

Decision Tree Models

①

Start w/ training data, labels

$$x_i \in \mathbb{R}^n, i=1 \dots N \quad \text{data}$$

$$y_i \in \mathbb{R}^N \quad \text{labels}$$

No index, just concat.

all labels into one big vector.

split

I.e. each ~~ele~~ element of \mathbb{R}^N

We want to partition the space of features ST. the x_i 's w/ the same y_i are grouped together.

Let's represent the data @ some node m w/ Q_m , w/ n_m samples.
 made up of a set of x_i, y labels

Then, represent each split as $\theta = (j, t_m)$
 \uparrow Feature (I.e. index of each x_i element)
 \uparrow Threshold value.

\Rightarrow Partition data into left & right subsets:

$$Q_m^{\text{left}}(\theta) = \{(x, y) \mid x_j \leq t_m\} \quad (\text{Below or equal to threshold})$$

$$Q_m^{\text{right}}(\theta) = Q_m \setminus Q_m^{\text{left}}(\theta) \quad (\text{everything else})$$

Then, define a quality score based on a loss function:

$$G(Q_m, \theta) = \frac{n_m^{\text{left}}}{n_m} H(Q_m^{\text{left}}(\theta)) + \frac{n_m^{\text{right}}}{n_m} H(Q_m^{\text{right}}(\theta))$$

Then, self select θ^* that min. G : $\theta^* = \underset{\theta}{\operatorname{argmin}} G(Q_m, \theta)$

To define H , define a measure for the proportion of class k observations in node m :

②

$$p_{mk} = \frac{1}{n_m} \sum_{y \in Q_m} \mathbb{I}(y=k)$$

\Rightarrow Gini loss: $H(Q_m) = \sum_k p_{mk} (1 - p_{mk})$

log loss: $H(Q_m) = -\sum_k p_{mk} \log(p_{mk})$

Random Forest: & XG Boost.

General idea: combine influence of several different decision tree models.

RF :

Build a forest where each tree uses a random subset of features \rightarrow low correlation b/w trees.

Sklearn implementation - use all features or random subset of features. "Random subspace method"

2) Each tree is built from a sample drawn w/ replacement from the training set.

Net effect: decrease variance of ensemble, reduce overfitting.
Combine diverse trees, at cost of increasing model bias.

If training data is $[x_1, x_2, x_3, x_4, x_5, x_6]$,
one tree gets $[x_1, x_1, x_3, x_4, x_4, x_4, \dots]$. same size
 $\uparrow \quad \quad \quad \uparrow \quad \uparrow$
 Sampled with replacement
 "Bagging" / "Bootstrapping"

Form a prediction w/ a consensus of trees.

XG Boost

Additive trees: classifier & regressor trees
& CART:

Group Members

Louis

Rob.

Sabari

Yeonjoon

Shree.

Interests

Likes GPUs

Yes

No.

Sabari
Yeonjoon
ShreeLouis
Rob.Pred. score
+2Pred score
0.5

ML.

Organic

Rugs XTB regularly

Yes

No.

Louis
Rob
YeonjoonSabari
ShreePred score
+1.5Pred
score
0.3

Organic

ML

$$f(\text{Sabari}) = 2 + 0.3 = 2.3$$

Mathematically, $\hat{y}_i = \sum_{k=1}^K f_k(x_i)$, $f_k \in \mathcal{F} \Rightarrow \mathcal{F}$ is set of all possible CARTs.
 $K \Rightarrow$ total number of trees.

$$\Rightarrow \min_{\theta} \text{obj} = \sum_{i=1}^N l(y_i, \hat{y}_i) + \sum_{k=1}^K \omega(f_k)$$

Loss fn.
over data.

Regularization
over trees

Controls
tree
complexity

③.

What are parameters of trees \Rightarrow what's "f"?

\hookrightarrow Structure of tree + leaf scores!

This isn't a gradient optimization problem:

Learning on space of all possible trees is totally intractable!

Instead, additive predictions: Fix what we have learned, add one tree at a time.

$$\hat{y}_i^0 = 0$$

$$\hat{y}_i^1 = f_1(x_i) := \hat{y}_i^0 + f_1(x_i) \quad \text{First CART tree.}$$

$$\hat{y}_i^2 = \frac{f_1(x_i) + f_2(x_i)}{2} = \hat{y}_i^{(2)} + f_2(x_i)$$

$$\hat{y}_i^{(t)} = \sum_{k=1}^t f_k(x_i) = \hat{y}_i^{(t-1)} + f_t(x_i)$$

Into objective fn:

$$\text{obj}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \omega(f_i)$$

$$= \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \omega(f_t) + c$$

We can only optimize structure of t^{th} tree!
Roll other ω 's into const.

