

法律声明

□ 本课件包括演示文稿、示例、代码、题库、视频和声音等内容，小象学院和主讲老师拥有完全知识产权的权利；只限于善意学习者在本课程使用，不得在课程范围外向任何第三方散播。任何其他人或机构不得盗版、复制、仿造其中的创意及内容，我们保留一切通过法律手段追究违反者的权利。

□ 课程详情请咨询

■ 微信公众号：小象

■ 新浪微博：ChinaHadoop



决策树和随机森林实践



小象学院
ChinaHadoop.cn

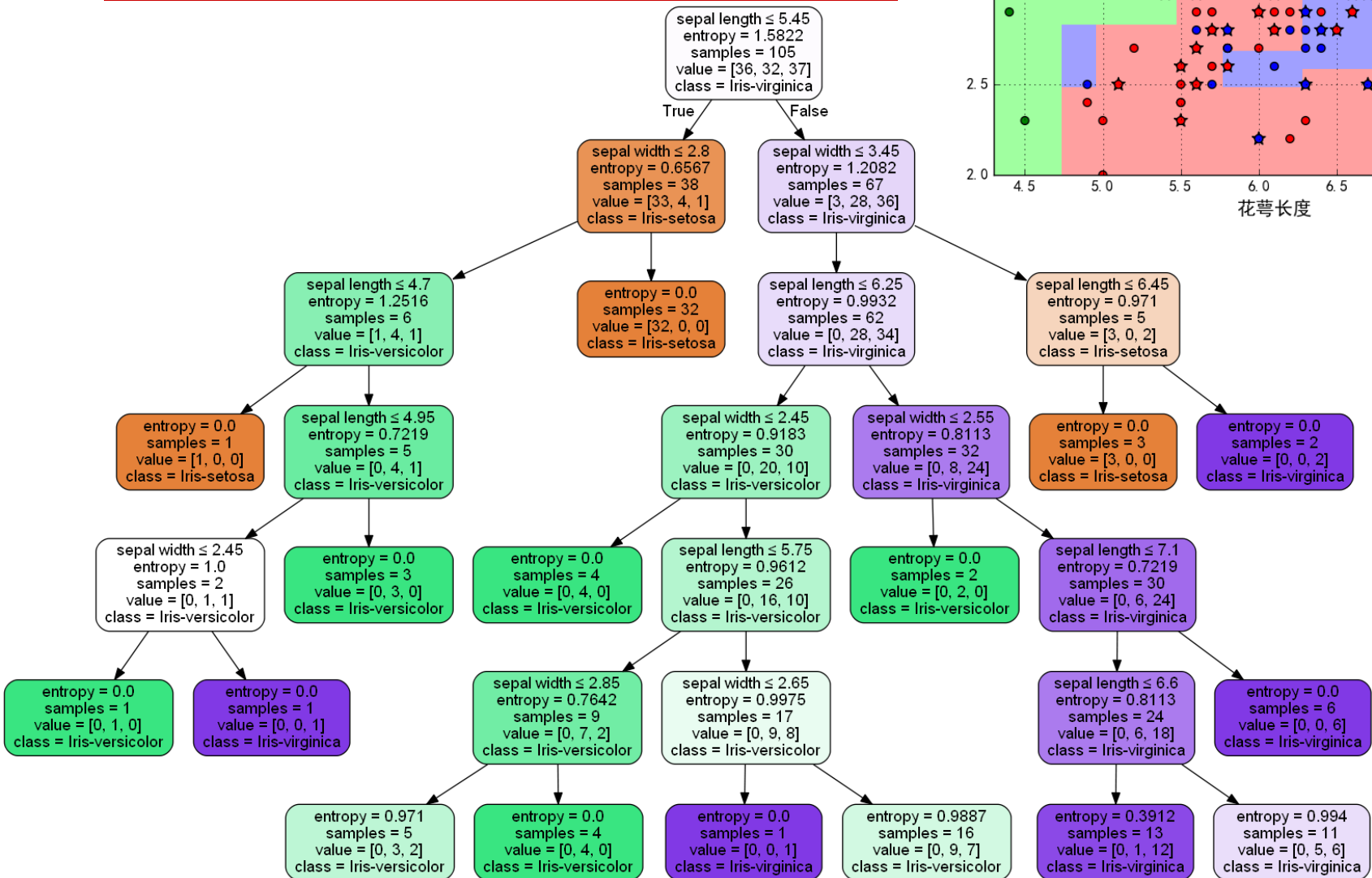
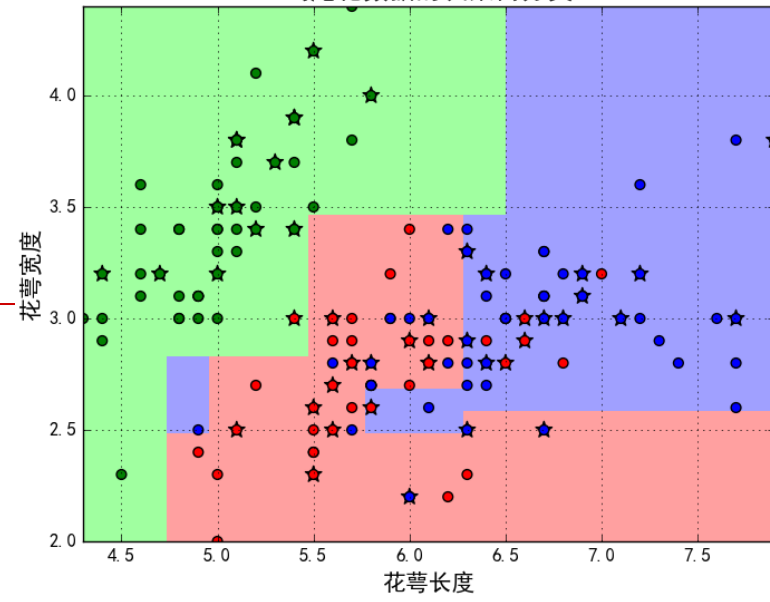
邹博

