# 法律声明

■ 课程详情请咨询

　◆ 微信公众号：北风教育

　◆ 官方网址：http://www.ibeifeng.com/

# 人工智能之机器学习

## 晚自习

主讲人：Gerry

上海育创网络科技有限公司

# 课程要求

- **课上课下"九字"真言**
  - ◆ 认真听，善摘录，勤思考
  - ◆ **多温故，乐实践**，再发散

- **四不原则**
  - ◆ 不懒散惰性，不迟到早退
  - ◆ 不请假旷课，不拖延作业

- **一点注意事项**
  - ◆ 违反"四不原则"，不包就业和推荐就业

# 严格是大爱

# 寄语

失败者找借口
成功者找方法

做别人不愿做的事，
做别人不敢做的事，
做别人做不到的事。

# 回归算法综合案例(二)：波士顿房屋租赁价格预测(作业)

■ 基于波士顿房屋租赁数据进行房屋租赁价格预测模型构建，分别使用Lasso回归、Ridge回两种回归算法构建模型；并分别构建2/3/4阶算法中的最优算法(参数)，并比较这两种回归算法的效果；另外使用lasso回归算法做特征选择(选择特征参数不为0的属性数据作为最终的特征属性，用这个选择出来的特征属性矩阵做Ridge回归)

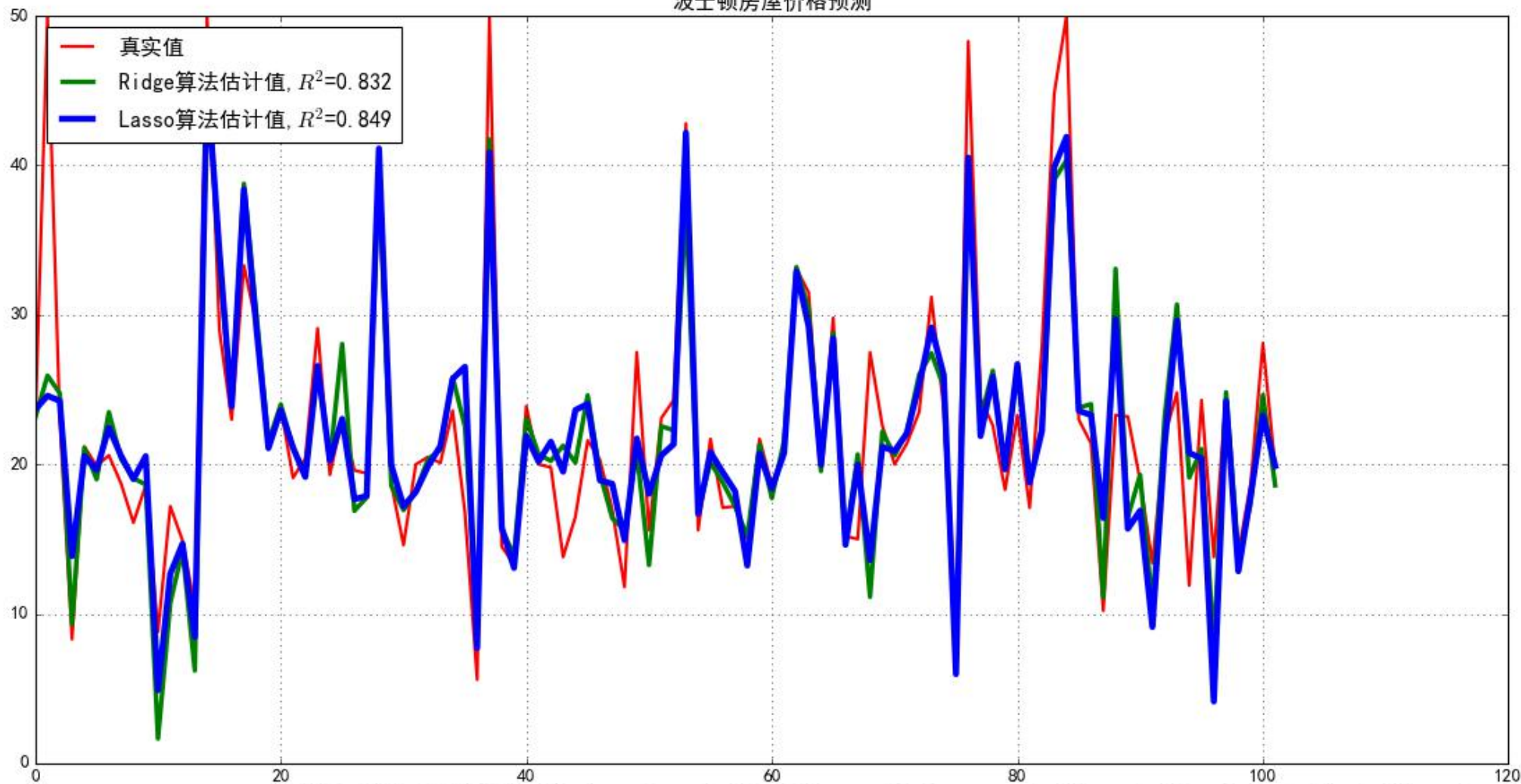◆ 数据下载url: http://archive.ics.uci.edu/ml/datasets/Housing(现在没法下载啦)

