

Tree Based Methods

Greg Strabel

July 1, 2020

1 Preliminaries

Definition: A *directed graph*, G , is a tuple (V, A) , consisting of a set V , called vertices, and a set $A \subseteq \{(v, w) \mid v, w \in V, v \neq w\}$, called arrows, consisting of order pairs of distinct elements of V .

Definition: Given a directed graph $G = (V, A)$, a *directed path* from v to w is a tuple of distinct arrows $\{(v_i, w_i)\}_{i=1}^P \in A$ such that $v_1 = v$, $w_P = w$ and $\forall i < P$, $w_i = v_{i+1}$

Definition: A *binary tree* is a directed graph $G = (V, A)$ such that:

1. There is a unique $r \in V$, called *root*, such that $\forall v \in V : (v, r) \notin A$.
2. For each vertex $v \in V$, there is a unique directed path from root to v .
3. For each vertex $v \in V$, $|\{w \in V : (v, w) \in A\}| \in \{0, 2\}$.

If $(v, w) \in A$, then we say that v is a *parent* of w and w is a *child* of v . If $v \in V$ has children, we call it an *internal vertex* or *internal node*. If $v \in V$ does not have any children, we call it a *leaf*.

Given a binary tree $G_0 = (V_0, A_0)$ and a leaf $l \in V$, we can construct another binary tree that is an extension of G_0 , by adding two new vertices v and w as children of l :

$$G_1 := (V_1, A_1) := (V_0 \cup \{v, w\}, A_0 \cup \{(l, v), (l, w)\}) \quad (1)$$

Starting with a tree consisting of just a root, $G_0 = (\{r\}, \emptyset)$, we can apply the extension procedure above recursively to build arbitrarily large binary trees.

2 Basics of Recursive Binary Partitions

Assume that we have a dataset $\mathbb{D} := \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ where \mathbf{x}_i is a n -tuple, $\mathbf{x}_i = (x_{i1}, \dots, x_{in}) \in \mathbf{X} = X_1 \times \dots \times X_n$, and $\mathbf{y}_i \in \mathbf{Y}$. Given a loss function L , we desire to find a function

$f : \mathbf{X} \rightarrow \mathbf{Y}$ of the form

$$f(\mathbf{x}) := \sum_{j \in J} c_j \mathbf{1}\{\mathbf{x} \in R_j\} \quad (2)$$

where $\{R_j\}_{j=1}^J$ is a partition of \mathbf{X} that approximately minimizes $\sum_{i=1}^N L(\mathbf{y}_i, f(\mathbf{x}_i))$. This approach is generally computationally infeasible. Instead, we search for a partition of \mathbf{X} in a recursive, greedy manner.

Let

$$\mathbb{D}_R := \{(\mathbf{x}_i, \mathbf{y}_i) \in \mathbb{D} : \mathbf{x}_i \in R\}, \quad (3)$$

$$\mathfrak{D} := \{(R, \mathbb{D}_R) : R \subseteq \mathbf{X}\} \quad (4)$$

and, for $\mathbf{D} = (R, \mathbb{D}_R) \in \mathfrak{D}$,

$$2^{\mathbf{D}} := \{((R_l, \mathbb{D}_{R_l}), (R_r, \mathbb{D}_{R_r})) \in \mathfrak{D}^2 : R_r \subseteq R, R_l = R \setminus R_r\} \quad (5)$$

For $\mathbf{D} = (R, \mathbb{D}_R) \in \mathfrak{D}$ the loss improvement from the refined binary partition $(\mathbf{D}_l, \mathbf{D}_r) = ((R_l, \mathbb{D}_{R_l}), (R_r, \mathbb{D}_{R_r})) \in 2^{\mathbf{D}}$ is

$$\text{lossImp}(\mathbf{D}, \mathbf{D}_l, \mathbf{D}_r) = \min_c \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathbb{D}_R} L(\mathbf{y}_i, c) - \min_{c_l, c_r} \left[\sum_{\mathbb{D}_{R_l}} L(\mathbf{y}_i, c_l) + \sum_{\mathbb{D}_{R_r}} L(\mathbf{y}_i, c_r) \right] \quad (6)$$

We require an algorithm which, for each $\mathbf{D} \in \mathfrak{D}$, searches over a representative subset of $2^{\mathbf{D}}$ for a binary partition of \mathbf{D} that approximately maximizes lossImp . One such algorithm is the Basic Binary Split of *Algorithm* (1) on the following page. Starting with $\mathbf{D} = (\mathbf{X}, \mathbb{D})$, we apply the Basic Binary Split recursively to each leaf until some stopping criteria is reached. Common stopping criteria include a minimum number of observations in a leaf or a minimum loss improvement from splitting.

