

Contents:

Lecture 17

§1 Regression tree

1. 模型设定
2. 模型预测
3. 模型拟合
4. causal tree

§2 Sampling splitting

1. 模型拟合与模型预测间的 dependency
2. double sample tree

Lecture 18

§1 Bootstrapping aggregation

1. Bagging (减小 variance)
2. Infinitesimal jackknife (估计 bagging estimator 的 variance)
3. 使用 subsample

Lecture 17 1 Random forest and its statistical inference

Random forest is considered to be an important tool for nonlinear regression and classification. However, its statistical inference has long been an issue until recently. This section introduces the construction of random forest and its application in causal inference study. Especially, we discuss how to estimate the variance of a random forest.

§1 回归树 1.1 Growing a regression tree

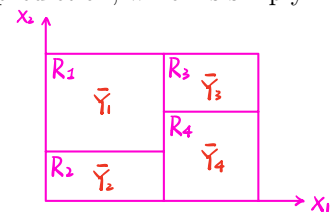
1. 模型设定 Suppose the i.i.d. observations are $(X_i, Y_i) \in \mathbf{R}^d \times \mathbf{R}$ that satisfy

$$Y_i = f(X_i) + \epsilon_i \text{ where } \epsilon_i \perp X_i, \text{ in such case } f(X_i) = \mathbf{E}Y_i|X_i,$$

2. 模型预测 we want to estimate $f(\cdot)$. The regression tree uses a set of “rectangles” to make prediction. Roughly speaking, the estimation follows two steps:

1. We divide the feature space into J distinct and non-overlapping regions R_1, \dots, R_J .
2. For every observation that falls into the region R_j , we make the same prediction, which is simply the mean of the response values for the training observations in R_j , i.e.,

$$\hat{f}(x) = \frac{1}{N_j} \sum_{X_i \in R_j} Y_i \text{ if } x \in R_j.$$



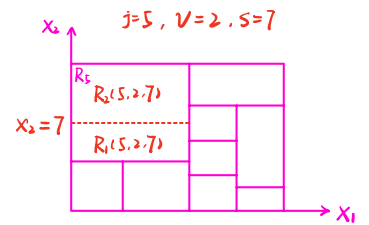
3. 模型拟合 The growing of a tree relies on the recursive binary splitting. We want to minimize the residual sum of square

$$\ell = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{f}(X_i))^2. \quad (RSS)$$

In general, this is too hard if we consider all possible R_j . However, we can do the following greedy algorithm to decrease the loss:

① Suppose we have rectangles R_1, R_2, \dots, R_{J_0} with loss

$$\ell = \sum_{j=1}^{J_0} \sum_{i \in R_j} (y_i - \hat{f}(X_i))^2.$$



② For a tuple $(j, v, s) \in \{1, 2, \dots, J_0\} \times \{1, 2, \dots, d\} \times \mathbf{R}$, we consider the half-planes

$$R_1(j, v, s) = \{X_i : X_i \in \mathbf{R}_j, X_{iv} < s\} \text{ and } R_2(j, v, s) = \{X_i : X_i \in \mathbf{R}_j, X_{iv} \geq s\}.$$

③ After that, we calculate the new loss (将 R_j 根据维度 v 和阈值 s 分割后, loss 的变化量)



$$\ell_1 = \left(\sum_{i \in R_1(j, v, s)} (y_i - \hat{f}(X_i))^2 + \sum_{i \in R_2(j, v, s)} (y_i - \hat{f}(X_i))^2 \right) - \sum_{i \in R_j} (y_i - \hat{f}(X_i))^2,$$

and find the minimizer of ℓ_1 . (调整 (j, v, s) 的选取, 使 loss 的下降量最大)

We separate the rectangle if $\ell_1 < -\alpha < 0$. (若 loss 的最大下降量小于 α , 停止迭代)

Remark If the outcome Y_i is a binary data, then the regression tree becomes decision tree. The difference here is,

1. In each leaf, we assign “yes” or “no” by checking the probability / ratio $\frac{1}{N_j} \sum_{X_i \in R_j} Y_i$.
2. For classification tree, the RSS is not a good criterion, instead, we grow the tree via the Gini index

or entropy

$$G = \sum_{j=1}^J \hat{p}_j(1 - \hat{p}_j) \text{ or } D = - \sum_{j=1}^J \hat{p}_j \log(\hat{p}_j).$$

In other words, we want the Y_i in each rectangle be in the same class.