主要内容

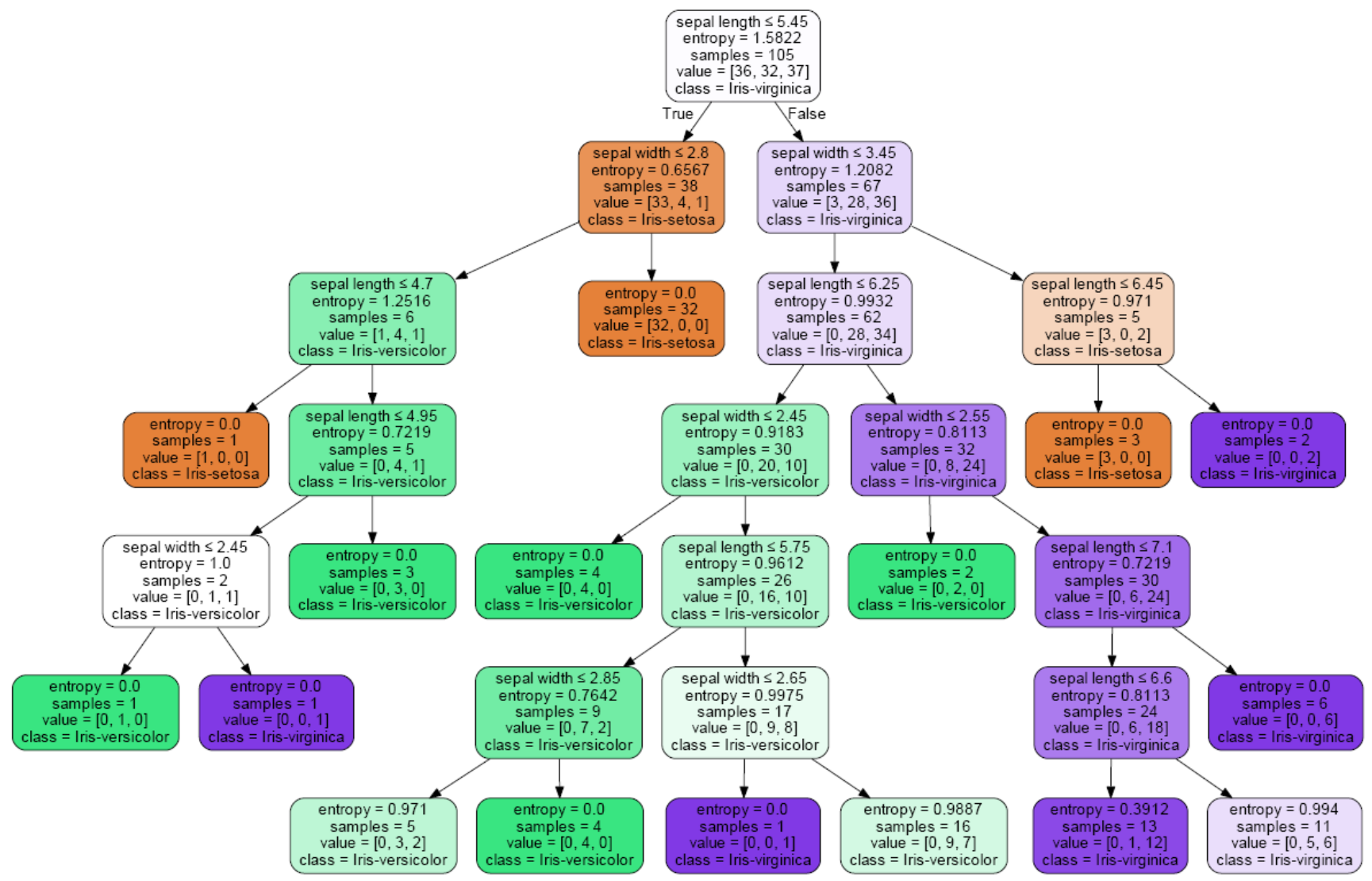
- 决策树分类鸢尾花数据
- 随机森林分类鸢尾花数据
- 决策树的可视化
- 决策树回归
- 多输出决策树