## Attribute Information:

1. CRIM: per capita crime rate by town
2. ZN: proportion of residential land zoned for lots over 25,000 sq.ft.
3. INDUS: proportion of non-retail business acres per town
4. CHAS: Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
5. NOX: nitric oxides concentration (parts per 10 million)
6. RM: average number of rooms per dwelling
7. AGE: proportion of owner-occupied units built prior to 1940
8. DIS: weighted distances to five Boston employment centres
9. RAD: index of accessibility to radial highways
10. TAX: full-value property-tax rate per $10,000
11. PTRATIO: pupil-teacher ratio by town
12. B: 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town
13. LSTAT: % lower status of the population
14. MEDV: Median value of owner-occupied homes in $1000's

| CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | MEDV |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.31533 | 0.00 | 6.200 | 0 | 0.5040 | 8.2660 | 78.30 | 2.8944 | 8 | 307.0 | 17.40 | 385.05 | 4.14 | 44.80 |
| 0.52693 | 0.00 | 6.200 | 0 | 0.5040 | 8.7250 | 83.00 | 2.8944 | 8 | 307.0 | 17.40 | 382.00 | 4.63 | 50.00 |
| 0.38214 | 0.00 | 6.200 | 0 | 0.5040 | 8.0400 | 86.50 | 3.2157 | 8 | 307.0 | 17.40 | 387.38 | 3.13 | 37.60 |
| 0.41238 | 0.00 | 6.200 | 0 | 0.5040 | 7.1630 | 79.90 | 3.2157 | 8 | 307.0 | 17.40 | 372.08 | 6.36 | 31.60 |
| 0.29819 | 0.00 | 6.200 | 0 | 0.5040 | 7.6860 | 17.00 | 3.3751 | 8 | 307.0 | 17.40 | 377.51 | 3.92 | 46.70 |
| 0.44178 | 0.00 | 6.200 | 0 | 0.5040 | 6.5520 | 21.40 | 3.3751 | 8 | 307.0 | 17.40 | 380.34 | 3.76 | 31.50 |
| 0.53700 | 0.00 | 6.200 | 0 | 0.5040 | 5.9810 | 68.10 | 3.6715 | 8 | 307.0 | 17.40 | 378.35 | 11.65 | 24.30 |
| 0.46296 | 0.00 | 6.200 | 0 | 0.5040 | 7.4120 | 76.90 | 3.6715 | 8 | 307.0 | 17.40 | 376.14 | 5.25 | 31.70 |
| 0.57529 | 0.00 | 6.200 | 0 | 0.5070 | 8.3370 | 73.30 | 3.8384 | 8 | 307.0 | 17.40 | 385.91 | 2.47 | 41.70 |
| 0.33147 | 0.00 | 6.200 | 0 | 0.5070 | 8.2470 | 70.40 | 3.6519 | 8 | 307.0 | 17.40 | 378.95 | 3.95 | 48.30 |
| 0.44791 | 0.00 | 6.200 | 1 | 0.5070 | 6.7260 | 66.50 | 3.6519 | 8 | 307.0 | 17.40 | 360.20 | 8.05 | 29.00 |
| 0.33045 | 0.00 | 6.200 | 0 | 0.5070 | 6.0860 | 61.50 | 3.6519 | 8 | 307.0 | 17.40 | 376.75 | 10.88 | 24.00 |
| 0.52058 | 0.00 | 6.200 | 1 | 0.5070 | 6.6310 | 76.50 | 4.1480 | 8 | 307.0 | 17.40 | 388.45 | 9.54 | 25.10 |
| 0.51183 | 0.000 | 6.200 | 0 | 0.5070 | 7.3580 | 71.60 | 4.1480 | 8 | 307.0 | 17.40 | 390.07 | 4.73 | 31.50 |
| 0.08244 | 30.00 | 4.930 | 0 | 0.4280 | 6.4810 | 18.50 | 6.1899 | 6 | 300.0 | 16.60 | 379.41 | 6.36 | 23.70 |
| 0.09252 | 30.00 | 4.930 | 0 | 0.4280 | 6.6060 | 42.20 | 6.1899 | 6 | 300.0 | 16.60 | 383.78 | 7.37 | 23.30 |

