

Lec 13

Tuesday, October 29, 2019 10:56

Recap: trees

$$\text{tree model: } \hat{f}(x) = \sum_{l=1}^L \mathbb{I}[x \in R_l] \hat{\beta}_l$$

R_1, \dots, R_L are tree partition
of feature space

Motivation: interpretable handles hi dim (like linear model)
but also flexible (like KNN/kernel)

How to fit? first focus on regression

Consider \hat{f} 's regression error on training data
(just like we did in OLS)

$$\begin{aligned} \sum_{i=1}^n (\hat{f}(x_i) - y_i)^2 &= \sum_{i=1}^n \sum_{l=1}^L \mathbb{I}[x_i \in R_l] (\hat{\beta}_l - y_i)^2 \\ &= \sum_{l=1}^L \underbrace{\sum_{i: x_i \in R_l} (\hat{\beta}_l - y_i)^2}_{\text{min over } \hat{\beta}_l} \end{aligned}$$

$$\begin{aligned} \rightarrow \hat{\beta}_l &= \frac{1}{|R_l|} \sum_{i: x_i \in R_l} y_i \\ |R_l| &= \sum_i \mathbb{I}[x_i \in R_l] \end{aligned}$$

So if we fix R_1, \dots, R_L & opt $\hat{\beta}$

We get

$$\sum_{i=1}^n (\hat{f}(x_i) - y_i)^2 = \sum_{l=1}^L \frac{1}{|R_l|} \left(\sum_{i: x_i \in R_l} y_i - |R_l| \hat{\beta}_l \right)^2$$

$$\frac{1}{|R_l|} \left(\sum_{i: x_i \in R_l} y_i - |R_l| \hat{\beta}_l \right)^2 = \frac{1}{|R_l|} \left(\sum_{i=1}^k y_i - \frac{1}{k} \sum_{i=1}^k y_i \right)^2$$

$$= k \cdot \text{Var}(\{y_1, \dots, y_k\})$$

Target for regression tree learning:

find a tree w/ leaf regions
that have minimal sum of
impurities SSE.

Hard optimization problem.

So: instead solve greedily

Recursive Partitioning Algo

Input: $(x_1, y_1), \dots, (x_n, y_n)$, impurity fn $I(\cdot)$

Find the best $j=1, \dots, p$ & $k'=1, \dots, k-1$

$$S_{jk'}^{L/R} = \left\{ (x_i, y_i) : \begin{array}{l} i=1, \dots, k \text{ s.t.} \\ x_{ij} < / > x_{k'+1, ij} \\ \underbrace{\hspace{2cm}} \\ L: < \\ R: > \end{array} \right\}$$

(looking at sorted values of x_{ij})

to minimize $I(S_{jk'}^L) + I(S_{jk'}^R)$

Then: recurse on $S_{jk'}^L, S_{jk'}^R$
or "stop"

For Classification: $(y_i \in \{1, \dots, m\})$

The gini impurity

$$I_{\text{gini}}(\{y_1, \dots, y_k\}) = k \sum_{j=1}^m \hat{p}_j (1 - \hat{p}_j)$$

$$\text{where } \hat{p}_j = \frac{1}{k} \sum_i \mathbb{1}[y_i = j]$$

10

$$\text{For } y_i \in \mathcal{Y}_{0,1}: I_{y_i} = 2K \hat{p}_i (1 - \hat{p}_i)$$

The entropy impurity

$$I_{\text{entropy}}(\{y_1, \dots, y_K\}) = -12 \sum_{j=1}^m \hat{p}_j \log(\hat{p}_j)$$

Motivation:

$$\sum_{l=1}^L I_{\text{entropy}}(\{y_i: i \in R_l\}) = \min_{\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}} \left(\begin{array}{l} \text{negative log like} \\ \text{of data under } \hat{f} \end{array} \right)$$

When to stop partitioning?

- Can stop at a given max depth
- Can stop at a given min leaf size

How to choose these params? CV!

