

10/21/2018

University of Illinois at Urbana-Champaign

IE598 - Machine Learning in Finance

# IE598 - MLF FINAL PROJECT

FALL 2018

*by Yuchen Duan, Ruozhong Yang, Fengkai Xu, Biao Feng, and Joseph Loss*

# TABLE OF CONTENTS

## Contents

|  |    |
|--|----|
| Chapter 1: Moody's Bond Rating Classifier      | 1  |
| Exploratory Data Analysis                      | 1  |
| Preprocessing & Feature Extraction/Selection   | 2  |
| Model Fitting & Evaluation(Binary&muticlasses) | 3  |
| Hyperparameter Tuning                          | 3  |
| Ensembling                                     | 4  |
| Conclusions                                    | 4  |
| Chapter 2: USPHCI Economic Activity Forecast   | 5  |
| Exploratory Data Analysis                      | 5  |
| Preprocessing & Feature Extraction/Selection   | 6  |
| Model Fitting & Evaluation                     | 8  |
| Hyperparameter Tuning                          | 8  |
| Ensembling                                     | 11 |
| Conclusions                                    | 11 |
| Appendix                                       | 12 |
| Github Repository                              | 12 |

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

## Chapter 1: Moody's Bond Rating Classifier

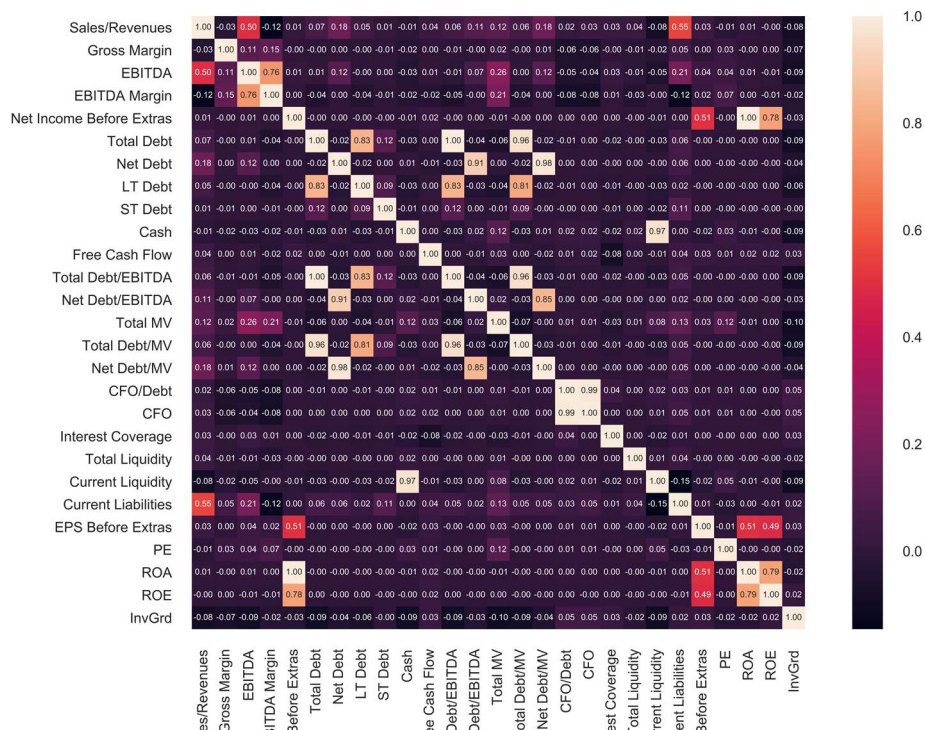
### EXPLORATORY DATA ANALYSIS

Here is what our data looks like:

```
RangeIndex: 1700 entries, 0 to 1699
Data columns (total 29 columns):
Sales/Revenues          1700 non-null float64
Gross Margin            1700 non-null float64
EBITDA                  1700 non-null float64
EBITDA Margin           1700 non-null float64
Net Income Before Extras 1700 non-null float64
Total Debt              1700 non-null float64
Net Debt                1700 non-null float64
LT Debt                 1700 non-null float64
ST Debt                 1700 non-null float64
Cash                    1700 non-null float64
Free Cash Flow          1700 non-null float64
Total Debt/EBITDA       1700 non-null float64
Net Debt/EBITDA         1700 non-null float64
Total MV                1700 non-null float64
Total Debt/MV           1700 non-null float64
Net Debt/MV             1700 non-null float64
CFO/Debt                1700 non-null float64
CFO                     1700 non-null float64
Interest Coverage        1700 non-null float64
Total Liquidity          1700 non-null float64
Current Liquidity        1700 non-null float64