- 思考：树的深度与过拟合

鸢尾花数据决策树



书籍
缩略图
页面
附件
注释



Graphviz

GVEdit File Edit Window Graph Help

View

```

graph TD
    Node0["X[0] <= 5.55  
entropy = 1.585  
samples = 150  
value = [50, 50, 50]"]
    Node1["X[1] <= 2.8  
entropy = 0.8128  
samples = 59  
value = [47, 11, 1]"]
    Node2["X[1] <= 3.7  
entropy = 1.1671  
samples = 91  
value = [3, 39, 49]"]
    Node3["X[0] <= 4.95  
entropy = 0.8167  
samples = 12  
value = [1, 10, 1]"]
    Node4["X[0] <= 5.35  
entropy = 0.1485  
samples = 47  
value = [46, 1, 0]"]
    Node5["X[0] <= 6.25  
entropy = 0.9937  
samples = 86  
value = [0, 39, 47]"]
    Node6["entropy = 0.971  
samples = 5  
value = [3, 0, 2]"]
    Node7["entropy = 1.585  
samples = 3  
value = [1, 1, 1]"]
    Node8["entropy = 0.0  
samples = 9  
value = [0, 9, 0]"]
    Node9["entropy = 0.0  
samples = 39  
value = [39, 0, 0]"]
    Node10["X[1] <= 3.45  
entropy = 0.5436  
samples = 8  
value = [7, 1, 0]"]
    Node11["X[0] <= 5.75  
entropy = 0.909  
samples = 37  
value = [0, 25, 12]"]
    Node12["entropy = 0.9183  
samples = 3  
value = [2, 1, 0]"]
    Node13["entropy = 0.0  
samples = 5  
value = [5, 0, 0]"]
    Node14["X[1] <= 2.85  
entropy = 0.65  
samples = 12  
value = [0, 10, 2]"]
    Node15["X[1] <= 2.95  
entropy = 0.971  
samples = 25  
value = [0, 15, 10]"]
    Node16["entropy = 0.0  
samples = 3  
value = [0, 2, 1]"]
    Node17["entropy = 0.8113  
samples = 4  
value = [0, 3, 1]"]
    Node18["entropy = 0.0  
samples = 5  
value = [0, 5, 0]"]
    Node19["X[1] <= 2.95  
entropy = 0.971  
samples = 25  
value = [0, 15, 10]"]
    Node20["X[1] <= 2.65  
entropy = 0.8631  
samples = 7  
value = [0, 5, 2]"]
    Node21["entropy = 0.0  
samples = 5  
value = [0, 5, 0]"]
    Node22["X[1] <= 2.85  
entropy = 0.9367  
samples = 17  
value = [0, 11, 6]"]
    Node23["X[1] <= 3.1  
entropy = 1.0  
samples = 8  
value = [0, 4, 4]"]
    Node24["X[0] <= 5.9"]

    Node0 -- True --> Node1
    Node0 -- False --> Node2
    Node1 --> Node3
    Node1 --> Node4
    Node2 --> Node5
    Node2 --> Node6
    Node3 --> Node7
    Node3 --> Node8
    Node4 --> Node9
    Node4 --> Node10
    Node5 --> Node11
    Node5 --> Node12
    Node6 --> Node13
    Node10 --> Node14
    Node10 --> Node15
    Node11 --> Node16
    Node11 --> Node17
    Node14 --> Node18
    Node14 --> Node19
    Node15 --> Node20
    Node15 --> Node21
    Node16 --> Node22
    Node16 --> Node23
    Node17 --> Node24
    Node17 --> Node25
    Node18 --> Node26
    Node18 --> Node27
    Node19 --> Node28
    Node19 --> Node29
    Node20 --> Node30
    Node20 --> Node31
    Node22 --> Node32
    Node22 --> Node33
    Node23 --> Node34
    Node23 --> Node35
    Node24 --> Node36
    Node24 --> Node37

