1.5)&(data21['int_tmp']>-2.5)]
```

data15.shape

(62391, 33)

data21.shape

(28053, 33)

过滤前

data15.shape

(59109, 32)

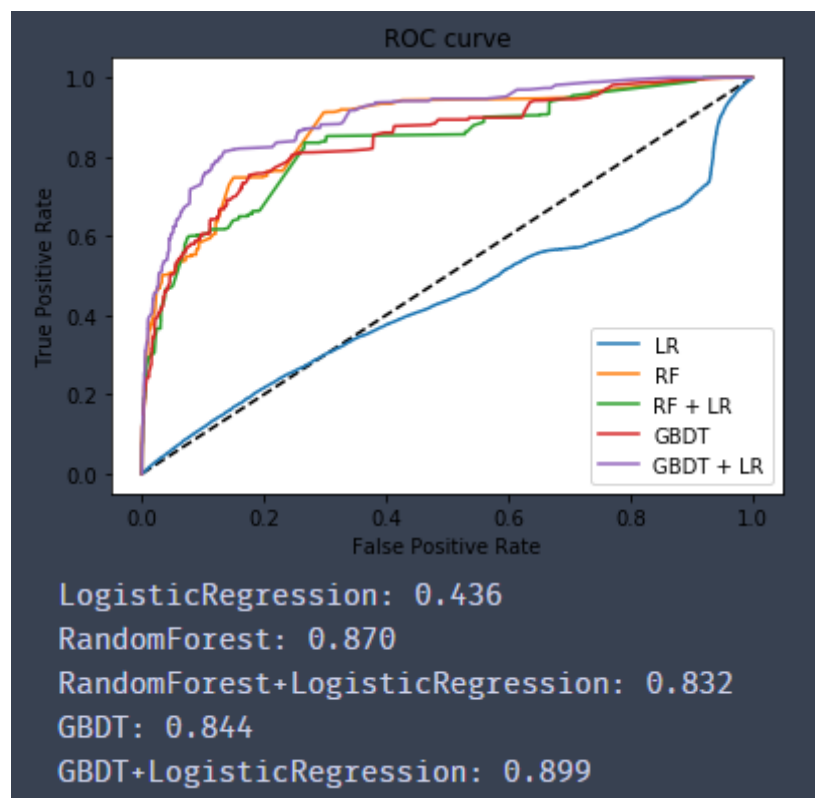
data21.shape

(24493, 32)

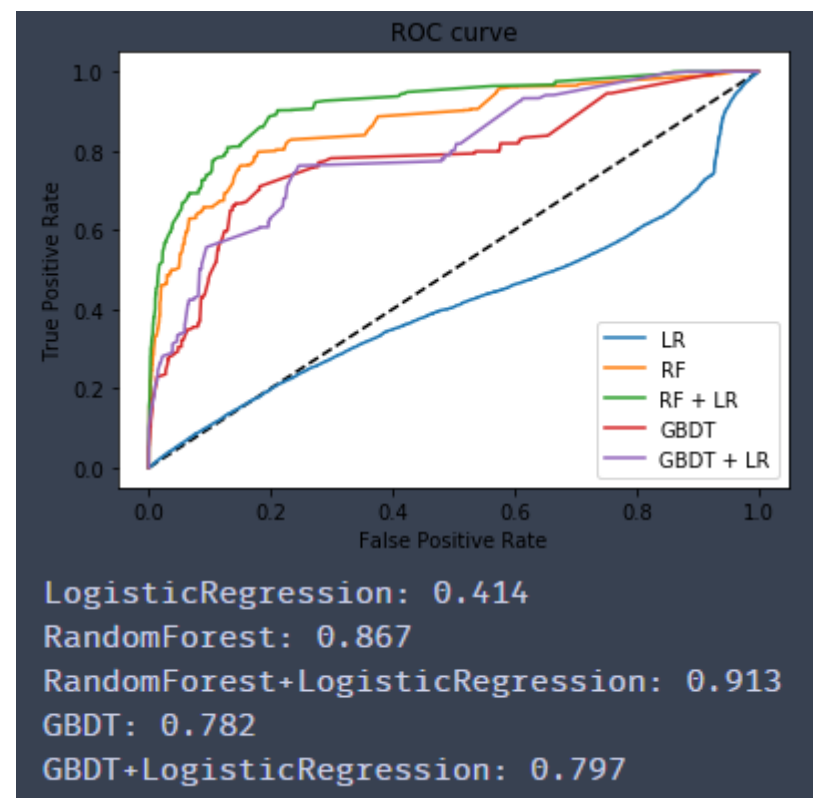
过滤后



## 强规则过滤优化效果

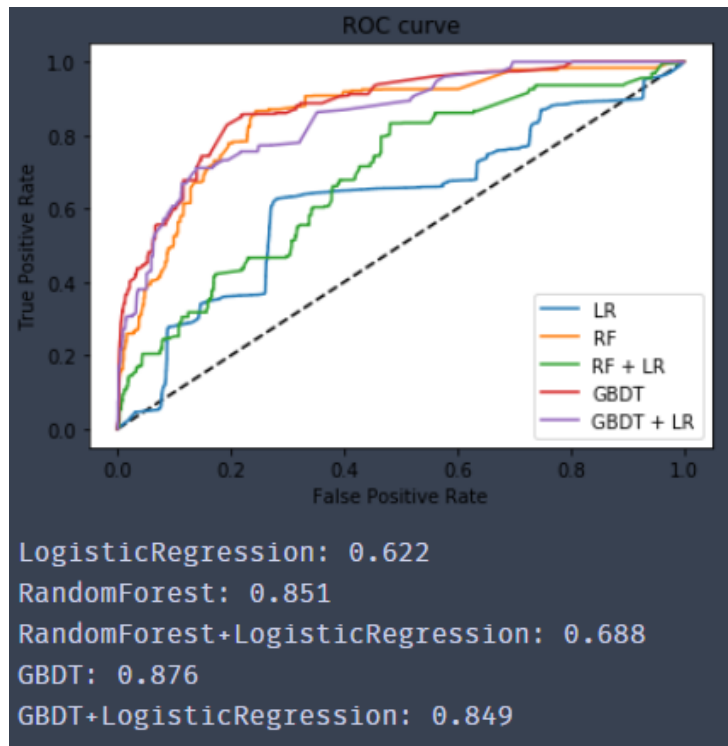