Unfortunately, the Basic Binary Split quickly becomes computationally infeasible when datasets contain categorical variables with even a small number of distinct values. For a categorical variable with K levels, there are $2^K - 1$ distinct partitions to consider when splitting. For instance, if $K = 30$, there are over 500 million possible partitions. We explore several alternative methods for splitting on categorical variables in the next section.

3 Splitting on a Categorical Variable

3.1 All possible partitions

3.2 Binary Target Variable

When \mathbf{Y} is binary, there is a shortcut for finding the best partition of a leaf $(R, \mathbb{D}_R) \in \mathfrak{D}$ on a categorical variable X_j . Without loss of generality, assume $\mathbf{Y} = \{0, 1\}$ and for each

Algorithm 1: Basic Binary Split

Input: Subset $\mathbf{D} = (R, \mathbb{D}_R) \in \mathfrak{D}$, loss function L

Output: Partition $(\mathbf{D}_l, \mathbf{D}_r) \in 2^{\mathbf{D}}$

```

1 for  $j \in \{1, \dots, n\}$  do
2   if  $X_j$  is categorical then
3      $S = \text{enum}(\{(\mathbf{D}_l, \mathbf{D}_r) \in 2^{\mathbf{D}} : (\exists A \subset X_j : R_l = A \cap R)\})$ 
4   else if  $X_j$  is ordinal then
5      $V \leftarrow \text{cutpoints}(\{\mathbf{x}_{ij}\})$  // e.g. midpoints of adjacent order
6        $S = \text{enum}(\{(\mathbf{D}_l, \mathbf{D}_r) \in 2^{\mathbf{D}} : R_l = R \cap \{x \leq d : d \in V\}\})$ 
7      $v = \vec{0} \in \mathbb{R}^{|S|}$ 
8     for  $t \in \{1, \dots, |S|\}$  do
9        $(\mathbf{D}_r, \mathbf{D}_l) = S(t)$ 
10       $v_t = \text{lossImp}(\mathbf{D}, \mathbf{D}_l, \mathbf{D}_r)$ 
11     $b_j = \text{argmax}_{t \in \{1, \dots, |S|\}} v_t$ 
12     $\mathbf{D}^{(j)} = S(b_j)$ 
13  $j^* = \text{argmax}_{j \in \{1, \dots, n\}} \{b_j\}$ 
14 return  $\mathbf{D}^{(j^*)}$ 

```

$A \subseteq X_j$ define

$$N_A = \sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R} \mathbf{1}\{\mathbf{x}_{ij} \in A\} \quad (7)$$

and

$$\bar{y}_A = \frac{1}{N_A} \sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R} y_i \mathbf{1}\{\mathbf{x}_{ij} \in A\} \quad (8)$$

Then the optimal partition (R_l, R_r) of R is of the form $R_l = R \cap A$ and $R_r = R \cap B$ where $A = \{z \in X_j : \bar{y}_{\{z\}} < p\}$ for some $p \in (0, 1)$, $B = X_j \setminus A$ and the corresponding predictions for R_l and R_r are p_A and p_B where $p_A < p_B$. To see why this is the case, suppose the optimal partition of R on X_j is (R_l, R_r) where $R_l = R \cap A$ for some $A \subset X_j$, $B = X_j \setminus A$ and $p_A < p_B$. Suppose that $\alpha \in A$ and $\beta \in B$. Since the partition is optimal

$$\bar{y}_{\{\alpha\}} L(1, p_A) + (1 - \bar{y}_{\{\alpha\}}) L(0, p_A) < \bar{y}_{\{\alpha\}} L(1, p_B) + (1 - \bar{y}_{\{\alpha\}}) L(0, p_B) \quad (9)$$

and

$$\bar{y}_{\{\beta\}} L(1, p_B) + (1 - \bar{y}_{\{\beta\}}) L(0, p_B) < \bar{y}_{\{\beta\}} L(1, p_A) + (1 - \bar{y}_{\{\beta\}}) L(0, p_A) \quad (10)$$