```

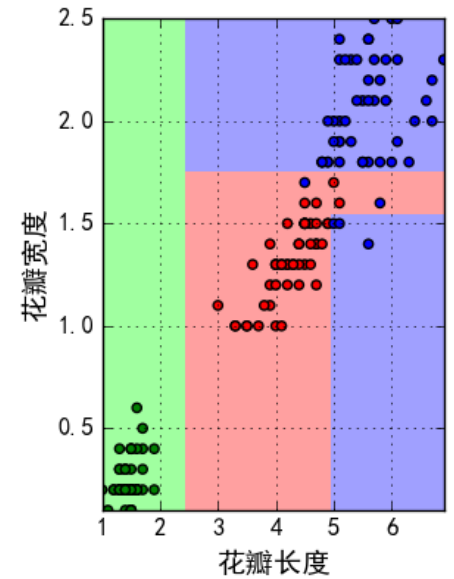
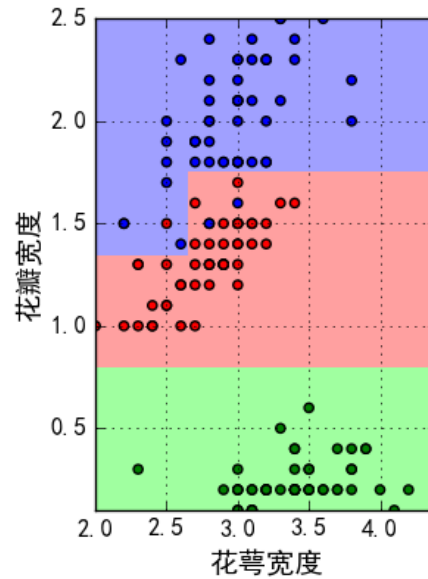
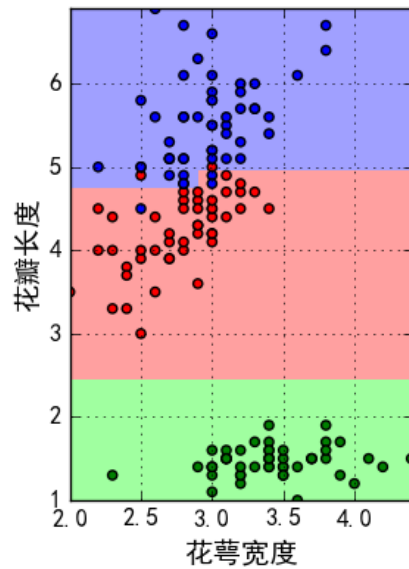
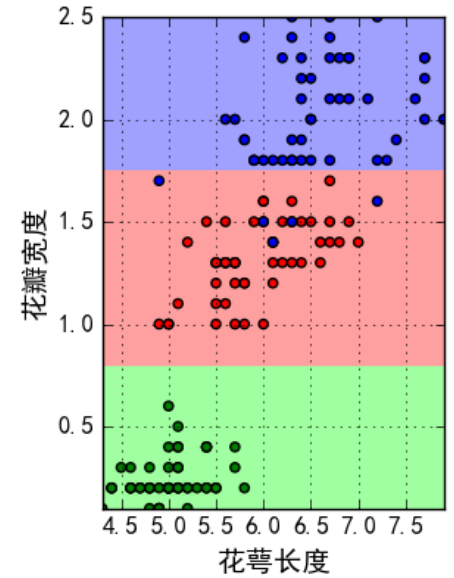
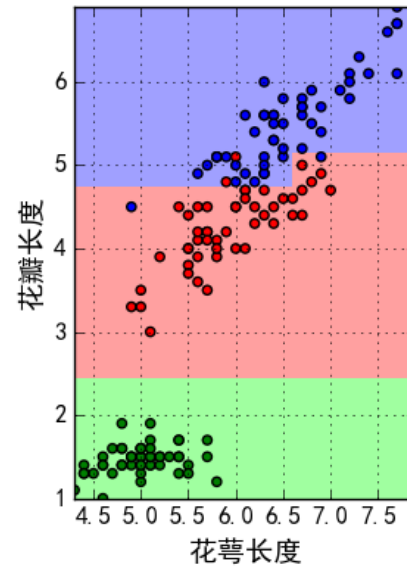
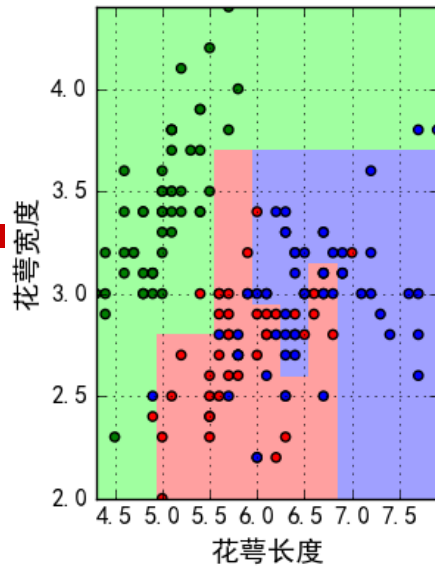
iris_tree.dot

```

digraph TD
    node [shape=box];
    0 [Label="X[0] <= 5.55\nentropy = 1.585\nsamples = 150\nvalue = [50, 50, 50]"];
    1 [Label="X[1] <= 2.8\nentropy = 0.8128\nsamples = 59\nvalue = [47, 11, 1]"];
    0 -- True --> 1;
    0 -- False --> 2;
    2 [Label="X[1] <= 3.7\nentropy = 1.1671\nsamples = 91\nvalue = [3, 39, 49]"];
    1 --> 3;
    1 --> 4;
    3 [Label="X[0] <= 4.95\nentropy = 0.8167\nsamples = 12\nvalue = [1, 10, 1]"];
    4 [Label="X[0] <= 5.35\nentropy = 0.1485\nsamples = 47\nvalue = [46, 1, 0]"];
    2 --> 5;
    2 --> 6;
    5 [Label="X[0] <= 6.25\nentropy = 0.9937\nsamples = 86\nvalue = [0, 39, 47]"];
    6 [Label="entropy = 0.971\nsamples = 5\nvalue = [3, 0, 2]"];
    3 --> 7;
    3 --> 8;
    7 [Label="entropy = 1.585\nsamples = 3\nvalue = [1, 1, 1]"];
    8 [Label="entropy = 0.0\nsamples = 9\nvalue = [0, 9, 0]"];
    4 --> 9;
    4 --> 10;
    9 [Label="entropy = 0.0\nsamples = 39\nvalue = [39, 0, 0]"];
    10 [Label="X[1] <= 3.45\nentropy = 0.5436\nsamples = 8\nvalue = [7, 1, 0]"];
    5 --> 11;
    5 --> 12;
    11 [Label="X[0] <= 5.75\nentropy = 0.909\nsamples = 37\nvalue = [0, 25, 12]"];
    12 [Label="entropy = 0.9183\nsamples = 3\nvalue = [2, 1, 0]"];
    6 --> 13;
    13 [Label="entropy = 0.0\nsamples = 5\nvalue = [5, 0, 0]"];