Other options: pruning

Ensembles

Warm up exercise for bagging

Suppose we had $\hat{f}_1, \dots, \hat{f}_B$ predictions and

that estimate $f^*(x) = \mathbb{E}[Y|X=x]$

$$\forall b=1, \dots, B \quad | \mathbb{E}[\hat{f}_b(x) - f^*(x)] | \leq \text{"bias"}$$

$$\text{Var} \hat{f}_b(x) \leq \text{"variance"}$$

$$\text{Cov}(\hat{f}_b(x), \hat{f}_{b'}(x)) = 0 \quad \forall b \neq b'$$

$$\text{Consider } \hat{f}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$$

$$\begin{aligned}
 \text{Then: } & | \mathbb{E}[\hat{f}(x) - f^*(x)] | \\
 &= | \mathbb{E} \left[\frac{1}{B} \sum_{b=1}^B \hat{f}_b(x) - f^*(x) \right] | \\
 &= | \frac{1}{B} \sum_{b=1}^B \mathbb{E}[\hat{f}_b(x) - f^*(x)] | \\
 &\leq \frac{1}{B} \sum_{b=1}^B | \mathbb{E}[\hat{f}_b(x) - f^*(x)] | \\
 &\leq \text{"bias"}
 \end{aligned}$$

$$\begin{aligned}
 \text{Var } \hat{f}(x) &= \text{Var} \left(\frac{1}{B} \sum_{b=1}^B \hat{f}_b(x) \right) \\
 &= \frac{1}{B^2} \text{Var} \left(\sum_{b=1}^B \hat{f}_b(x) \right) \\
 &= \frac{1}{B^2} \sum_{b=1}^B \text{Var}(\hat{f}_b(x)) \\
 &\leq \frac{1}{B} \text{"Variance"}
 \end{aligned}$$

When do we have $\text{Cov}(\hat{f}_b(x), \hat{f}_b(x)) = 0$?

Can happen, e.g., if \hat{f}_b is fit on separate, independent datasets.

Let's try to simulate this!

Bagging (bootstrap aggregating)

Given a sample $(x_1, y_1), \dots, (x_n, y_n)$

a bootstrap subsample is a

sample of n pts from this sample w/ replacement

For $j=1, \dots, n$:

- draw i_j at random from $\{1, \dots, n\}$

Return $(X_{i_1}, Y_{i_1}), \dots, (X_{i_n}, Y_{i_n})$

Note: - will have duplicates
- will miss some datapoints

Bagged trees:

For $b=1, \dots, B$:

- Train a CART on a new bootstrap subsample of our data to get \hat{f}_b

Return $\hat{f}(x) = \frac{1}{B} \sum_b \hat{f}_b(x)$



Usually let
depth limit = ∞
min # pts per leaf = 1

Random Forests

Idea: introduce extra randomness, in addition to bagging, to get $\text{Cov}(\hat{f}_b, \hat{f}_{b'})$ even smaller

Recall the recursive partitioning algo:

Input $X_1, Y_1, \dots, X_n, Y_n$
Find the best split $X_j \leq t$ $\forall j=1, \dots, p$ $t \in \mathbb{R}$

New randomized version:

Input: -/-
Pick p' features $J \subseteq \{1, \dots, p\}$ at random from $\{1, \dots, p\}$
Find the best split $X_j \leq t$ $\forall j \in J$ $t \in \mathbb{R}$

$L = \sum_{j=1}^n \lambda_j \sum_{i=1}^n \lambda_i \sum_{k=1}^n \lambda_k$

$RF = \left(\begin{array}{l} \text{bagging of } B \text{ trees w/} \\ \text{random subset of features} \\ \text{considered at each node} \end{array} \right)$

Usually $p' = \sqrt{p}$

Note: one tree in the ensemble may use all features at different nodes

Regression: we take avg of tree predictions

Classification: we take plurality vote over tree class predictions
or: avg of predicted probabilities