

# 随机森林建模

## 简介

随机森林是由一堆决策树组成的，每一个决策树有一个结果，看有多少个决策树对同一个 $y$ 进行投票来确定 $y$ 。**分类就是少数服从多数，回归就是各个决策树取平均值**。随机森林是平均多个深决策树以降低方差的一种方法，其中，决策树是在一个数据集上的不同部分进行训练的。这是以偏差的小幅增加和一些可解释性的丧失为代价的，但是在最终的模型中通常会大大提高性能。

## 建模过程

- 1. 用有抽样放回的方法（bootstrap）从样本集中选取 $n$ 个样本作为一个训练集
- 2. 用抽样得到的样本集生成一棵决策树。在生成的每一个结点随机不重复地选择 $d$ 个特征，利用这 $d$ 个特征分别对样本集进行划分，找到最佳的划分特征（可用基尼系数、增益率或者信息增益判别）
- 3. 重复步骤1到步骤2共 $k$ 次， $k$ 即为随机森林中决策树的个数
- 4. 用训练得到的随机森林对测试样本进行预测，并用票选法决定预测的结果

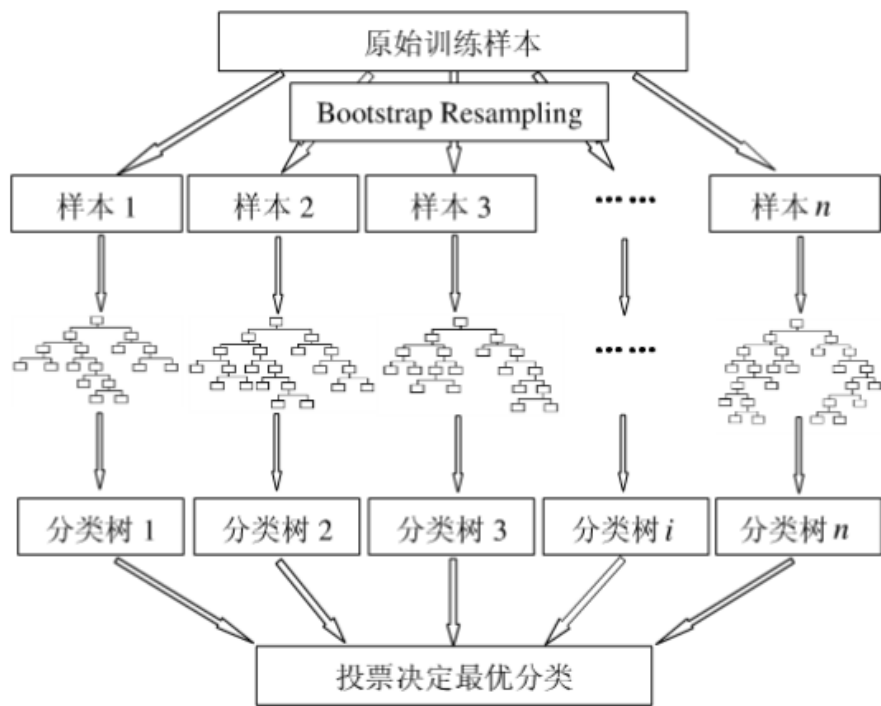


图1：随机森林算法示意图

## 实例分析

```
# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
print(accuracy_score(y_test, y_pred_rf))
print(recall_score(y_test, y_pred_rf))
print(f1_score(y_test, y_pred_rf))
print(sklearn.metrics.roc_auc_score(y_test, y_pred_rf))
```

这是银行客户流失数据的随机森林，使用了随机森林分类算法，其所有基评估器都是决策树

## n\_estimators:

这是森林中树木的数量，即基评估器的数量。这个参数对随机森林模型的精确性影响是单调的，**n\_estimators 越大，模型的效果往往越好**。但是相应的，任何模型都有决策边界，n\_estimators 达到一定的程度之后，随机森林的精确性往往不再上升或开始波动，并且，n\_estimators 越大，需要的计算量和内存也越大，训练的时间也会越来越长。对于这个参数，我们是渴望在训练难度和模型效果之间取得平衡。

## max\_depth:

树的最大深度，超过最大深度的树枝都会被剪掉

## criterion:

不纯度的衡量指标，有基尼系数和信息熵两种选择

## random\_state:

对于随机森林这个模型，它本质上是随机的，设置不同的随机状态（或者不设置 random\_state 参数）可以彻底改变构建的模型，固定 random\_state 后，每次构建的模型是相同的、生成的数据集是相同的、每次的拆分结果也是相同的

## 改变参数分析性能变化

### 1. 改变 n\_estimators，观察 accuracy\_score 变化

```
# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=8, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=9, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
print(sklearn.metrics.f1_score(y_test, y_pred_rf))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=11, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
```

随机森林		Predicted Class	0	1
Actual Class		0	659	188
		1	179	604
0	0.7748466257668711			
0	0.7713920817369093			
0	0.766984126984127			
0	0.7747161117067073			

随机森林		Predicted Class	0	1
Actual Class		0	665	182
		1	182	601
0	0.7766871165644171			
0	0.7675606641123882			
0	0.7675606641123882			
0	0.7763423155272685			

随机森林		Predicted Class	0	1
Actual Class		0	668	179
		1	184	599
0	0.7773006134969325			
0	0.7650063856960408			
0	0.7674567584881485			
0	0.7768361326355058			