Current Liabilities      1700 non-null float64
EPS Before Extras        1700 non-null float64
PE                       1700 non-null float64
ROA                     1700 non-null float64
ROE                     1700 non-null float64
InvGrd                  1700 non-null int64
Rating                  1700 non-null object
Class                   1700 non-null int64
dtypes: float64(26), int64(2), object(1)
memory usage: 385.2+ KB
```

Also, we have a correlation matrix:

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。



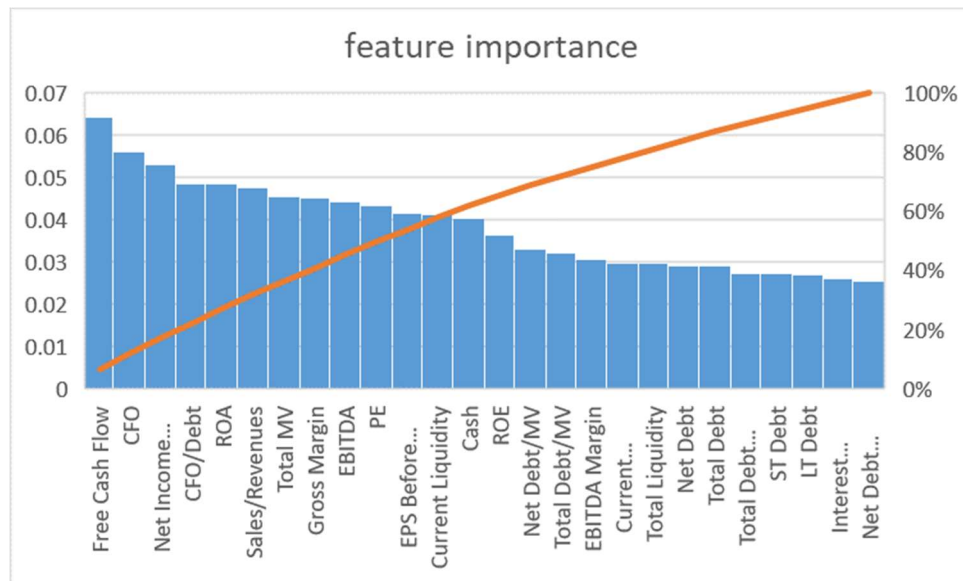
## PREPROCESSING & FEATURE EXTRACTION/SELECTION

The preprocessing part combine some steps that need to be done before we try to fit our model:

1. Split the test and train database via `train_test_split`(with `test_size = 0.1` and `random_state=42`)
2. Standardize features via `StandardScaler` for better model performance.

We also calculate the importance of each feature and select 13 of them for our models.

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。



## MODEL FITTING & EVALUATION(BINARY&MUTICLASSES)

### 1. Model 1

The first model is the KNN model.

### 2. Model 2

The second model is the Random Forest model.

### 3. Model 3

The third model is the Decision Tree model.

### 4. Model 4

The forth model is the Logistic Regression model.

We will discuss those models in the hyperparameter tuning and ensemble parts.

## HYPERPARAMETER TUNING

We deal with different parameters via GridSearchCV function, the range of each model's parameter is form 1 to 100. Here is the best result for each model:

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

| binary              |             | muticlasses         |             | percents |
|---------------------|-------------|---------------------|-------------|----------|
| KNN                 | 0.8         | KNN                 | 0.458823529 | 0.573529 |
| RandomForest        | 0.858823529 | RandomForest        | 0.676470588 | 0.787671 |
| Decision tree       | 0.794117647 | Decision tree       | 0.447058824 | 0.562963 |
| Logistic Regression | 0.741176471 | Logistic Regression | 0.247058824 | 0.333333 |

From the table, it is easy to find the muticlasses task lead to poor prediction(mutl\_lr score is about 1/3 compare to the binary one). There are several improvements can be done for better models, we will discuss them at the conclusion.

