

# 第十一章 随机森林

## 11.1 简介

## 11.2 随机森林基本概况

## 11.3 随机森林基本理论

### 11.3.1 回归树基本理论

### 10.3.2 分类树基本理论

## 11.4 随机森林实践

## 11.1 简介

# 11.1 简介

- 回归树和分类树的优点是解释性强且易于计算, 然而就像一个硬币具有两面, 基于单棵树的回归或者分类模型也具有明显的缺陷:
  - ▶ 1. 由于模型的简单性及仅用单一的常数作为最终区域的预测值, 从而使得单棵树的回归或分类模型很难具有最优的预测能力<sup>[71]</sup>;
  - ▶ 2. 单棵树的预测是不稳定的<sup>[125][70]</sup>, 数据有较小的变动就可能会导致完全不一样的分裂变量和分裂点;
  - ▶ 3. 单棵树具有容易遭受选择偏差影响的问题, 即取值多的分类自变量比取值少的分类自变量更容易被选择为分裂变量<sup>[126][127][128]</sup>.
- 为了克服单一统计模型 (单棵树) 稳定性差 (方差大) 的缺陷, 集成学习方法得到了广泛的研究和发展. 所谓集成学习, 就是指分类 (回归) 器的集成. 集成学习通过构建并结合多个弱学习器来完成学习任务, 一般的方法是先产生一组个体学习器, 再用某种策略将它们结合起来, 常见的结合策略有平均法、投票法和学习法等. 例如 Bagging 方法利用 Bootstrap 方法 (有放回抽样) 对训练集进行抽样, 得到一系列新的训练集, 对每个训练集都构建一棵树, 最后通过平均法、投票法组合所有预测器得到最终的预测模型.

# 11.1 简介

- 后来, 为了克服单棵树的缺陷及降低每次 Bootstrap 抽样之间的相关性, Breiman<sup>[129]</sup> 综合以往集成学习的优缺点提出了一种新的集成学习方法——随机森林. 接下来我们将详细介绍随机森林的基本思想、算法步骤及变量重要性评价等内容.

## 11.2 随机森林基本概况

## 11.2 随机森林基本概况

- 在引入随机森林这一算法之前,先考虑与之相关的概念——Bagging. Bagging 是 Bootstrap Aggregating 的缩写.
- 特别地,我们考虑每个预测器都是决策树模型的 Bagging 算法. 在 Bagging 算法中,我们注意到在对训练集  $N_0 = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$  进行 Bootstrap 抽样 (样本量为  $n$ ) 以获得  $M$  个新的训练集, 记为  $\{N_m, m = 1, 2, \dots, M\}$ , 鉴于 Bootstrap 抽样的性质, 可以证明新训练集  $N_b, b \in \{1, 2, \dots, M\}$  大约只包含原训练集  $N_0$  的三分之二 (因为一个样本, 经  $n$  次有放回抽样, 仍未被抽中的概率是  $\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = e^{-1} \approx \frac{1}{3}$ ). 这些未被使用的观测值称为此树的袋外观测值 (out-of-bag, OOB). 袋外观测值构成的集合称为袋外示例.
- 可以将袋外示例作为对应训练集生成的树的测试集来评估训练的结果, 即可以用所有将第  $i$  个观测值作为 OOB 的树来预测第  $i$  个观测值的响应值. 我们可以对这些预测响应值求平均 (回归情况下) 或执行多数投票 (分类情况下), 已得到第  $i$  个观测值的一个 OOB 预测. 用这种方法可以求出每个观测值的 OOB 预测, 根据这些就可以计算总体的 OOB 均方误差 (对回归问题) 或分类误差 (对分类问题). 由此得到的 OOB 误差是对 Bagging 模型测试误差的有效估计. 我们将应用 OOB 误差来评估随机森林变量的重要性.