# 回归算法综合案例(二)：波士顿房屋租赁价格预测

```
## 模型训练 ===> 单个Lasso模型 (一阶特征选择) <2参数给定1阶情况的最优参数>
model = Pipeline([
        ('ss', StandardScaler()),
        ('poly', PolynomialFeatures(degree=1, include_bias=True, interaction_only=True)),
        ('linear', LassoCV(alphas=np.logspace(-3,1,20), fit_intercept=False))
    ])
# 模型训练
model.fit(x_train, y_train)


# 模型评测
## 数据输出
print "参数:", zip(names, model.get_params('linear')['linear'].coef_)
print "截距:", model.get_params('linear')['linear'].intercept_
```

参数: [('CRIM', 22.600592809201991), ('ZN', -0.93534557687414488), ('INDUS', 1.0202352850146854), ('CHAS', -0.0), ('NOX', 0.5948313841546149), ('RM', -1.8002644875942369), ('AGE', 2.5861907995357281), ('DIS', -0.064956108249539249), ('RAD', -2.8017533936656509), ('TAX', 1.9343329692037559), ('PTRATIO', -1.7218677875512203), ('B', -2.2762334623842988), ('LSTAT', 0.70288003005515387)]
截距: 0.0

CHAS列的数据对于LassoCV模型而言无用，所以在进行实际模型构建的时候，可以不考虑该特征

# 决策树案例二：波士顿房屋租赁价格预测(作业)

■ 使用决策树算法API对波士顿房屋租赁数据进行回归操作，预测房屋的价格信息，并理解及进行决策树API的相关参数优化

■ 数据来源：波士顿房屋租赁数据

## Housing Data Set

Download: Data Folder, Data Set Description

**Abstract**: Taken from StatLib library

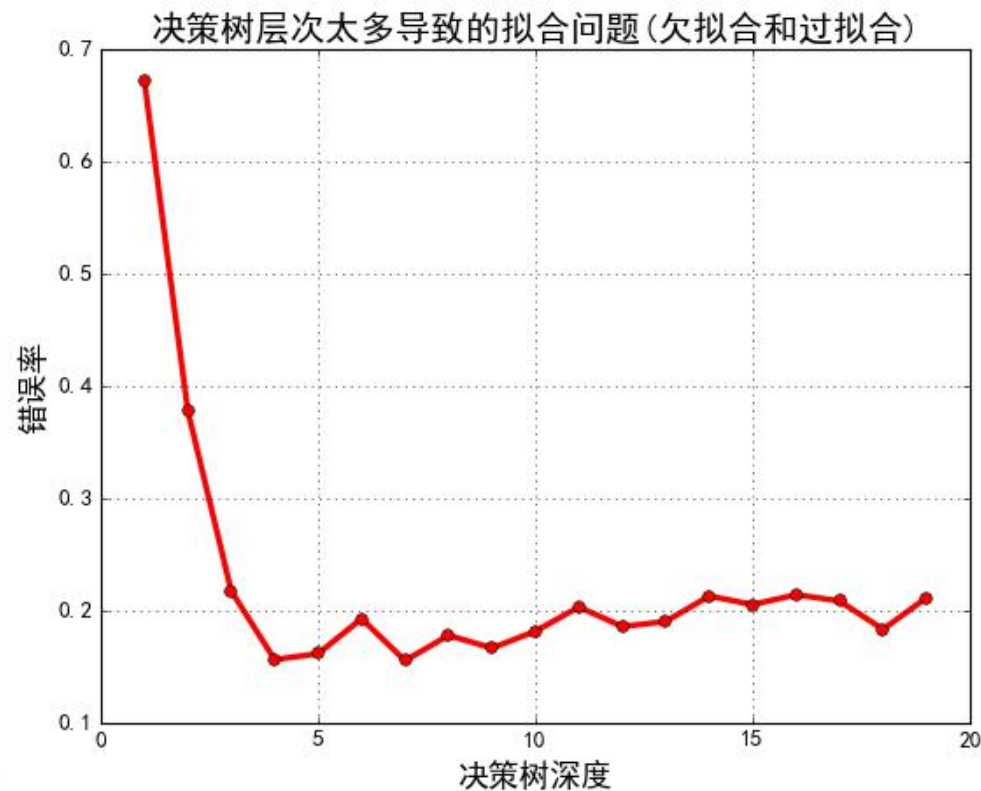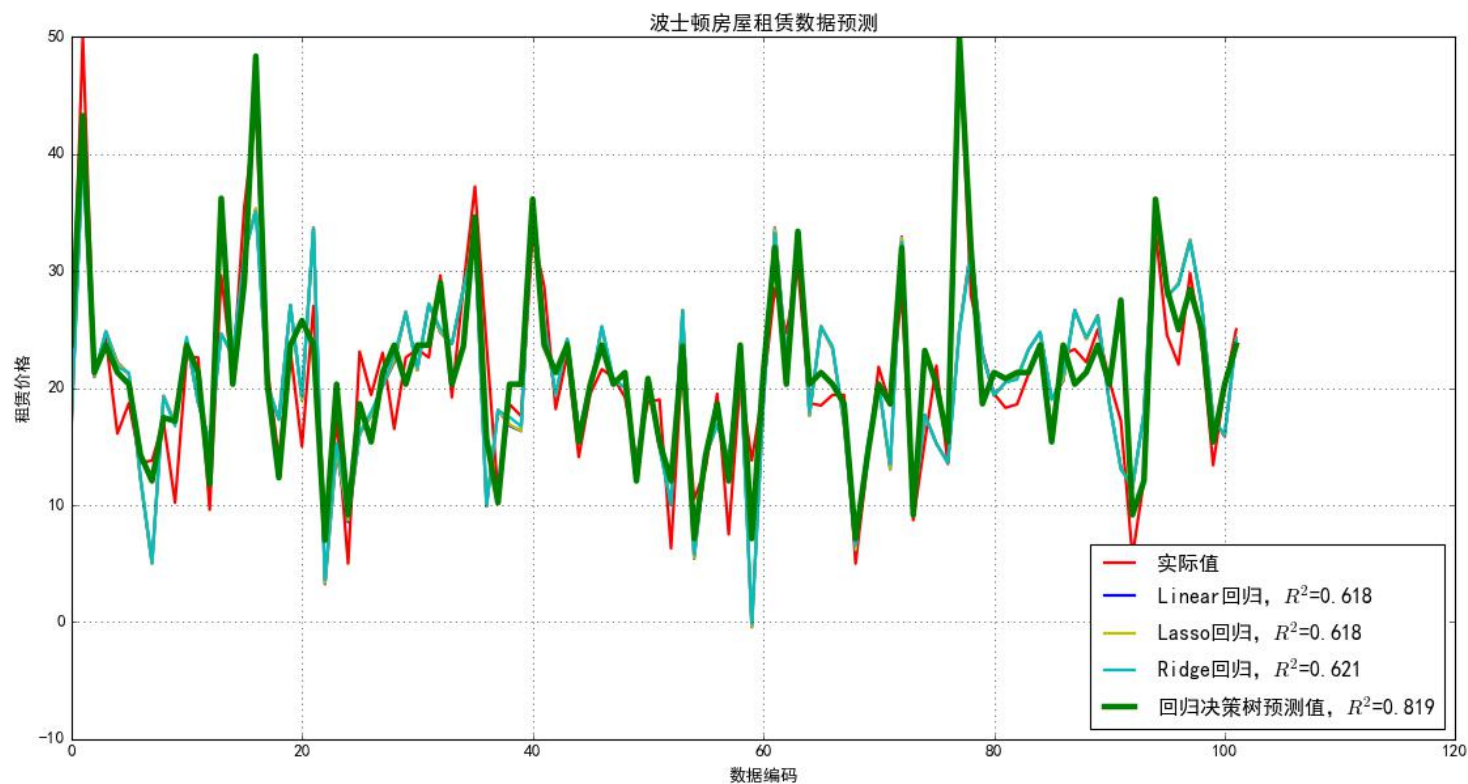| Data Set Characteristics: | Multivariate | Number of Instances: | 506 | Area: | N/A |
| --- | --- | --- | --- | --- | --- |
| Attribute Characteristics: | Categorical, Integer, Real | Number of Attributes: | 14 | Date Donated | 1993-07-07 |
| Associated Tasks: | Regression | Missing Values? | No | Number of Web Hits: | 328263 |