## ENSEMBLING

Our team use the ensembling method for binary classification does not support muticlassification task. Result showed below:

| binary   |               |                   |
|----------|---------------|-------------------|
| ROC AUC: | 0.73(+/-0.05) | [KNN]             |
| ROC AUC: | 0.9(+/-0.02)  | [RandomForest]    |
| ROC AUC: | 0.75(+/-0.05) | [Decision tree]   |
| ROC AUC: | 0.89(+/-0.02) | [Majority voting] |

## CONCLUSIONS

The best result for binary model is 0.89(after ensembling) and the best for multiclass is 0.67. There are several things we can do to improve our model:

### 1. Dimension reduction

We can reduce the dimension of our model for better prediction, but may let miss some important information.

### 2. Internal relationships

Some features are highly correlated, we can find them and just use one of them. Besides, many features have internal relationships, thus, some of them may actually talk about the same thing.

### 3. Weight adjustment

Although those models adjust weights of each feature automatically, people from accounting major may hold different view of those weights.

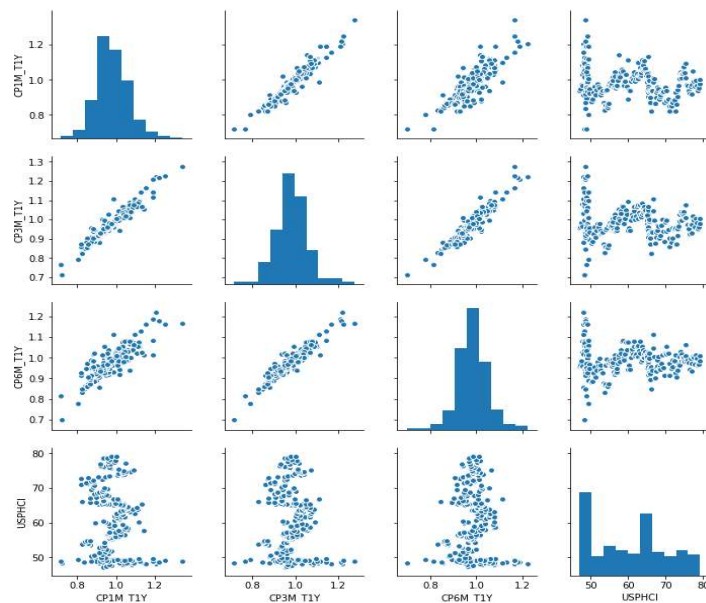
错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

## Chapter 2: USPHCI Economic Activity Forecast

### EXPLORATORY DATA ANALYSIS

Our data looks like this:

```
Int64Index: 223 entries, 0 to 222
Data columns (total 16 columns):
T1Y Index      223 non-null float64
T2Y Index      223 non-null float64
T3Y Index      223 non-null float64
T5Y Index      223 non-null float64
T7Y Index      223 non-null float64
T10Y Index     223 non-null float64
CP1M           223 non-null float64
CP3M           223 non-null float64
CP6M           223 non-null float64
CP1M_T1Y       223 non-null float64
CP3M_T1Y       223 non-null float64
CP6M_T1Y       223 non-null float64
USPHCI         223 non-null float64
PCT 3MO FWD    223 non-null float64
PCT 6MO FWD    223 non-null float64
PCT 9MO FWD    223 non-null float64
dtypes: float64(16)
memory usage: 29.6 KB
```



Also, we need to see the relation between each features:

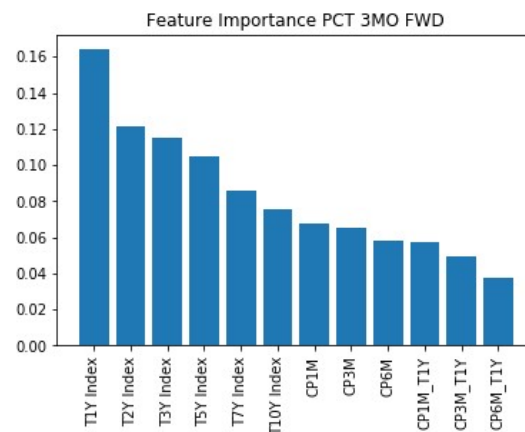


错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

|             |           |           |           |           |           |            |      |      |      |          |          |          |        |             |             |             |
|-------------|-----------|-----------|-----------|-----------|-----------|------------|------|------|------|----------|----------|----------|--------|-------------|-------------|-------------|
| T1Y Index   | 1.0       | 1.0       | 1.0       | 1.0       | 0.9       | 0.9        | 1.0  | 1.0  | 1.0  | 0.2      | 0.2      | 0.0      | -0.8   | -0.4        | -0.5        | -0.5        |
| T2Y Index   | 1.0       | 1.0       | 1.0       | 1.0       | 1.0       | 1.0        | 0.9  | 0.9  | 1.0  | 0.1      | 0.1      | -0.0     | -0.8   | -0.4        | -0.4        | -0.4        |
| T3Y Index   | 1.0       | 1.0       | 1.0       | 1.0       | 1.0       | 1.0        | 0.9  | 0.9  | 0.9  | 0.1      | 0.1      | -0.1     | -0.8   | -0.4        | -0.4        | -0.4        |
| T5Y Index   | 1.0       | 1.0       | 1.0       | 1.0       | 1.0       | 1.0        | 0.9  | 0.9  | 0.9  | 0.1      | 0.0      | -0.1     | -0.8   | -0.4        | -0.4        | -0.4        |
| T7Y Index   | 0.9       | 1.0       | 1.0       | 1.0       | 1.0       | 1.0        | 0.9  | 0.9  | 0.9  | 0.0      | 0.0      | -0.1     | -0.8   | -0.3        | -0.4        | -0.4        |
| T10Y Index  | 0.9       | 1.0       | 1.0       | 1.0       | 1.0       | 1.0        | 0.9  | 0.9  | 0.9  | 0.0      | -0.0     | -0.1     | -0.8   | -0.3        | -0.4        | -0.4        |
| CP1M        | 1.0       | 0.9       | 0.9       | 0.9       | 0.9       | 0.9        | 1.0  | 1.0  | 1.0  | 0.5      | 0.4      | 0.2      | -0.7   | -0.4        | -0.5        | -0.5        |
| CP3M        | 1.0       | 0.9       | 0.9       | 0.9       | 0.9       | 0.9        | 1.0  | 1.0  | 1.0  | 0.4      | 0.4      | 0.2      | -0.7   | -0.4        | -0.5        | -0.5        |
| CP6M        | 1.0       | 1.0       | 0.9       | 0.9       | 0.9       | 0.9        | 1.0  | 1.0  | 1.0  | 0.4      | 0.4      | 0.2      | -0.8   | -0.4        | -0.5        | -0.5        |
| CP1M_T1Y    | 0.2       | 0.1       | 0.1       | 0.1       | 0.0       | 0.0        | 0.5  | 0.4  | 0.4  | 1.0      | 1.0      | 0.8      | -0.1   | -0.2        | -0.2        | -0.3        |
| CP3M_T1Y    | 0.2       | 0.1       | 0.1       | 0.0       | 0.0       | -0.0       | 0.4  | 0.4  | 0.4  | 1.0      | 1.0      | 0.9      | -0.1   | -0.1        | -0.2        | -0.3        |
| CP6M_T1Y    | 0.0       | -0.0      | -0.1      | -0.1      | -0.1      | -0.1       | 0.2  | 0.2  | 0.2  | 0.8      | 0.9      | 1.0      | 0.0    | 0.0         | -0.1        | -0.2        |
| USPHCI      | -0.8      | -0.8      | -0.8      | -0.8      | -0.8      | -0.8       | -0.7 | -0.7 | -0.8 | -0.1     | -0.1     | 0.0      | 1.0    | 0.2         | 0.2         | 0.2         |
| PCT 3MO FWD | -0.4      | -0.4      | -0.4      | -0.4      | -0.3      | -0.3       | -0.4 | -0.4 | -0.4 | -0.2     | -0.1     | 0.0      | 0.2    | 1.0         | 0.9         | 0.9         |
| PCT 6MO FWD | -0.5      | -0.4      | -0.4      | -0.4      | -0.4      | -0.4       | -0.5 | -0.5 | -0.5 | -0.2     | -0.2     | -0.1     | 0.2    | 0.9         | 1.0         | 1.0         |
| PCT 9MO FWD | -0.5      | -0.4      | -0.4      | -0.4      | -0.4      | -0.4       | -0.5 | -0.5 | -0.5 | -0.3     | -0.3     | -0.2     | 0.2    | 0.9         | 1.0         | 1.0         |
|             | T1Y Index | T2Y Index | T3Y Index | T5Y Index | T7Y Index | T10Y Index | CP1M | CP3M | CP6M | CP1M_T1Y | CP3M_T1Y | CP6M_T1Y | USPHCI | PCT 3MO FWD | PCT 6MO FWD | PCT 9MO FWD |