## 11.2 随机森林基本概况

- 随机森林的**基本思想**是：为了降低单棵树的缺陷及各次抽样间的相关性,随机森林采用有放回抽样,每次抽取  $n$  个样本,再无放回随机抽取  $p$  个属性(自变量)中的  $k$ (一般取  $k$  为  $\sqrt{p}$ ) 个属性,把这  $k$  个属性当成新的特征,并结合因变量和抽取的  $n$  个样本,生成一棵回归树或者分类树,重复这一过程  $M$ 次得到  $M$  棵回归树或分类树,随机森林是通过集成上述  $M$  棵回归树或分类树而成. 通过集成  $M$  棵树可以有效避免单棵树的不稳定性,而每一棵树只用  $k$  个属性代替  $p$  个属性来建模,不仅可以有效降低树之间的相关性,还能提高计算速度和节省计算机内存.
- 随机森林算法步骤如下:
  - ▶ (1) 从数据集  $(X_1, Y_1), \dots, (X_n, Y_n)$  中进行 Bootstrap 抽样 (有放回抽样), 抽取  $n$  个样本, 得到样本集  $N_m$ ;
  - ▶ (2) 利用  $N_m$  建立一棵决策树, 对于树上的每个节点, 重复以下步骤, 直到节点的样本数达到指定的最小限定值  $n_{\min}$  :
    - ▶ a) 从全部  $p$  个随机变量中随机取  $k(k < p)$  个;
    - ▶ b) 从这  $k$  个变量中选取最优分裂变量, 将此节点分裂成两个子节点.
  - ▶ 注: 对于分类问题, 构造每棵树时默认使用  $k = \sqrt{p}$  个随机变量, 节点最小样本数为 1; 对于回归问题, 构造每棵树时默认使用  $k = \frac{p}{3}$  个随机变量, 节点最小样本数为 5.

## 11.2 随机森林基本概况

- ▶ (3) 重复以上过程  $M$  次, 得到  $M$  棵树构成一个随机森林.
- ▶ (4) 当对新样本进行预测时, 由每个决策树得到一个预测结果, 再进行平均或“投票”得出最后的结果. a) 对于回归问题, 最后的预测结果为所有决策树预测值的平均数; b) 对于分类问题, 最终的预测结果为所有决策树预测结果中最多的那类, 即采用“投票”得出最后的分类结果.

■ 从随机森林的基本思想和算法可以看出, 随机森林具有以下**特点**:

- ▶ (1) 与其他的集成学习如 Bagging 相比, 由于每次只选取  $k(k < p)$  个预测, 能够有效降低树间的相关性, 从而最大限度地减少预测方差, 提高预测的精度;
- ▶ (2) 由于每次只用到  $k$  个自变量, 因此能有效节省计算时间和计算机内存;
- ▶ (3) 与 Bagging 相比, 随机森林最大的不同就在于自变量子集的规模  $k$ . 若取  $k=p$  建立随机森林, 则等同于建立 Bagging 树. 因此, Bagging 是随机森林的特例.

■ **变量重要性评估**: 与所有集成学习方法一样, 随机森林很难得到自变量(特征) 与因变量间的一个直接的显式表达关系. 因而, 很难评估自变量的重要性. 考虑到随机森林只用到了部分自变量, Breiman<sup>[129]</sup> 建议通过如下方式来度量某个特征  $X_j$  的重要性:

## 11.2 随机森林基本概况

- ▶ (1) 根据未被抽取样本 OOB 计算随机森林中第  $i$  棵回归树的袋外误差  $e_i$ ;
- ▶ (2) 随机打乱训练集在变量  $X_j$  所在列的取值顺序, 并计算新的袋外误差  $e_i^j$ ;
- ▶ (3) 重复步骤直至计算出所有决策树的误差变化, 最后变量  $X_j$  预测误差的平均变化, 即重要性指标:  $V(X_j) = \sum_{i=1}^M (e_i^j - e_i)^2 / M$ .

■ 这里袋外误差是指我们使用针对某一棵树的袋外数据得到的预测误差的均值. 因为共有  $M$  棵树, 故而有  $M$  个袋外示例. 由上述袋外示例会生成  $M$  个袋外误差  $e_i (i = 1, 2, \dots, M)$ . 由此可知, 若特征变量  $X_j$  的变化引起重要性指标增加越大, 精度减少得越多, 则说明该变量越重要.

## 11.3 随机森林基本理论

## 11.3 随机森林基本理论

- 在上一节已经详细介绍了随机森林算法的基本思想、算法及与其他集成学习相比的优缺点等,这一节我们将详细介绍随机森林相关的基本理论. 接下来我们将分别介绍回归树和分类树的基本理论.