**Attribute Information:**

1. CRIM: per capita crime rate by town
2. ZN: proportion of residential land zoned for lots over 25,000 sq.ft.
3. INDUS: proportion of non-retail business acres per town
4. CHAS: Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
5. NOX: nitric oxides concentration (parts per 10 million)
6. RM: average number of rooms per dwelling
7. AGE: proportion of owner-occupied units built prior to 1940
8. DIS: weighted distances to five Boston employment centres
9. RAD: index of accessibility to radial highways
10. TAX: full-value property-tax rate per $10,000
11. PTRATIO: pupil-teacher ratio by town
12. B: 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town
13. LSTAT: % lower status of the population
14. MEDV: Median value of owner-occupied homes in $1000's

```
class sklearn.tree.DecisionTreeRegressor(criterion='mse', splitter='best', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None,
max_leaf_nodes=None)¶                                                                    [source]
```

9

# 决策树案例二：波士顿房屋租赁价格预测



波士顿房屋租赁数据预测

实际值
Linear回归，$R^2$=0.618
Lasso回归，$R^2$=0.618
Ridge回归，$R^2$=0.621
回归决策树预测值，$R^2$=0.819

数据编码

租赁价格



决策树层次太多导致的拟合问题（欠拟合和过拟合）

错误率

决策树深度

# GBDT回归案例：波士顿房屋租赁价格预测(作业)

- 基于波士顿房屋租赁数据进行房屋租赁价格预测模型构建，使用集成学习的算法方式对模型进行构建，比较基于GBDT的模型效果和单模型(单个线性回归、单个决策树)情况下的R2的评估值的比较。

  - 数据下载url: http://archive.ics.uci.edu/ml/datasets/Housing(现在没法下载啦)