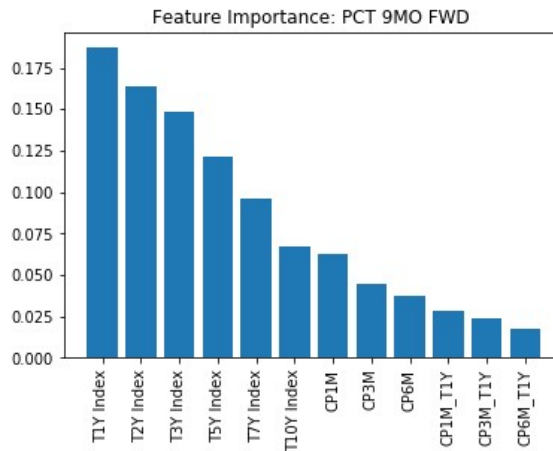
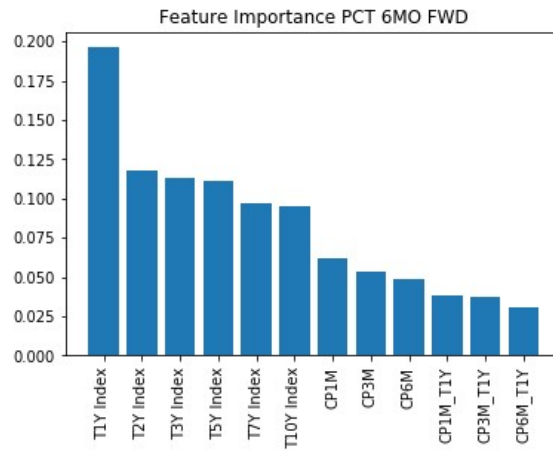
## PREPROCESSING & FEATURE EXTRACTION/SELECTION

We can see the importance of each feature in all three situations:





错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。



#### 3MO FWD RATE - Feature Importance

|               |          |
|---------------|----------|
| 1) T1Y Index  | 0.163890 |
| 2) CP1M_T1Y   | 0.121097 |
| 3) T10Y Index | 0.114853 |
| 4) T3Y Index  | 0.104408 |
| 5) T2Y Index  | 0.085722 |
| 6) CP1M       | 0.075296 |
| 7) CP3M       | 0.067459 |
| 8) CP6M_T1Y   | 0.065129 |
| 9) T5Y Index  | 0.057837 |
| 10) T7Y Index | 0.057455 |
| 11) CP6M      | 0.049308 |
| 12) CP3M_T1Y  | 0.037545 |

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

| 6MO FWD RATE - Feature Importance |          |
|-----------------------------------|----------|
| 1) T1Y Index                      | 0.195985 |
| 2) CP1M                           | 0.117591 |
| 3) CP3M                           | 0.112635 |
| 4) T10Y Index                     | 0.110856 |
| 5) CP1M_T1Y                       | 0.096567 |
| 6) CP6M                           | 0.095394 |
| 7) T3Y Index                      | 0.061621 |
| 8) T5Y Index                      | 0.053676 |
| 9) T7Y Index                      | 0.048926 |
| 10) T2Y Index                     | 0.038705 |
| 11) CP6M_T1Y                      | 0.037129 |
| 12) CP3M_T1Y                      | 0.030916 |

| 9MO FWD RATE - Feature Importance |          |
|-----------------------------------|----------|
| 1) CP1M                           | 0.186953 |
| 2) CP3M                           | 0.164119 |
| 3) CP6M                           | 0.148786 |
| 4) T10Y Index                     | 0.121164 |
| 5) T1Y Index                      | 0.095936 |
| 6) CP1M_T1Y                       | 0.067408 |
| 7) T7Y Index                      | 0.063060 |
| 8) T5Y Index                      | 0.044385 |
| 9) T3Y Index                      | 0.037779 |
| 10) CP6M_T1Y                      | 0.028688 |
| 11) CP3M_T1Y                      | 0.024268 |
| 12) T2Y Index                     | 0.017453 |

## MODEL FITTING & EVALUATION

### 1. Model 1

We use Linear Regression for 3-month situation.