4. causal tree

In causal inference (e.g., Wager and Athey, 2017), suppose we observe a set of data $(X_i, Y_i, W_i) \in \mathbf{R}^p \times \mathbf{R} \times \mathbf{R}$, and for each test point x , we want to estimate the conditional average treatment effect

$$\tau(x) = \mathbf{E}Y_i(1)|X_i = x - \mathbf{E}Y_i(0)|X_i = x.$$

To achieve the goal, we consider what happens in ℓ . Notice that

$$\ell = \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{f}(X_i))^2 = \sum_{j=1}^J \sum_{i \in R_j} (y_i^2 + \hat{f}(X_i)^2 - 2y_i \hat{f}(X_i)) = \sum_{j=1}^J \sum_{i \in R_j} y_i^2 - \sum_{j=1}^J \sum_{i \in R_j} 2y_i \hat{f}(X_i).$$

Therefore, minimizing ℓ is equivalent to maximizing the sum of square

$$\sum_{j=1}^J \sum_{i \in R_j} \hat{f}(X_i)^2. \quad (\text{一个重要的 observation})$$

This observation is the key to the validity of the causal tree. We consider the similar criterion

$$\ell = \sum_{j=1}^J \sum_{i \in R_j} (Y_i(1) - Y_i(0) - \hat{\tau}(X_i))^2,$$

where

$$\hat{\tau}(x) = \frac{1}{\sum_{W_i=1, X_i \in R(x)} 1} \sum_{W_i=1, X_i \in R(x)} Y_i - \frac{1}{\sum_{W_i=0, X_i \in R(x)} 1} \sum_{W_i=0, X_i \in R(x)} Y_i.$$

Unfortunately, this idea is infeasible because we cannot simultaneously observe $Y_i(0)$ and $Y_i(1)$. However, if we simulate the above procedure, this can be transformed to maximizing the sum of square



$$\sum_{j=1}^J \sum_{i \in R_j} \hat{\tau}(X_i)^2.$$

Remark There are some issues for this criterion, especially if the leaf contains too few samples.

Therefore, in practice, we may consider restricting by forcing each node containing at least k nodes.

§2 数据划分 1.2 Sampling splitting and kernel representation of random forest

1. 模型拟合与模型预测问的 dependency One of the issue in analyzing the regression tree is the dependency between the tree-growing process and the prediction. Roughly speaking, after growing a tree, the predictor for a feature x is actually

$$\hat{f}(x) = \sum_{i=1}^n K(x, X_i) Y_i,$$

(kernel representation)

($K(\cdot)$ 的求解, 即 $\{R_j\}$ 的求解, 要用上 Y)
(而 $\hat{f}(x)$ 的求解也要用上 Y)

where $K(x, X_i) = \frac{1}{\sum_{X_i \in R_j} 1}$ if x in R_j and 0 otherwise. Therefore, if the tree-growing process is independent of Y_i , then the analysis can be taken conditional on X , which motivates the following algorithm.

2. double sample tree **Algorithm 1: Double sample trees**

1. Divide the samples into two disjoint sets I and J .
2. Grow a tree via recursive partitioning. The splits are chosen using any data from I samples, and X_i from J samples, but do not use Y_i in J samples when calculating the residual sum of squares.
3. Estimate the leaf-wise nodes using data in J samples.

Lecture 18 2 Bootstrapping aggregation

§1 Bootstrap A decision tree suffers from large variance. To settle the issue, statisticians may consider leveraging the bootstrap aggregation, whose algorithm is as follows:

1. Bagging (减小 variance) **Algorithm 2: Bagging**

- The input of data $Z_i = (X_i, Y_i) \in \mathbf{R}^d \times \mathbf{R}, i = 1, 2, \dots, n$, predictor $x \in \mathbf{R}^d$.
1. Draw with replacement from Z_1, \dots, Z_n and form data Z_1^*, \dots, Z_n^* .
 2. Calculate the estimator $Y_b^* = f(x, Z_1^*, \dots, Z_n^*)$.
 3. Repeat experiment for B times and derive the estimator

$$\hat{Y} = \frac{1}{B} \sum_{b=1}^B Y_b^*.$$

The problem we are interested in is, how to analyze \hat{Y} .

First let $B \rightarrow \infty$, in such case, from law of large number, we have

$$\frac{1}{B} \sum_{b=1}^B Y_b^* \rightarrow \mathbf{E}f(x, Z_1^*, \dots, Z_n^*) | Z_1, \dots, Z_n = \frac{1}{n^n} \sum_{i_1=1}^n \sum_{i_2=1}^n \dots \sum_{i_n=1}^n f(x, Z_{i_1}, \dots, Z_{i_n}).$$

In other words, for sufficiently large B , the bagging approximately estimates the sample average across all possible combinations.

2. Infinitesimal
jackknife
(估计 bagging
estimator 的
variance)

The next problem involves estimating the variance of the bagging estimator. Efron(2014) established an estimator for the variance of $\mathbf{E}f(x, Z_1^*, \dots, Z_n^*) | Z_1, \dots, Z_n$.