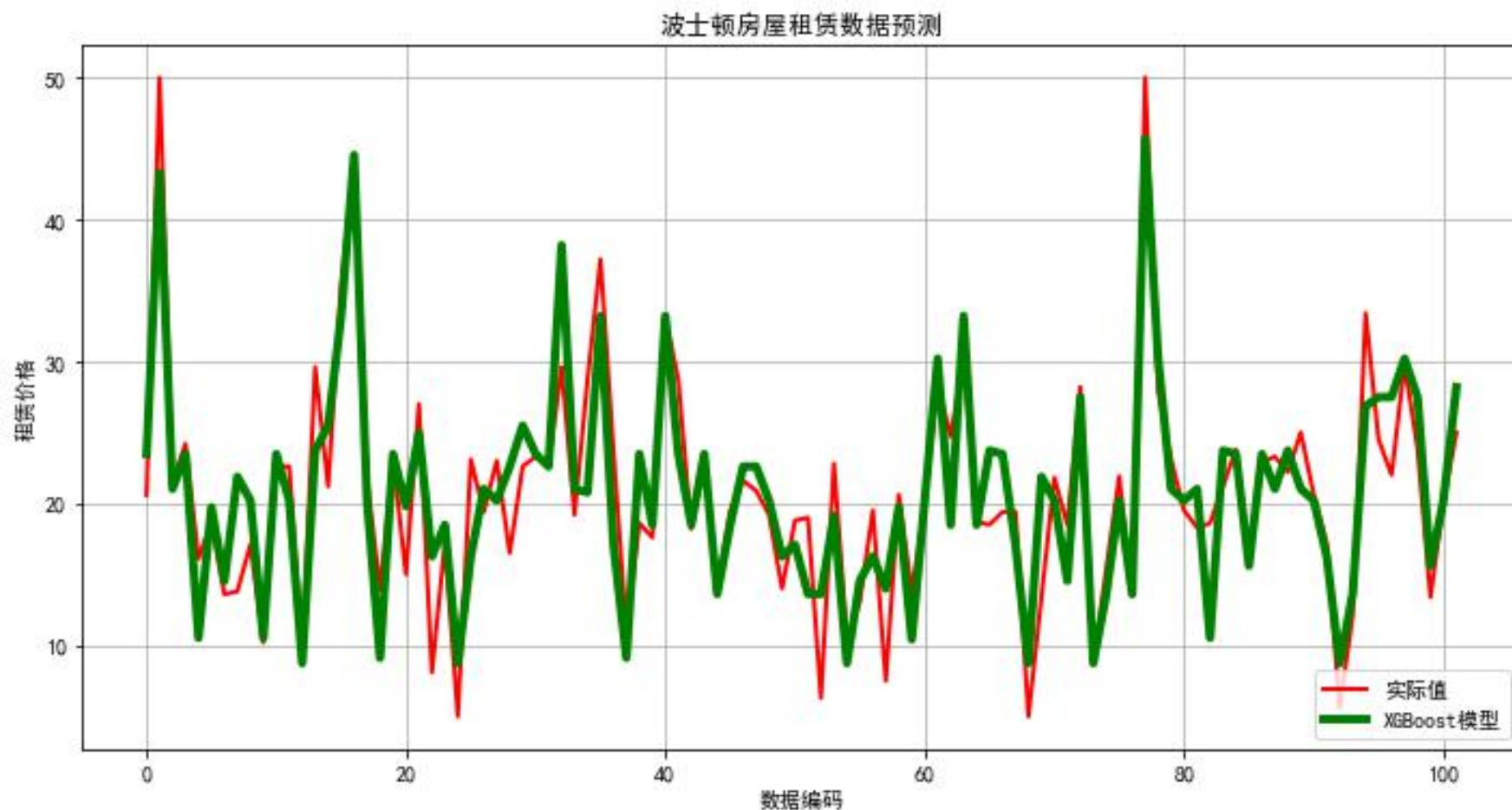**Attribute Information:** boston_housing.data

1. CRIM: per capita crime rate by town
2. ZN: proportion of residential land zoned for lots over 25,000 sq.ft.
3. INDUS: proportion of non-retail business acres per town
4. CHAS: Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
5. NOX: nitric oxides concentration (parts per 10 million)
6. RM: average number of rooms per dwelling
7. AGE: proportion of owner-occupied units built prior to 1940
8. DIS: weighted distances to five Boston employment centres
9. RAD: index of accessibility to radial highways
10. TAX: full-value property-tax rate per $10,000
11. PTRATIO: pupil-teacher ratio by town
12. B: 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town
13. LSTAT: % lower status of the population
14. MEDV: Median value of owner-occupied homes in $1000's