    10 --> 14;
    10 --> 15;
    14 [Label="X[1] <= 2.85\nentropy = 0.65\nsamples = 12\nvalue = [0, 10, 2]"];
    15 [Label="X[1] <= 2.95\nentropy = 0.971\nsamples = 25\nvalue = [0, 15, 10]"];
    11 --> 16;
    11 --> 17;
    16 [Label="entropy = 0.0\nsamples = 3\nvalue = [0, 2, 1]"];
    17 [Label="entropy = 0.8113\nsamples = 4\nvalue = [0, 3, 1]"];
    14 --> 18;
    14 --> 19;
    18 [Label="entropy = 0.971\nsamples = 25\nvalue = [0, 15, 10]"];
    19 [Label="entropy = 0.8631\nsamples = 7\nvalue = [0, 5, 2]"];
    15 --> 20;
    15 --> 21;
    20 [Label="entropy = 0.0\nsamples = 5\nvalue = [0, 5, 0]"];
    21 [Label="entropy = 0.9367\nsamples = 17\nvalue = [0, 11, 6]"];
    17 --> 22;
    17 --> 23;
    22 [Label="entropy = 1.0\nsamples = 8\nvalue = [0, 4, 4]"];
    23 [Label="entropy = 1.0\nsamples = 8\nvalue = [0, 4, 4]"];
    18 --> 24;
    18 --> 25;
    24 [Label="entropy = 0.0\nsamples = 5\nvalue = [0, 2, 3]"];
    25 [Label="entropy = 0.9183\nsamples = 3\nvalue = [0, 2, 1]"];
    19 --> 26;
    19 --> 27;
    26 [Label="entropy = 0.0\nsamples = 5\nvalue = [0, 2, 3]"];
    27 [Label="entropy = 0.9183\nsamples = 3\nvalue = [0, 2, 1]"];
    20 --> 28;
    20 --> 29;
    28 [Label="entropy = 0.9183\nsamples = 3\nvalue = [0, 2, 1]"];
    29 [Label="entropy = 0.9183\nsamples = 3\nvalue = [0, 2, 1]"];
    22 --> 30;
    22 --> 31;
    30 [Label="entropy = 0.9183\nsamples = 3\nvalue = [0, 2, 1]"];
    31 [Label="entropy = 0.9183\nsamples = 3\nvalue = [0, 2, 1]"];

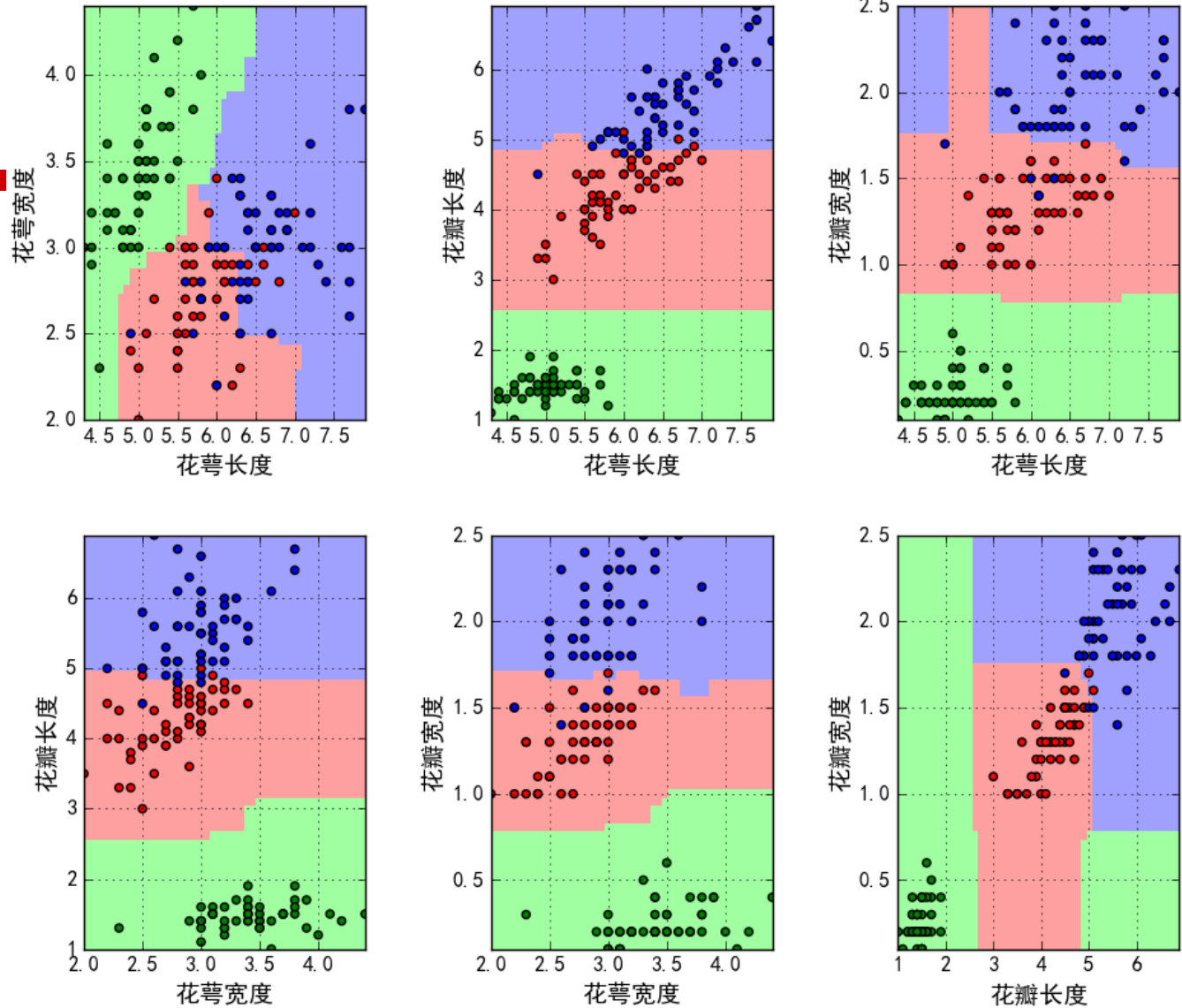
```