随机森林		Predicted Class	0	1
Actual Class		0	673	174
		1	187	596
0	0.7785276073619631			
0	0.7611749680715197			
0	0.7675466838377334			
0	0.7778720176839299			

```
# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=12, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=13, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
```

随机森林		
Predicted Class	0	1
Actual Class		
0	673	174
1	189	594

随机森林		
Predicted Class	0	1
Actual Class		
0	675	172
1	189	592

可以发现，随着 `n_estimators` 从 8 增加到 13，`accuracy_score` 先增大，之后就基本不变了，说明 `n_estimators` 增大，一定范围内可以改善 `accuracy_score`，达到一定程度后 `accuracy_score` 就基本不变了，达到了决策边界

## 2. 改变 `max_depth`，观察 `accuracy_score` 变化

```
# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=6, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=7, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
print(accuracy_score(y_test, y_pred_rf))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=9, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=10, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=11, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
print(accuracy_score(y_test, y_pred_rf))
```

随机森林		
Predicted Class	0	1
Actual Class		
0	663	184
1	228	555

随机森林		
Predicted Class	0	1
Actual Class		
0	689	158
1	218	565

随机森林		
Predicted Class	0	1
Actual Class		
0	668	179
1	184	599

随机森林		
Predicted Class	0	1
Actual Class		
0	678	169
1	189	594

随机森林		
Predicted Class	0	1
Actual Class		
0	670	177
1	177	606

随机森林		
Predicted Class	0	1
Actual Class		
0	685	162
1	166	617

可以看出，随着 `max_depth` 从 6 增加到 11，`accuracy_score` 一直在增加，说明 `max_depth` 增大，`accuracy_score` 也在变大

## 3. 改变 `criterion`，观察 `accuracy_score` 变化

```
# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=8, criterion='entropy', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
print(accuracy_score(y_test, y_pred_rf))

# 随机森林
from sklearn.ensemble import RandomForestClassifier

classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=8, criterion='gini', random_state=42)
classifier_rf.fit(X_train, y_train)
# Predicting the Test set results
y_pred_rf = classifier_rf.predict(X_test)
print('随机森林')
print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class']))
```

随机森林		
Predicted Class	0	1
Actual Class		
0	668	179
1	184	599

随机森林		
Predicted Class	0	1
Actual Class		
0	679	168
1	190	593

<pre># 随机森林 from sklearn.ensemble import RandomForestClassifier  classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=9, criterion='entropy', random_state=42) classifier_rf.fit(X_train, y_train) # Predicting the Test set results y_pred_rf = classifier_rf.predict(X_test) print('随机森林') print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class'])) print(accuracy_score(y_test, y_pred_rf))</pre>	<pre>随机森林 Predicted Class    0    1 Actual Class 0                678  169 1                189  594 0.7803680981595092 0.7586206896551724 0.7684346701164295 0.7795464723364441 Logistics AUC: 0.759 SVM AUC: 0.695 Random Forest AUC: 0.780</pre>
<pre># 随机森林 from sklearn.ensemble import RandomForestClassifier  classifier_rf = RandomForestClassifier(n_estimators=10, max_depth=9, criterion='gini', random_state=42) classifier_rf.fit(X_train, y_train) # Predicting the Test set results y_pred_rf = classifier_rf.predict(X_test) print('随机森林') print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class'])) print(accuracy_score(y_test, y_pred_rf)) print(recall_score(y_test, y_pred_rf))</pre>	<pre>随机森林 Predicted Class    0    1 Actual Class 0                669  178 1                176  607 0.7828220858895706 0.7752234993614304 0.7742346938775511 0.7825350082403374 Logistics AUC: 0.759 SVM AUC: 0.695 Random Forest AUC: 0.783</pre>
<pre># 随机森林 from sklearn.ensemble import RandomForestClassifier  classifier_rf = RandomForestClassifier(n_estimators=11, max_depth=8, criterion='entropy', random_state=42) classifier_rf.fit(X_train, y_train) # Predicting the Test set results y_pred_rf = classifier_rf.predict(X_test) print('随机森林') print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class'])) print(accuracy_score(y_test, y_pred_rf))</pre>	<pre>随机森林 Predicted Class    0    1 Actual Class 0                673  174 1                187  596 0.7785276073619631 0.7611749680715197 0.7675466838377334 0.7778720176839299 Logistics AUC: 0.759</pre>
<pre># 随机森林 from sklearn.ensemble import RandomForestClassifier  classifier_rf = RandomForestClassifier(n_estimators=11, max_depth=8, criterion='gini', random_state=42) classifier_rf.fit(X_train, y_train) # Predicting the Test set results y_pred_rf = classifier_rf.predict(X_test) print('随机森林') print(pd.crosstab(y_test, y_pred_rf, rownames=['Actual Class'], colnames=['Predicted Class'])) print(accuracy_score(y_test, y_pred_rf)) print(recall_score(y_test, y_pred_rf))</pre>	<pre>随机森林 Predicted Class    0    1 Actual Class 0                682  165 1                187  596 0.7840490797546013 0.7611749680715197 0.772020725388601 0.7831848866331625 Logistics AUC: 0.759 SVM AUC: 0.695 Random Forest AUC: 0.783</pre>

通过改变 `criterion`，分别改变 `n_estimators`、`max_depth`，发现在相同条件下，`gini` 比 `entropy` 的 `accuracy_score` 更大