过滤前

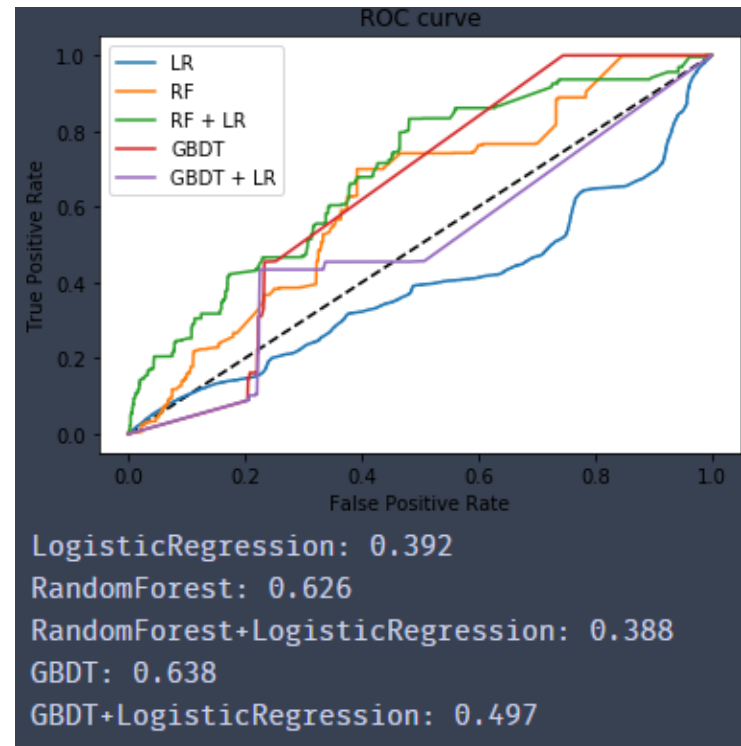


过滤后

## 优化效果（切割后预测）

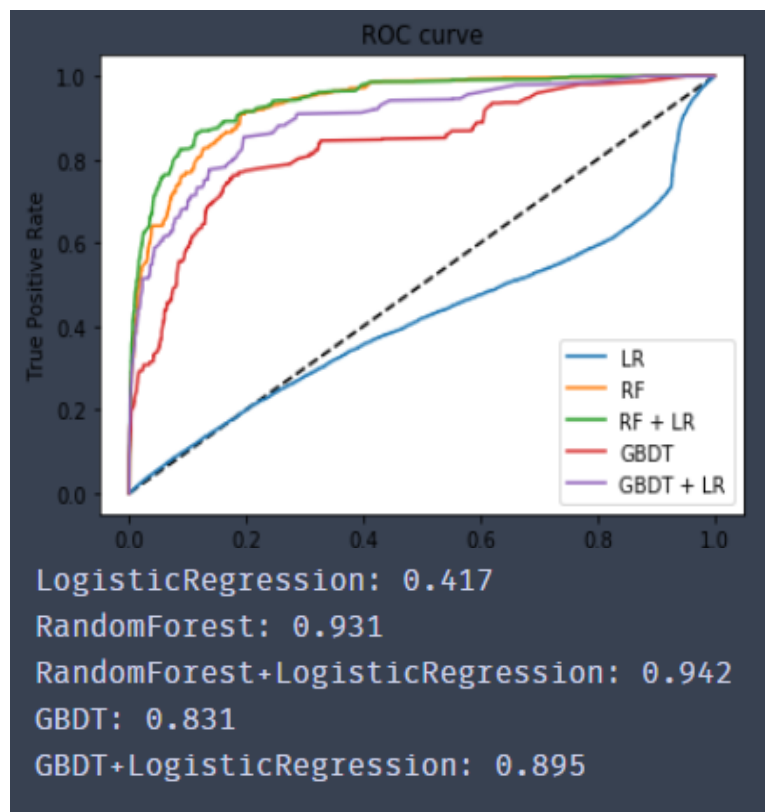


Period=0

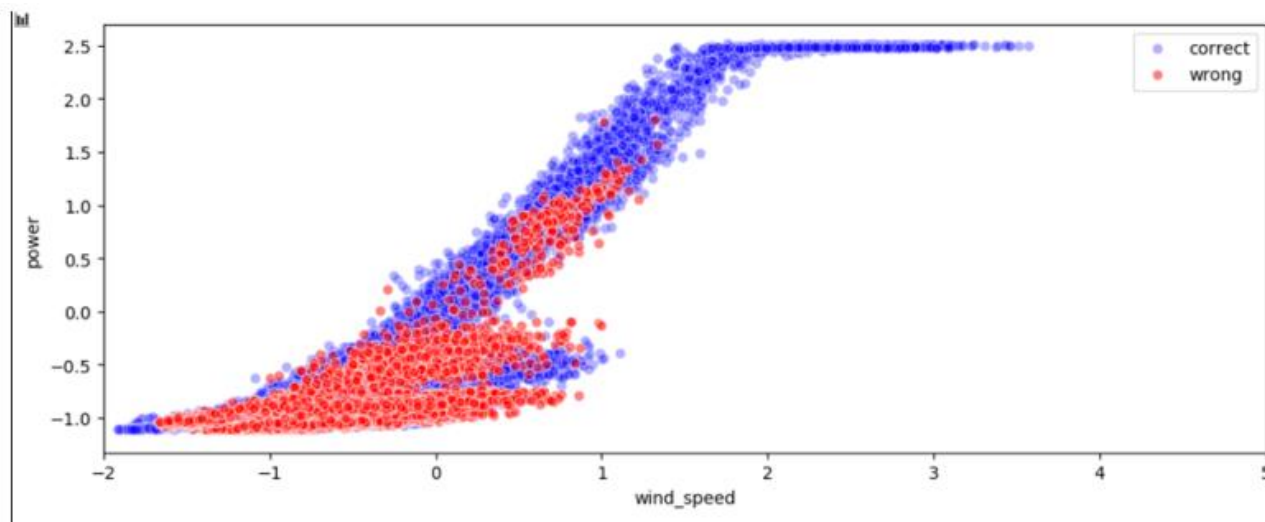


Period=1

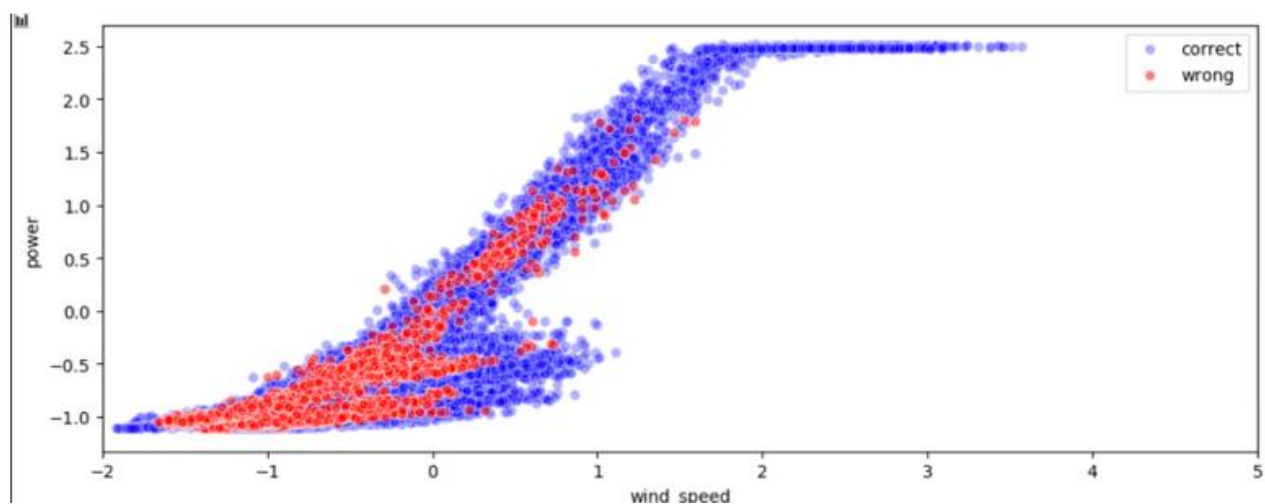
优化效果（把时段作为特征属性）



从数据分割前后模型在训练集和测试集上的结果来看，通过强规则过滤和引入新的二分类变量可以一定程度上提高模型的预测能力，提升模型的泛化能力。