决策树

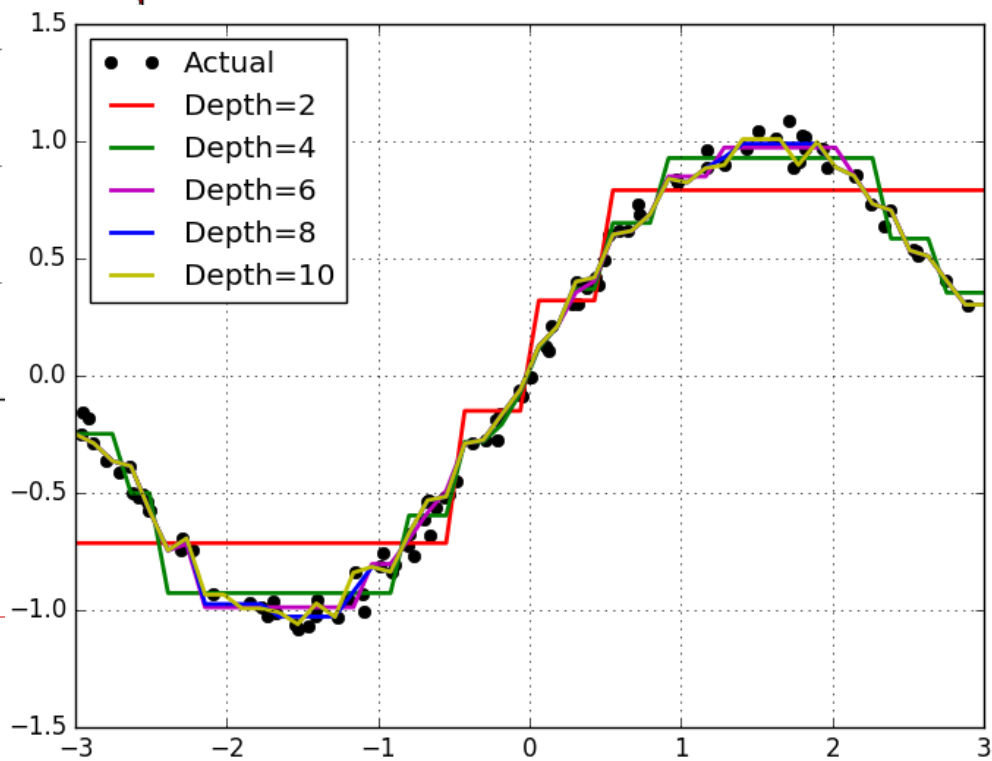
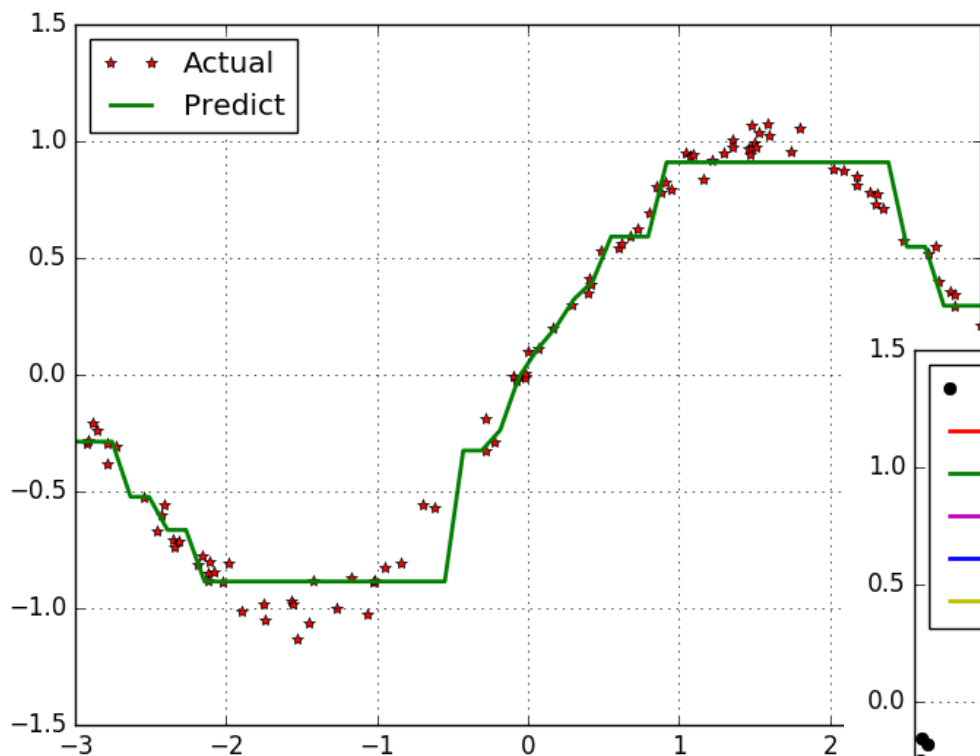