## 11.3.1 回归树基本理论

- 假设随机森林是通过树的预测  $h(\mathbf{X}, \boldsymbol{\theta})$  生成的, 其中,  $\boldsymbol{\theta}$  是  $q$  维随机向量,  $h(\mathbf{X}, \boldsymbol{\theta})$  是  $\mathbf{R}^{p+q} \rightarrow \mathbf{R}$  上的实值函数. 不妨假设  $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$  独立同分布地来自  $(\mathbf{X}, Y)$ , 并定义均方误差为  $E_{\mathbf{X}, Y} (Y - h(\mathbf{X}, \boldsymbol{\theta}))^2$ .
- 通过对  $M$  棵单一回归树  $h(\mathbf{X}, \boldsymbol{\theta}_i)$  取平均来生成随机森林的预测, 当  $M \rightarrow \infty$  时, 我们可以得到定理 11.1 的结论, 即随机森林预测是均方收敛的.
- **定理 11.1** 随着随机森林中树的数目趋向无穷, 如下结论几乎处处成立:

$$\lim_{M \rightarrow \infty} E_{\mathbf{X}, Y} \left( Y - \frac{1}{M} h(\mathbf{X}, \boldsymbol{\theta}_i) \right)^2 \rightarrow E_{\mathbf{X}, Y} \left( Y - E_{\boldsymbol{\theta}} h(\mathbf{X}, \boldsymbol{\theta}) \right)^2.$$

- 由定理 11.1 可知当树的数目  $M \rightarrow \infty$  时, 随机森林的预测误差趋向于总体预测误差. 接下来我们将推导出随机森林预测误差的上界, 其结果在定理 11.2 中给出.

## 11.3.1 回归树基本理论

■ 定理 11.2 假设对所有  $\Theta$ ,  $E_Y = E_X(h(X, \Theta))$ , 则如下结论成立:

$$E_{X,Y} \left[ E_{\Theta} (Y - h(X, \Theta)) \right]^2 \leq \rho E_{\Theta} \left[ E_{X,Y} (Y - h(X, \Theta))^2 \right],$$

▶ 其中,  $\rho$  是  $Y - h(X, \Theta)$  和  $Y - h(X, \Theta')$  间的加权相关系数,  $\Theta$  和  $\Theta'$  相互独立.

■ 证明

$$E_{X,Y} \left[ E_{\Theta} (Y - h(X, \Theta)) \right]^2 = E_{\Theta} E_{\Theta'} E_{X,Y} (Y - h(X, \Theta))(Y - h(X, \Theta')), \quad (11.3.1)$$

▶ 式 (11.3.1) 的右边是一个协方差并且可以写成:

$$E_{\Theta} E_{\Theta'} (\rho(\Theta, \Theta') \text{sd}(\Theta) \text{sd}(\Theta')), \quad (11.3.2)$$

▶ 其中,  $\text{sd}(\Theta) = \sqrt{E_{X,Y} (Y - h(X, \Theta))^2}$ . 加权相关系数的定义为

$$\rho = E_{\Theta} E_{\Theta'} (\rho(\Theta, \Theta') \text{sd}(\Theta) \text{sd}(\Theta')) / (E_{\Theta} \text{sd}(\Theta))^2. \quad (11.3.3)$$

## 11.3.1 回归树基本理论

► 因此,

$$\begin{aligned} E_{X,Y} \left[ E_{\Theta} (Y - h(X, \Theta)) \right]^2 &= \rho (E_{\Theta} \text{sd}(\Theta))^2 \\ &\leq \rho E_{\Theta} \left[ E_{X,Y} (Y - h(X, \Theta))^2 \right]. \end{aligned} \tag{11.3.4}$$

► 证毕.

- 由定理 11.2 可以看出随机森林预测的准确性取决于单棵树的预测能力及树之间相关性的强弱. 补偿一点, 树  $\Theta$  和  $\Theta'$  的建立是相互独立的, 但两颗树在拟合数据, 即误差表现  $Y - h(X, \Theta)$  和  $Y - h(X, \Theta')$  有相关性.

## 11.3.2 分类树基本理论

- 假设随机森林是树型分类器  $\{h(\mathbf{X}, \boldsymbol{\theta}_i), i = 1, 2, \dots, M\}$  的集合, 其中  $\mathbf{X}$  是预测向量;  $\boldsymbol{\theta}_i$  是独立同分布的随机向量, 决定了单棵分类树的生成过程; 元分类器  $h(\mathbf{X}, \boldsymbol{\theta}_i)$  是用 CART<sup>[67]</sup> 算法构建的无剪枝的分类决策树. 则当  $M$  趋向无穷时, 对分类树也是收敛的, 即有如下结论:

- **定理 11.3** 随着随机森林中树的数目趋向无穷, 如下结论成立:

$$\begin{aligned} \lim_{M \rightarrow \infty} P_{\mathbf{X}, Y} \left( \left[ \frac{1}{M} \sum_{i=1}^M I(h(\mathbf{X}, \boldsymbol{\theta}_i) = Y) - \max_{j \neq Y} \frac{1}{M} I(h(\mathbf{X}, \boldsymbol{\theta}_i) = j) \right] < 0 \right) \\ \rightarrow P_{\mathbf{X}, Y} \left( \left[ P_{\boldsymbol{\theta}}(h(\mathbf{X}, \boldsymbol{\theta}) = Y) - \max_{j \neq Y} P_{\boldsymbol{\theta}}(h(\mathbf{X}, \boldsymbol{\theta}) = j) \right] < 0 \right), \end{aligned}$$

► 其中,  $I(\cdot)$  是示性函数.