Before

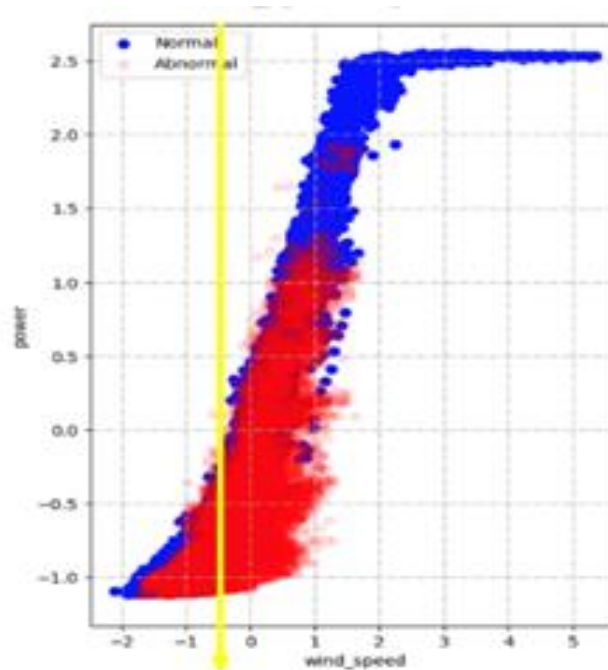


After

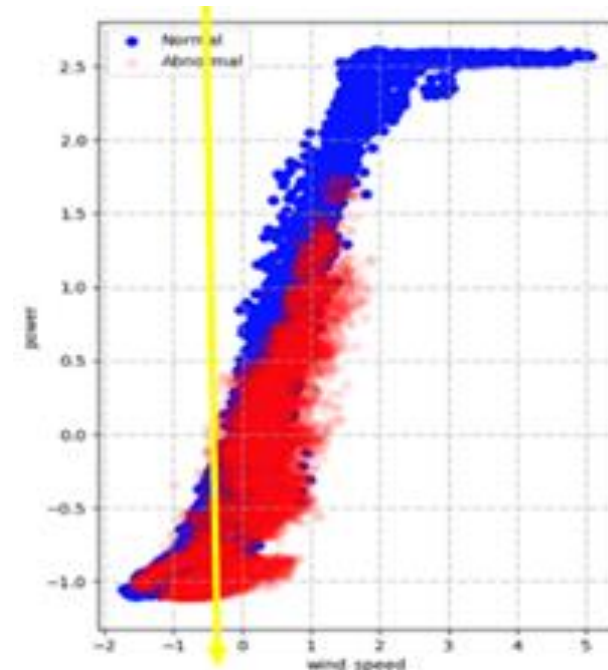
通过数据分割，我们的模型在15号风机 $\text{wind\_speed} < -0.5$ 时取得了85.6%的准确率， $\text{wind\_speed} > -0.5$ 时取得了83.1%的准确率；在21号测试集 $\text{wind\_speed} < -0.5$ 时取得了79.3%的准确率， $\text{wind\_speed} > -0.5$ 时取得了73.8%的准确率。

模型的预测能力有了明显的下降，我们猜测是因为不同风机的数据强分割界线可能有所不同，我们要选取不同的分割标准。

另外，我们接下来还可以对其他特征考虑进行数据分割，并进一步结合强规则过滤、特征选择来综合提升模型的能力。



15号风机 $\text{wind\_speed}$ 与 $\text{power}$ 关系



21号风机 $\text{wind\_speed}$ 与 $\text{power}$ 关系



工程管理学院  
SCHOOL OF MANAGEMENT & ENGINEERING

# 非常感谢您的观看