随机森林

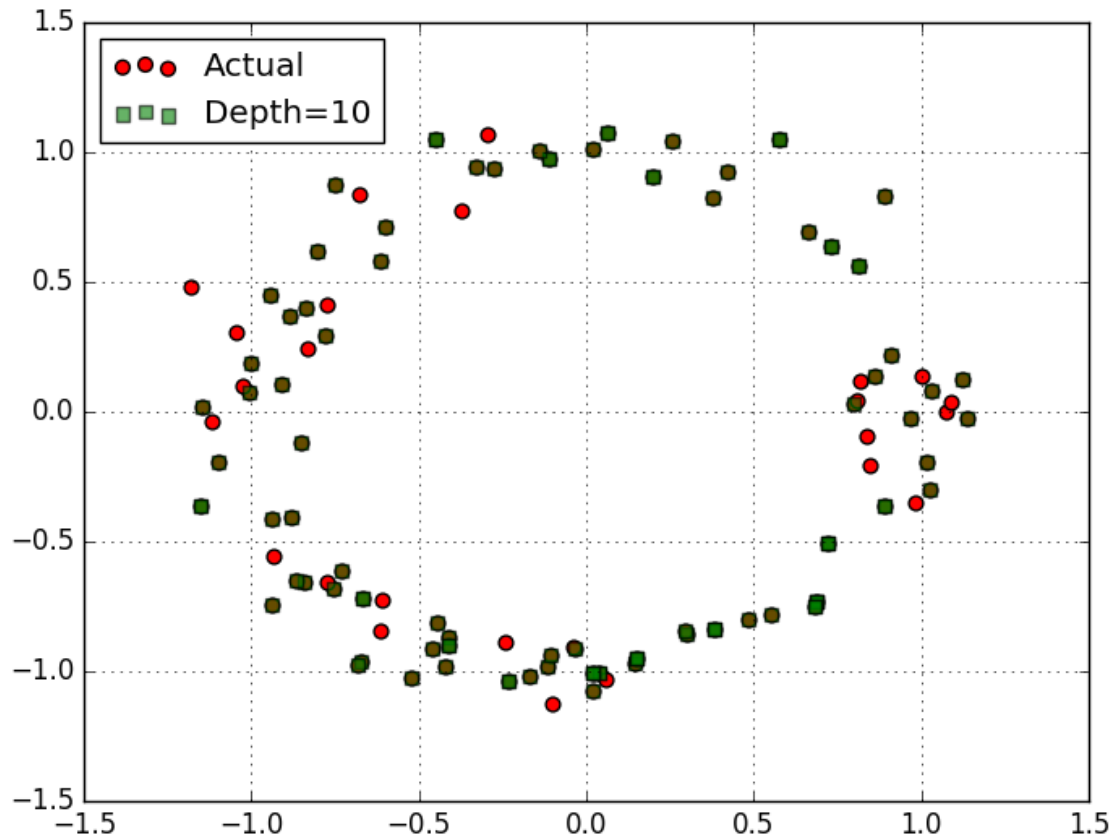
□ 50



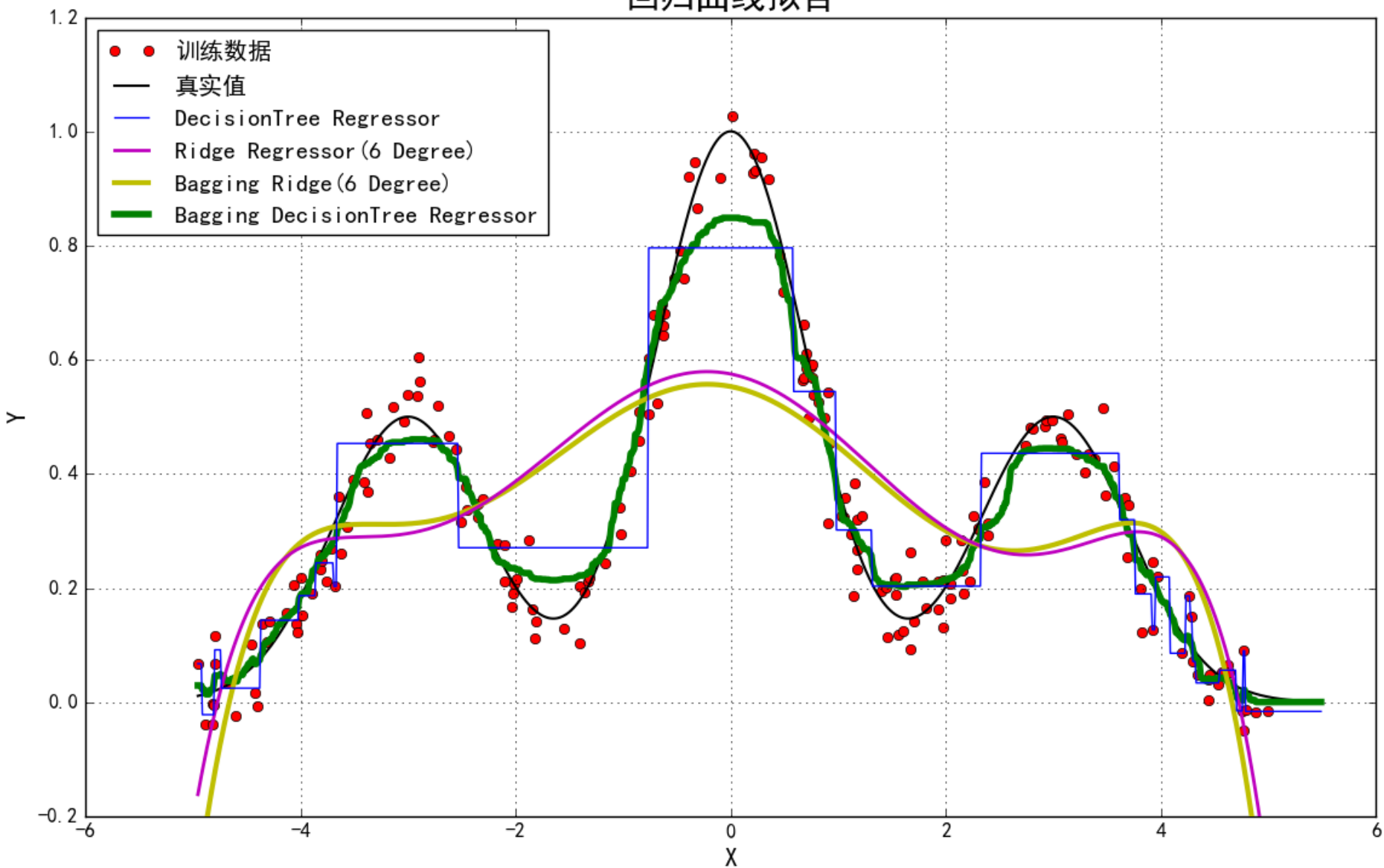
决策树用于拟合



多输出的决策树回归



回归曲线拟合



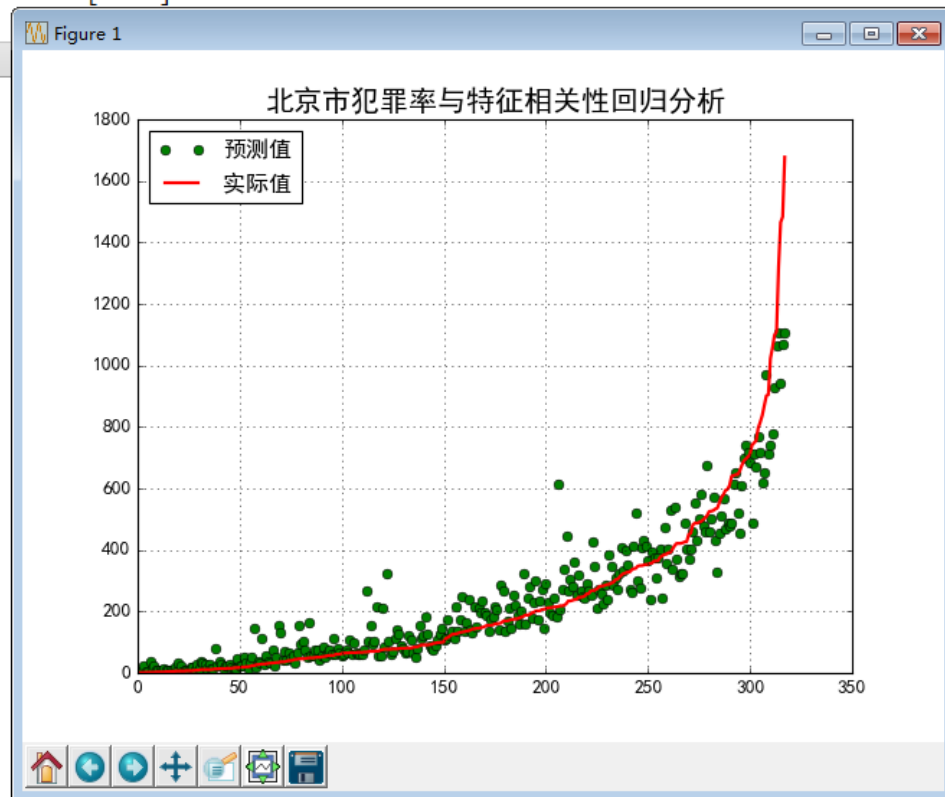
再谈北京市区域犯罪率分析

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	
	地区	盗窃案件数	批发和零售业数量	交通运输仓储邮政业数量	房地产业数量	住宿和餐饮业数量	卫生和社会工作数量	居民服务修理服务业数量	大型单位数量	中型单位数量	小微单位数量	金融业单位数量	液化石油气	能源合计	从业人员	销售费用	营业收入	营业税及附加	总产值	利润总额	人员支出	
1	安定门街道办事处																					
2	安定镇																					
3	安贞街道办事处																					
4	奥运村街道办事处																					
5	八宝山街道办事处																					
6	八达岭镇																					
7	八角街道办事处																					
8	八里庄街道办事处(朝阳)																					
9	八里庄街道办事处(海淀)																					
10	白纸坊街道办事处																					
11	百泉街道办事处																					
12	百善镇																					
13	宝山镇																					
14	北房镇																					
15	北京经济技术开发区																					
16	北七家镇																					
17	北石槽镇																					
18	北太平庄街道办事处																					
19	北务镇																					
20	北下关街道办事处																					
21	北小营镇																					
22	北新桥街道办事处																					
23																						

Code

```
model = RandomForestRegressor(n_estimators=100, criterion='mse', max_depth=10, min_samples_split=5,
                             max_features=0.6, oob_score=True)
model.fit(x, y)
print 'OOB Score = ', model.oob_score_
y_hat = model.predict(x)
rmse = np.sqrt(np.mean((y_hat - y)**2))
print 'RMSE = ', rmse, 'Predict Score = ', rmse / np.mean(y)
feature_importances = np.array(zip(columns, model.feature_importances_))
feature_importances[:, 1] = feature_importances[:, 1].astype(np.float)
feature_importances.sort(axis=0)