Theorem 2.1 (Infinitesimal jackknife). For b th bootstrapped sample $Z_{b1}^*, Z_{b2}^*, \dots, Z_{bn}^*$, define the number $S_{bj}^* = \#\{Z_{bk}^* = Z_j\}$ for $b = 1, \dots, B$ and $j = 1, \dots, n$. An estimator for the variance of \hat{Y} is

每个 sample 在 resampling 时被 (重复) 抽出的次数

$$\hat{V} = \sum_{j=1}^n \widehat{cov}_j^2, \text{ where } \widehat{cov}_j = \frac{1}{B} \sum_{b=1}^B (S_{bj}^* - \bar{S}_{\cdot j}^*) \times (f(x, Z_{b1}^*, \dots, Z_{bn}^*) - \hat{Y}). \quad (1)$$

Remark If we assess the distribution of S_{bj}^* , we have

$$Prob(S_{bj}^* = s | Z_1, \dots, Z_n) = \binom{n}{s} \times (1/n)^s \times (1 - 1/n)^{n-s}$$

and

$$Prob(S_{b1}^* = s_1, \dots, S_{bn}^* = s_n | Z_1, \dots, Z_n) = \frac{n!}{s_1! \dots s_n!} \left(\frac{1}{n}\right)^{s_1 + \dots + s_n}.$$

In other words, S_{bj}^* in such case satisfies a binomial distribution. $S_{bj}^* \sim \text{Binomial}(n, \frac{1}{n})$

Again, suppose $B \rightarrow \infty$, then

$$\Rightarrow (S_{b1}^*, \dots, S_{bn}^*) \sim \text{multinomial}(n, \frac{1}{n}, \dots, \frac{1}{n})$$

$$\widehat{cov}_j \rightarrow Cov(S_{bj}^*, f(x, Z_{b1}^*, \dots, Z_{bn}^*) | Z_1, \dots, Z_n).$$

With such technique, suppose we want to quantify the precision of the estimator, we may consider building the confidence interval by

$$[\hat{Y} - q_{\alpha/2} \times \sqrt{\hat{V}}, \hat{Y} + q_{\alpha/2} \times \sqrt{\hat{V}}].$$

Notably, this case is for general bagging estimator. If we choose $f()$ to be a decision / regression tree, then the algorithm can be used to estimate the variance of the bagged trees.

Remark [Bias-corrected variance estimation] In practice, for B is finite, we can consider the bias-corrected version

$$\hat{V}_{debias} = \hat{V} - \frac{n}{B^2} \sum_{b=1}^B (f(x, Z_{b1}^*, \dots, Z_{bn}^*) - \hat{Y})^2$$

to decrease the bias. If we do not do that, then a $B \approx n^{1.5}$ is needed to eliminate the bias.

Remark Suppose a random vector $Q \sim \text{multinomial}(n, p_1, \dots, p_n)$, define function

$$S(p_1, \dots, p_n) = \mathbf{E}f(Q), \text{ with } (p_1^\dagger, \dots, p_n^\dagger) = (1/n, 1/n, \dots, 1/n).$$

Define the directional derivative

$$A_j = \lim_{\epsilon \rightarrow 0} \frac{S(p_1^\dagger - \frac{\epsilon}{n}, p_2^\dagger - \frac{\epsilon}{n}, \dots, p_j^\dagger + \epsilon - \frac{\epsilon}{n}, \dots, p_n^\dagger - \frac{\epsilon}{n}) - S(p_1^\dagger, \dots, p_n^\dagger)}{\epsilon},$$

then the variance is given by

$$\frac{1}{n^2} \sum_{j=1}^n A_j^2.$$

To illustrate what goes on. Suppose the estimator we are interested is

$$\hat{T} = X_1.$$

For the variance and covariances of multinomial distribution is

$$\text{Var}(X_i) = np_i(1 - p_i) \text{ and } \text{Cov}(X_i, X_j) = -np_i p_j,$$

we have

$$\text{Var}(\hat{T}) = 1 - \frac{1}{n}.$$

On the other hand, $\mathbf{E}\hat{T} = np_1$. Therefore,

$$A_1 = \lim_{\epsilon \rightarrow 0} \frac{n(p_1^\dagger + \epsilon - \frac{\epsilon}{n}) - np_1^\dagger}{\epsilon} = n - 1, \quad A_j = \lim_{\epsilon \rightarrow 0} \frac{-\epsilon}{\epsilon} = -1.$$

Therefore,

$$\frac{1}{n^2} (n-1)^2 + \frac{1}{n^2} \sum_{j=2}^n 1 = 1 - 1/n.$$

3. 使用 subsample

The idea of variance estimation in Wager and Athey (2017) for causal tree is similar. Suppose the subsample size is s and we draw without replacement, then the variance of the conditional average treatment effect can be estimated by

$$\hat{V} = \frac{n-1}{n} \left(\frac{n}{n-s} \right)^2 \sum_{i=1}^n \widehat{\text{cov}}(\hat{\tau}_b(x)^*, N_{ib}^*)^2, \text{ where } \widehat{\text{cov}}(\hat{\tau}_b(x)^*, N_{ib}^*) = \frac{1}{B} \sum_{b=1}^B (N_{ib}^* - \bar{N}_{i.}^*) \times (\hat{\tau}_b(x)^* - \hat{Y}).$$

Here $N_{ib}^* \in \{0, 1\}$ indicate whether or not the i -th training example was used for the b -th tree. Further-

more, we have

$$(\widehat{\tau}(x) - \tau(x))/\sqrt{\widehat{V}} \rightarrow_d N(0, 1),$$

so the confidence interval for the conditional average treatment effect can be

$$[\widehat{\tau}(x) - q_{\alpha/2} \times \sqrt{\widehat{V}}, \widehat{\tau}(x) + q_{\alpha/2} \times \sqrt{\widehat{V}}].$$