Rearranging and combining these equations and using the assumptions that $L(1, p)$ is decreasing in p and $L(0, p)$ is increasing in p , we have

$$\begin{aligned}\bar{y}_{\{\alpha\}} [L(1, p_A) - L(1, p_B)] &< (1 - \bar{y}_{\{\alpha\}}) [L(0, p_B) - L(0, p_A)] \\ &< (1 - \bar{y}_{\{\beta\}}) [L(0, p_B) - L(0, p_A)] \\ &< \bar{y}_{\{\beta\}} [L(1, p_A) - L(1, p_B)]\end{aligned}\tag{11}$$

which implies that $\bar{y}_{\{\alpha\}} < \bar{y}_{\{\beta\}}$. The result follows immediately from the fact that α and β were arbitrary elements of A and B , respectively.

Therefore, finding the optimal split of a categorical variable X_j in a binary classification problem reduces to a problem of order $O(|X_j|)$.

3.3 Numeric Target Variable and Squared Loss

When \mathbf{Y} is numeric and $L(y, f) = (f - y)^2$, there is a shortcut for finding the best partition of a leaf $(R, \mathbb{D}_R) \in \mathfrak{D}$ on a categorical variable X_j analogous to that for a binary response. For each $A \subseteq X_j$ define

$$N_A = \sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R} \mathbf{1}\{\mathbf{x}_{ij} \in A\}\tag{12}$$

and

$$\bar{y}_A = \frac{1}{N_A} \sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R} y_i \mathbf{1}\{\mathbf{x}_{ij} \in A\}\tag{13}$$

Then the optimal partition (R_l, R_r) of R is of the form $R_l = R \cap A$ and $R_r = R \setminus A$ where $A = \{z \in X_j : \bar{y}_{\{z\}} < p\}$ for some $p \in \mathbb{R}$ and the corresponding predictions for R_l and R_r are \bar{y}_A and $\bar{y}_{X_j \setminus A}$. To see why this is the case, suppose the optimal partition of R on X_j is (R_l, R_r) where $R_l = R \cap A$ for some $A \subset X_j$, $B = X_j \setminus A$ and $\bar{y}_A < \bar{y}_B$. Since the split is optimal and we are using squared error loss, the predictions on R_l and R_r are \bar{y}_A and \bar{y}_B , respectively. Suppose that $\alpha \in A$ and $\beta \in B$. Optimality implies that

$$\begin{aligned}\sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R : \mathbf{x}_{ij} = \alpha} \left[(y_i - \bar{y}_{\{\alpha\}})^2 + (\bar{y}_{\{\alpha\}} - \bar{y}_A)^2 \right] &= \sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R : \mathbf{x}_{ij} = \alpha} (y_i - \bar{y}_A)^2 \\ &< \sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R : \mathbf{x}_{ij} = \alpha} (y_i - \bar{y}_B)^2 \\ &= \sum_{(\mathbf{x}_i, y_i) \in \mathbb{D}_R : \mathbf{x}_{ij} = \alpha} \left[(y_i - \bar{y}_{\{\alpha\}})^2 + (\bar{y}_{\{\alpha\}} - \bar{y}_B)^2 \right]\end{aligned}\tag{14}$$

It follows that

$$(\bar{y}_{\{\alpha\}} - \bar{y}_A)^2 < (\bar{y}_{\{\alpha\}} - \bar{y}_B)^2\tag{15}$$

A similar argument shows that

$$\left(\bar{y}_{\{\beta\}} - \bar{y}_B\right)^2 < \left(\bar{y}_{\{\beta\}} - \bar{y}_A\right)^2 \quad (16)$$

It follows that

$$\bar{y}_{\{\alpha\}} < \frac{1}{2} (\bar{y}_A + \bar{y}_B) < \bar{y}_{\{\beta\}} \quad (17)$$

establishing the desired results. Hence, finding the optimal split of a categorical variable X in a regression problem with squared error loss reduces to a problem of order $O(|X|)$.