- **证明** 容易证明, 参数空间  $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_M$  上存在一个零概率集合  $C$ , 在  $C$  之外, 对于所有的  $\mathbf{X}$ , 有下式成立

$$\frac{1}{M} \sum_{i=1}^M I(h(\mathbf{X}, \boldsymbol{\theta}_i) = g) \rightarrow P_{\boldsymbol{\theta}}(h(\mathbf{X}, \boldsymbol{\theta}) = g).$$

## 11.3.2 分类树基本理论

- ▶ 在一个固定的训练集和参数空间  $\Theta$  上, 所有满足  $h(\Theta, X) = g$  的  $X$  构成的集合是一个超矩形单元. 对于所有  $h(\Theta, X)$  只有有限的  $K$  个这种超矩阵单元, 记作  $S_1, \dots, S_K$ . 若  $\{X: h(\Theta, X) = g\} = S_g$ , 此时定义  $\phi(\Theta) = g$ , 并令  $N_g$  为前  $N$  次实验中  $\phi(\Theta_i) = g$  的次数. 那么有

$$\frac{1}{M} \sum_{i=1}^M I(h(X, \Theta_i) = g) = \frac{1}{M} \sum_{g=1}^M N_g I(X \in S_g).$$

- ▶ 再由大数定理可得

$$N_g = \frac{1}{M} \sum_{i=1}^M I(\phi(\Theta_i) = g),$$

- ▶ 会收敛到  $P_{\Theta}(\phi(\Theta) = g)$ . 对于  $g$  的某个值, 所有集合的并集都不会发生收敛, 得到一个概率为零的集合  $C$ , 因此在  $C$  之外有

$$\frac{1}{M} \sum_{i=1}^M I(h(X, \Theta_i) = g) \rightarrow \sum_{g=1}^M P_{\Theta}(\phi(\Theta) = g) N_g I(X \in S_g).$$

- ▶ 上式右边即是  $P_{\Theta}(h(X, \Theta) = g)$ . 这样我们就证明了定理.

## 11.3.2 分类树基本理论

■ 由定理 11.3可知, 随着随机森林中树数量增加, 模型的分类误差上限趋于一个固定值. 即随机森林不会随着分类树数目的增加而产生过度拟合的问题, 将对未知实例预测提供较好的参考思路和应用性. 类似定理 11.2, 我们将给出随机森林分类误差的一个上界, 为叙述方便, 先给出如下定义:

■ 给定一组分类器  $h(\mathbf{X}, \boldsymbol{\Theta}_1), h(\mathbf{X}, \boldsymbol{\Theta}_2), \dots, h(\mathbf{X}, \boldsymbol{\Theta}_M)$ , 并使用从随机向量  $(\mathbf{X}, Y)$  的分布中随机抽取的训练集, 将边缘函数定义为

$$\text{mg}(\mathbf{X}, Y) = \frac{1}{M} \sum_{i=1}^M I(h(\mathbf{X}, \boldsymbol{\Theta}_i) = Y) - \max_{j \neq Y} \frac{1}{M} \sum_{i=1}^M I(h(\mathbf{X}, \boldsymbol{\Theta}_i) = j),$$

▶ 其中,  $I(\cdot)$  是示性函数.

■ 边缘函数衡量的是正确分类在  $(\mathbf{X}, Y)$  的平均投票数超过任何其他分类的平均投票数的程度. 差距越大, 对分类的信心就越大. 一般泛化误差由下式给出:

$$\text{PE}^* = P_{\mathbf{X}, Y}(\text{mg}(\mathbf{X}, Y) < 0).$$

## 11.3.2 分类树基本理论

- ▶ 随机森林的边际函数为

$$\text{mg}(\mathbf{X}, Y) = P_{\theta} (h(\mathbf{X}, \Theta) = Y) - \max_{j \neq Y} P_{\theta} (h(\mathbf{X}, \Theta) = j).$$