If we apply MSE as the loss fn:

$$\text{obj}^{(t)} = \sum_{i=1}^n (y_i - (\hat{y}_i^{(t-1)} + f_t(x_i)))^2 + \omega(f_t) + c$$

$$= \sum_{i=1}^n (y_i^2 - 2y_i(\hat{y}_i^{(t-1)} + f_t(x_i)) + (\hat{y}_i^{(t-1)} + f_t(x_i))^2) + \omega(f_t) + c$$

$$= \sum_{i=1}^n [2(\hat{y}_i^{(t-1)} - y_i)f_t(x_i) + (f_t(x_i))^2] + \omega(f_t) + c$$

$$\text{obj}^{(t)} = \sum_{i=1}^n \left[2(\hat{y}_i^{(t-1)} - y_i) f_{*}(x_i) + (f_{*}(x_i))^2 \right] + \omega(f_t) + c. \quad (4)$$

\uparrow \uparrow
 1st 2nd
 order order
 term term
 (residual)

We got a nice form in the case of MSE; but for other loss fns this isn't the case! So, in general, use Taylor expansion:

$$\text{obj}^{(t)} = \sum_{i=1}^n \left[l(y_i, \hat{y}_i^{(t-1)}) + g_i f_{*}(x_i) + \frac{1}{2} h_i f_{*}^2(x_i) \right] + \omega(f_t) + c$$

where $g_i = \frac{\partial l(y_i, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}}; h_i = \frac{\partial^2 l(y_i, \hat{y}_i^{(t-1)})}{\partial^2 \hat{y}_i^{(t-1)}} \}$ "gradient boosting"

Constants don't affect minimization procedure:

$$\Rightarrow \text{obj}^{(t)} = \sum_{i=1}^n \left[g_i f_{*}(x_i) + \frac{1}{2} h_i f_{*}^2(x_i) \right] + \omega(f_t)$$

XGBoost defines $\omega(f)$ as:

$$\omega(f) = \gamma T + \frac{1}{2} \sum_{j=1}^T \omega_j^2$$

\nwarrow Total num of trees \nwarrow Vector of scores on each leaf
 \uparrow Gamma \uparrow hyperparam

Rewrite & compress: $\text{obj}^{(t)} = \sum_{j=1}^T [G_j \omega_j + \frac{1}{2} (H_j + \lambda) \omega_j^2] + \gamma T.$

SHAP Analysis

①

SHAP \Rightarrow "Shapley Additive Explanations"

Comes from cooperative game theory:

A number of players cooperate to achieve an objective, that leads to some overall gain.

But, some players may contribute more

One might have more bargaining power

Another might threaten to destroy everything

\Rightarrow Question: what amount of the overall gain should be assigned to each player?

I.e. how important is each player to the overall gain, and how much payoff should each player expect

John: Imagine the group being formed one ^{feature} person at a time. Each person demands their contribution as "fair compensation".

Then, for each person, average contribution over the different ways we can form the group.

(2)

For ML: Assume only some features are present, while others aren't

>

Shapley values are only attribution method that satisfies:

1) Efficiency: Feature contributions sum to the diff of a prediction x and the average:

$$\sum_{j=1}^p \phi_j = \hat{f}(x) - E_x \left(\hat{f}(x) \right) \quad \begin{array}{l} \text{Expectation} \\ \text{value.} \end{array}$$

2) Symmetry: Contributions of features j & k are the same if they contribute equally to all possible groups.

$$\text{If } \text{val.}(S \cup \{j\}) = \text{val.}(S \cup \{k\}).$$

$$\forall S \subseteq \{1, \dots, p\} \setminus \{j, k\}$$

↑ Group w/o j, k

$$\text{Then } \phi_j = \phi_k.$$

3) Dummy: A feature j that doesn't change predicted value, regardless of group, gets a value of 0.

$$\text{If } \text{val}(S \cup \{j\}) = \text{val}(S) \quad \forall S \subseteq \{1, \dots, p\}.$$

$$\phi_j = 0$$

Additivity: For a objective w/ combined payouts, ~~the~~ SHAP values are additive. ③

$$\phi_j + \phi_j^+ \Rightarrow \text{combined value.}$$

This is important, since for RF/ensemble models, additivity guarantees that if we calc SHAP val. for each indiv. tree & average \rightarrow SHAP for forest