feature_importances = feature_importances[::-1]
for fi in feature_importances:
```

```
房地产业数量 0.0240475164437
总产值 0.0196562211564
居民服务修理服务业数量 0.0194228549579
小微单位数量 0.0187625408241
大型单位数量 0.0184318382301
地铁线路 0.0182858907058
地铁站 0.016660182329
卫生和社会工作数量 0.0144664713616
劳务费 0.0137140251708
利润总额 0.0128121358596
公用支出 0.0113859310541
公交线路 0.0113687880457
公交站 0.0112938309575
住宿和餐饮业数量 0.0107283288012
从业人员 0.00870104496645
人员支出 0.00541540683171
交通运输仓储邮政业数量 0.0053253206699
中型单位数量 0.00385674089549
```



作业

- 使用决策树做任意数据集的分类。
 - 离散变量
- 使用随机森林做数据回归。
 - 连续变量

我们在这里

□ <http://wenda.ChinaHadoop.cn>

■ 视频/课程/社区

□ 微博

■ @ChinaHadoop

■ @邹博_机器学习

□ 微信公众号

■ 小象

■ 大数据分析挖掘

The screenshot shows the website wenda.chinahadoop.cn/explore/. The page features a navigation bar with '发现' (Discover) highlighted in a red circle. Below the navigation bar, there are tabs for '全部', '招聘求职', '机器学习', '大数据平台技术', 'DCon', '大数据行业应用', 'NoSQL数据库', '数据科学', and '江湖救急'. The main content area displays a list of questions and answers, including:

- yarn运行时一直重复这个info...好像没找到资源, 应该从哪里检查呢?
- 两种不同的相关推荐列表
- 如何在Linux下配java的JDK?
- sqoop把mysql数据导入Hbase报如图错误
- 泛化误差公式推导
- kafkaOffsetMonitor打开页面以后无法显示内容?
- markdown公式编辑\$符号不起作用
- hadoop-2.6.2-src源码编译成功之后找不到native下如图一所示文件, 执行图三所示搜索命令也没有找到, 进入源码编译之后的目录如图二! 这个文件找不到怎么解决呢? 是编译没产生?
- opentsdb安装时出现72个warning, 是正常的么?
- 关于在线广告和个性化推荐区别的一点浅见

On the right side, there are sections for '专题' (Topics) including '招聘求职', '大数据行业应用', '数据科学', '系统与编程', and '云计算技术'. There are also '热门话题' (Popular Topics) like '机器学习', 'spark', '算法', 'linux', and 'hbase', and '热门用户' (Popular Users) like 'gongfc', 'Hagrid', 'yanglei', '天然下雨', and 'hiveman'.

感谢大家!

恳请大家批评指正!