- ▶ 分类器集  $h(\mathbf{X}, \Theta)$  的强度定义为

$$s = E_{\mathbf{X}, Y} \text{mg}(\mathbf{X}, Y).$$

- ▶ 不妨设  $s \geq 0$ , 由切比雪夫不等式可得下式成立:

$$\text{PE}^* \leq \text{Var}(\text{mg}(\mathbf{X}, Y)) / s^2. \quad (11.3.5)$$

- ▶  $\text{mg}(\mathbf{X}, Y)$  方差的一个显式表达为:

$$\mathfrak{Q}(\mathbf{X}, Y) = \arg \max_{j \neq Y} P_{\theta} (h(\mathbf{X}, \Theta) = j),$$

- ▶ 故

$$\begin{aligned} \text{mg}(\mathbf{X}, Y) &= P_{\theta} (h(\mathbf{X}, \Theta) = Y) - P_{\theta} (h(\mathbf{X}, \Theta) = \mathfrak{Q}(\mathbf{X}, Y)) \\ &= E_{\theta} \left[ I(h(\mathbf{X}, \Theta) = Y) - I(h(\mathbf{X}, \Theta) = \mathfrak{Q}(\mathbf{X}, Y)) \right] \end{aligned}$$

## 11.3.2 分类树基本理论

- ▶ 原始边际函数定义为

$$\text{rmg}(\boldsymbol{\theta}, \mathbf{X}, Y) = I(h(\mathbf{X}, \boldsymbol{\theta}) = Y) - I(h(\mathbf{X}, \boldsymbol{\theta}) = \mathcal{Q}(\mathbf{X}, Y)),$$

- ▶ 因此,  $\text{mg}(\mathbf{X}, Y)$  是  $\text{rmg}(\boldsymbol{\theta}, \mathbf{X}, Y)$  关于  $\boldsymbol{\theta}$  的期望. 对于任意函数  $f$ , 当  $\boldsymbol{\theta}, \boldsymbol{\theta}'$  独立同分布时, 有下式成立:

$$\left[ E_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) \right]^2 = E_{\boldsymbol{\theta}, \boldsymbol{\theta}'} f(\boldsymbol{\theta}) f(\boldsymbol{\theta}'),$$

- ▶ 即有

$$\text{mg}(\mathbf{X}, Y) = E_{\boldsymbol{\theta}, \boldsymbol{\theta}'} \text{rmg}(\boldsymbol{\theta}, \mathbf{X}, Y) \text{rmg}(\boldsymbol{\theta}', \mathbf{X}, Y) \quad (11.3.6)$$

- ▶ 再由式 (11.3.2) 可得

$$\begin{aligned} \text{Var}(\text{mg}(\mathbf{X}, Y)) &= E_{\boldsymbol{\theta}, \boldsymbol{\theta}'} \left( \text{Cov}_{\mathbf{X}, Y}(\text{rmg}(\boldsymbol{\theta}, \mathbf{X}, Y), \text{rmg}(\boldsymbol{\theta}', \mathbf{X}, Y)) \right) \\ &= E_{\boldsymbol{\theta}, \boldsymbol{\theta}'} \left( \rho(\boldsymbol{\theta}, \boldsymbol{\theta}') \text{sd}(\boldsymbol{\theta}) \text{sd}(\boldsymbol{\theta}') \right), \end{aligned}$$

- ▶ 其中,  $\rho(\boldsymbol{\theta}, \boldsymbol{\theta}')$  是  $\text{rmg}(\boldsymbol{\theta}, \mathbf{X}, Y)$  和  $\text{rmg}(\boldsymbol{\theta}', \mathbf{X}, Y)$  间的相关系数,  $\text{sd}(\boldsymbol{\theta})$  是  $\text{rmg}(\boldsymbol{\theta}, \mathbf{X}, Y)$  的标准差.

## 11.3.2 分类树基本理论

▶ 那么,

$$\text{Var}(\text{mg}(\mathbf{X}, Y)) = \rho (E_{\theta} \text{sd}(\Theta))^2 \leq \rho E_{\theta} \text{Var}(\Theta). \quad (11.3.7)$$

▶ 其中,  $\rho$  是相关系数的平均值, 即

$$E_{\theta} \text{Var}(\Theta) \leq E_{\theta} (E_{\mathbf{X}, Y} \text{rmg}(\Theta, \mathbf{X}, Y))^2 - s^2 \leq 1 - s^2. \quad (11.3.8)$$

▶ 由式 (11.3.7) 及 (11.3.8) 可得到随机森林分类误差的上界.

■ **定理 11.4** 随机森林泛化误差的上界由下式给出

$$\text{PE}^* \leq \frac{\rho(1-s^2)}{s^2}.$$

■ 由定理 11.4 可知, 分类错误的上界受随机森林中单个分类器强弱及分类间的相关性影响.

## 11.4 随机森林实践



实践代码