3.4 Numeric Response and Quadratic Taylor Approximation to the Loss Function

3.5 Start Right, Pull Left

Versions of this algorithm appeared in Buntine and Caruana (1991) and Mehta et al. (1996)

Algorithm 2: Start Right, Pull Left

Input: Subset $\mathbf{D} = (R, \mathbb{D}_R) \in \mathfrak{D}$, loss function L , categorical feature X_j

Output: Partition $(\mathbf{D}_l, \mathbf{D}_r) \in 2^{\mathbf{D}}$

```

1  $A = \emptyset, B = X_j$ 
2 for  $j \in \{1, \dots, |X_j| - 1\}$  do
3   for  $\alpha \in B$  do
4      $W = A \cup \{\alpha\}, U = B \setminus \{\alpha\}$ 
5      $R_W = R \cap \{x \in \mathbf{X} : x_j \in W\}, R_U = R \setminus R_W$ 
6      $(\mathbf{D}_W, \mathbf{D}_U) = ((R_W, \mathbb{D}_{R_W}), (R_U, \mathbb{D}_{R_U}))$ 
7      $b_j(\alpha) = \text{lossImp}(\mathbf{D}, \mathbf{D}_{R_W}, \mathbf{D}_{R_U})$ 
8    $\alpha_j^* = \text{argmax}_{\alpha \in B} b_j(\alpha)$ 
9    $b_j^* = b_j(\alpha_j^*)$ 
10   $A_j^* = A \cup \{\alpha_j^*\}, B_j^* = B \setminus \{\alpha_j^*\}$ 
11   $A = A_j^*, B = B_j^*$ 
12  $j^* = \text{argmax}_j \{b_j^*\}$ 
13  $R_l = R \cap \{x \in \mathbf{X} : x_j \in A_{j^*}^*\}, R_r = R \setminus R_l$ 
14 return  $(\mathbf{D}_l, \mathbf{D}_r)$ 
```

3.6 Principal Component Scoring

Principal Component Scoring is a heuristic method for splitting categorical variables when the target is a C -class categorical variable. The method was developed by Coppersmith,

Hong and Hosking (Partitioning Nominal Attributes in Decision Trees). Suppose $|X| = K$. For ease of notation, enumerate both the K levels of X and the C levels of Y . Define the $K \times C$ matrix \mathbf{N} by $\mathbf{N}_{k,c} = \sum_i \mathbf{1}\{x_i = k, y_i = c\}$. For $k \in \{1, \dots, K\}$, let $N_k = \sum_c \mathbf{N}_{k,c}$. Define the $K \times C$ matrix \mathbf{P} by $\mathbf{P}_{k,c} = \mathbf{N}_{k,c}/N_k$. Let the vector of mean class probabilities be

$$\bar{\mathbf{p}} = \frac{1}{N} \sum_k \mathbf{N}_k.$$

Let

$$\Sigma = \frac{1}{N-1} \sum_k \mathbf{N}_k (\mathbf{P}_{k\cdot} - \bar{\mathbf{p}}) (\mathbf{P}_{k\cdot} - \bar{\mathbf{p}})^T$$

Let \mathbf{v} be the first principal component of Σ . The principal component score of $\mathbf{P}_{k\cdot}$ is

$$S_k = \mathbf{v} \cdot \mathbf{P}_{k\cdot}.$$

If $\{S_{(k)}\}$ are the order statistics for $\{S_k\}$ then we consider all partitions of the form

$$A_j = \{k : S_k < S_{(j)}\}, \quad B_j = X \setminus A_j \quad (18)$$

Coppersmith et al. recommend also considering partitions of the form

$$A_j = \{k : S_k < S_{(j)} \text{ or } S_k = S_{(j+1)}\}, \quad B_j = X \setminus A_j \quad (19)$$

Experiments in Coppersmith et al. suggests that this heuristic performs comparatively well

3.7 Work

We require a Partition Refinement Algorithm (PRA):

$$C : \mathbf{D} \in \mathfrak{D} \mapsto C(\mathbf{D}) \in 2^R \quad (20)$$

which specifies a two-partition refinement of any $\mathbf{D} \in \mathfrak{D}$. An obvious PRA is the Basic Binary Split