### 2. Model 2

We use Ridge Regression for 6-month situation.

### 3. Model 3

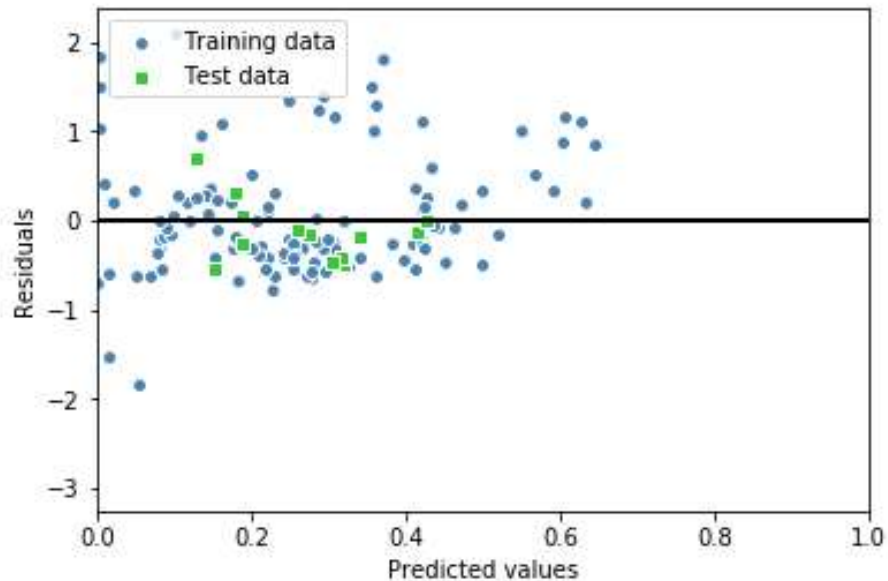
We use Lasso Regression for 9-month situation.

## HYPERPARAMETER TUNING

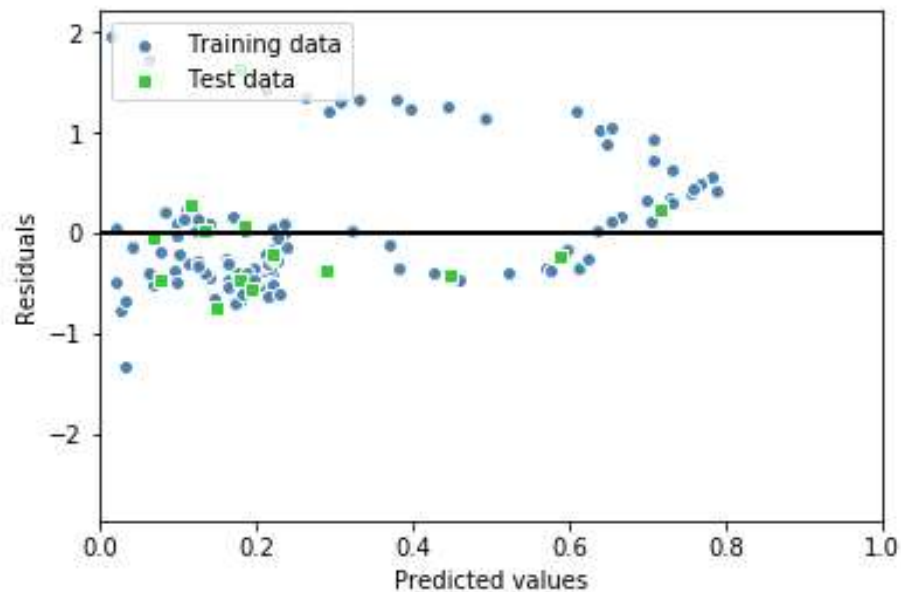
In the first case(linear regression), we cannot change the parameter, in the second and third cases, we change the alpha(ridge from  $10^{-3}$  to  $10^0$ , lasso from  $10^{-6}$  to  $10^{-3}$ ).We only show those images for the best model of each case and show the rest of them in a table.