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.31533 | 0.00 | 6.200 | 0 | 0.5040 | 8.2660 | 78.30 | 2.8944 | 8 | 307.0 | 17.40 | 385.05 | 4.14 | 44.80 |
| 0.52693 | 0.00 | 6.200 | 0 | 0.5040 | 8.7250 | 83.00 | 2.8944 | 8 | 307.0 | 17.40 | 382.00 | 4.63 | 50.00 |
| 0.38214 | 0.00 | 6.200 | 0 | 0.5040 | 8.0400 | 86.50 | 3.2157 | 8 | 307.0 | 17.40 | 387.38 | 3.13 | 37.60 |
| 0.41238 | 0.00 | 6.200 | 0 | 0.5040 | 7.1630 | 79.90 | 3.2157 | 8 | 307.0 | 17.40 | 372.08 | 6.36 | 31.60 |
| 0.29819 | 0.00 | 6.200 | 0 | 0.5040 | 7.6860 | 17.00 | 3.3751 | 8 | 307.0 | 17.40 | 377.51 | 3.92 | 46.70 |
| 0.44178 | 0.00 | 6.200 | 0 | 0.5040 | 6.5520 | 21.40 | 3.3751 | 8 | 307.0 | 17.40 | 380.34 | 3.76 | 31.50 |
| 0.53700 | 0.00 | 6.200 | 0 | 0.5040 | 5.9810 | 68.10 | 3.6715 | 8 | 307.0 | 17.40 | 378.35 | 11.65 | 24.30 |
| 0.46296 | 0.00 | 6.200 | 0 | 0.5040 | 7.4120 | 76.90 | 3.6715 | 8 | 307.0 | 17.40 | 376.14 | 5.25 | 31.70 |
| 0.57529 | 0.00 | 6.200 | 0 | 0.5070 | 8.3370 | 73.30 | 3.8384 | 8 | 307.0 | 17.40 | 385.91 | 2.47 | 41.70 |
| 0.33147 | 0.00 | 6.200 | 0 | 0.5070 | 8.2470 | 70.40 | 3.6519 | 8 | 307.0 | 17.40 | 378.95 | 3.95 | 48.30 |
| 0.44791 | 0.00 | 6.200 | 1 | 0.5070 | 6.7260 | 66.50 | 3.6519 | 8 | 307.0 | 17.40 | 360.20 | 8.05 | 29.00 |
| 0.33045 | 0.00 | 6.200 | 0 | 0.5070 | 6.0860 | 61.50 | 3.6519 | 8 | 307.0 | 17.40 | 376.75 | 10.88 | 24.00 |
| 0.52058 | 0.00 | 6.200 | 1 | 0.5070 | 6.6310 | 76.50 | 4.1480 | 8 | 307.0 | 17.40 | 388.45 | 9.54 | 25.10 |
| 0.51183 | 0.00 | 6.200 | 0 | 0.5070 | 7.3580 | 71.60 | 4.1480 | 8 | 307.0 | 17.40 | 390.07 | 4.73 | 31.50 |
| 0.08244 | 30.00 | 4.930 | 0 | 0.4280 | 6.4810 | 18.50 | 6.1899 | 6 | 300.0 | 16.60 | 379.41 | 6.36 | 23.70 |
| 0.09252 | 30.00 | 4.930 | 0 | 0.4280 | 6.6060 | 42.20 | 6.1899 | 6 | 300.0 | 16.60 | 383.78 | 7.37 | 23.30 |

# XGBoost案例(作业)

- 使用XGBoost相关算法API对波士顿房价进行预测，并最终输出R^2值；比较一下和GBDT的执行速度。



波士顿房屋租赁数据预测

# THANK YOU

上海育创网络科技有限公司