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

Linear Regression:

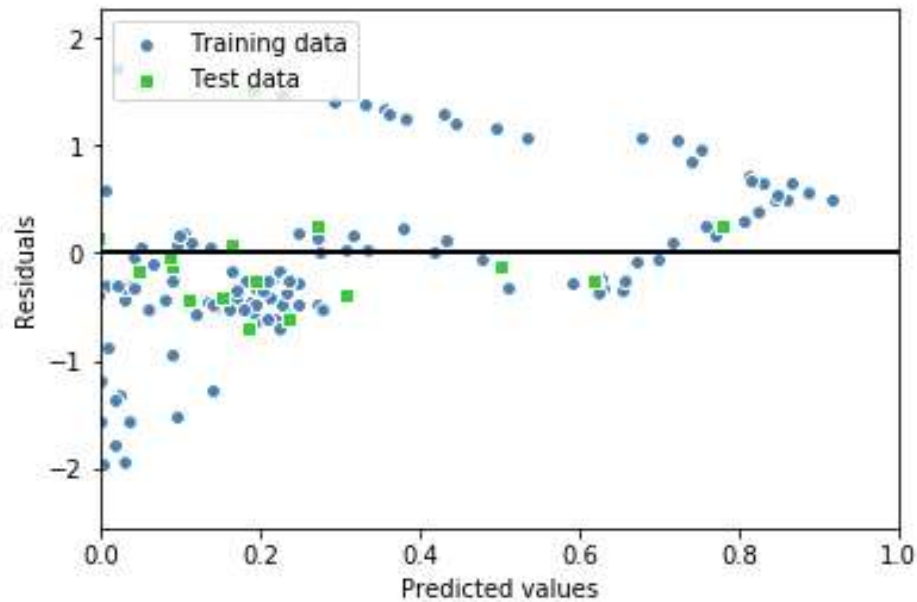


Ridge Regression:( Ridgealpha: 0.010)



错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

Lasso Regression:( Lassoalpha: 0.000100)



Here is the table:

| ridge     |           |           |           |          |           |           |
|-----------|-----------|-----------|-----------|----------|-----------|-----------|
| alpha     | MSE train | MSE test  | R^2 Train | R^2 test | Slope     | Intercept |
| 0.001     | 0.769     | 0.477     | 0.248     | 0.398    | -1.022    | -0.018    |
| 0.01      | 0.774     | 0.471     | 0.243     | 0.405    | -0.599    | -0.017    |
| 0.1       | 0.788     | 0.496     | 0.229     | 0.374    | -0.202    | -0.016    |
| 1         | 0.808     | 0.543     | 0.209     | 0.314    | -0.087    | -0.015    |
| Lasso     |           |           |           |          |           |           |
| alpha     | MSE train | MSE test  | R^2 Train | R^2 test | Slope     | Intercept |
| 0.000001  | 0.715     | 0.416     | 0.296     | 0.509    | -0.825    | -0.014    |
| 0.00001   | 0.715     | 0.416     | 0.296     | 0.509    | -0.816    | -0.014    |
| 0.0001    | 0.716     | 0.414     | 0.295     | 0.512    | -0.714    | -0.014    |
| 0.001     | 0.726     | 0.425     | 0.286     | 0.499    | -0.264    | -0.013    |
| Linear    |           |           |           |          |           |           |
| MSE train | MSE test  | R^2 Train | R^2 test  | Slope    | Intercept |           |
| 0.823     | 0.619     | 0.194     | 0.239     | -3.219   | -0.02     |           |

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

## ENSEMBLING

After the ensembling process, we get better result:

Test set MSE: 0.31

Test set R-Squared: 0.64

## CONCLUSIONS

INSERT TEXT HERE

错误!使用“开始”选项卡将 **HEADING 1** 应用于要在此处显示的文字。

## Appendix

GITHUB REPOSITORY

[IE598 F18 MLF GROUP PROJECT \(LINK\)](#)