

Decision Tree

2020年11月16日 13:40

Pros: good explanations (vs. Deep NNs)

- Cons:
- Big & Deep trees
 - Overfit easily.
 - Hard switch paths

Formalization:

$$f: \{-1, 1\}^n \rightarrow [-1, 1]$$

Two functions are close if: $\mathbb{E}_{x \sim D} [f(x) - g(x)]^2 \leq \epsilon$.

Fourier Basis: $\chi_S(x) = \prod_{i \in S} x_i$ for $S \subset [n]$

2ⁿ of them

construct an orthonormal basis for Boolean function

$$\Rightarrow f(x) = \sum_{S \subset [n]} \hat{f}_S \chi_S(x).$$

$$\text{where } \hat{f}_S = \langle f(x), \chi_S(x) \rangle \stackrel{\text{def}}{=} \mathbb{E}_{x \sim D} [f(x) \cdot \chi_S(x)]$$

Def L_1 norm $L_1(f) = \sum_S |\hat{f}_S|$.

Sparsity $L_0(f) = |\{S \mid \hat{f}_S \neq 0\}|$

Parseval's identity $\mathbb{E}_{x \sim D} [f(x)^2] = \sum_S \hat{f}_S^2$

Theorem

For any decision tree T with s leaf nodes, there exists a degree- $\log(\frac{s}{\epsilon})$ and sparsity- $(\frac{s}{\epsilon})$ function h that 2ϵ -approximates T .

Step 1. Truncate the tree at depth $\log(\frac{s}{\epsilon})$.

- There are $2^{\log \frac{s}{\epsilon}} = \frac{s}{\epsilon}$ nodes

- By union bound, at most $\frac{\epsilon}{s} \cdot s = \epsilon$ leaves cut

\Rightarrow Now assume T height $\leq \log(\frac{s}{\epsilon})$

Step 2. T with s leaves represented by $f: L_1(f)=s$ and degree $\log(\frac{s}{\epsilon})$

- A tree with s leaf nodes can be represented by 最小项之和

- Each 最小项 $L_i \leq 1$

$$\Rightarrow L_1(f) \leq s$$

$$\Rightarrow \text{Degree} \leq \text{height} = \log(\frac{s}{\epsilon})$$

Step 3. For f s.t. $L_1 \leq s$, $\text{deg} \leq \log(\frac{s}{\epsilon})$. $\exists h$ with $L_0 \leq \frac{s^2}{\epsilon}$, $\text{deg} = \log(\frac{s}{\epsilon})$ s.t.

$$\sum_{x \in D} (h(x) - f(x))^2 \leq \epsilon$$

Let h include all terms in

$$\Lambda = \left\{ s \mid |f_s| \geq \frac{\epsilon}{L_1(f)} \right\}$$

h has terms at most

$$\frac{L_1(f)}{\frac{\epsilon}{L_1(f)}} = \frac{L_1(f)^2}{\epsilon} \leq \frac{s^2}{\epsilon}$$

Meanwhile,

$$\sum_{s \notin \Lambda} (\hat{f}_s)^2 \leq \max_{s \notin \Lambda} \hat{f}_s \cdot \sum_{s \notin \Lambda} |f_s| \leq \frac{\epsilon}{L_1(f)} \cdot L_1(f) = \epsilon$$

KM algorithm.

- prune less promising branches
- \underline{f}_α the sum of all Fourier basis started by α
 e.g. $f_{110} = \hat{f}_{x_1 x_2 x_4} \cdot x_1 x_2 x_4 + \hat{f}_{x_1 x_2} \cdot x_1 x_2$

Algo.

Def Coef(α):

If $\mathbb{E}(f_\alpha^2) \geq \theta^2$:

If $|\alpha| = n$:

outputs α

Else:

Coef(α_0), Coef(α_1)

LMN algorithm

LMN algorithm

Take m random uniform samples for f

Estimate:

$$\hat{f}_s \triangleq \langle f, \chi_s \rangle = \mathbb{E}_{x \sim D} [f(x) \chi_s(x)] \simeq \frac{\sum_{i=1}^m f(x_i) \chi_s(x_i)}{m}$$

Compressed Sensing

$$y = Ax + e$$

measurement \rightarrow error

x is sparse in χ_s basis ($\frac{s}{2}$).

How to pick random variables?

- Gini index

If labels in n classes.

$$Gini = 1 - \sum_{i=1}^n P_i^2 \in [0, 1)$$

$$P_i = \Pr[\text{left branch \& result} = i]$$

Gini index = 0 \rightarrow all in one class

$$Gini = \Pr[\text{left}] \cdot Gini(\text{left}) + \Pr[\text{right}] \cdot Gini(\text{right})$$

pick the variable with smallest Gini

Overfitting \Rightarrow Random Forest (Most practical)

• Bagging

$(x_1, y_1), \dots, (x_n, y_n)$ instead of use it directly,

we sample it n times with replacement $\Rightarrow X_b \quad b=1, 2, \dots, B$

Each Element: miss probability $(1 - \frac{1}{n})^n \simeq \frac{1}{e}$

Get B trees (forest) using B dataset

$$\Rightarrow \hat{f}(X_{\text{test}}) = \frac{1}{B} \sum_{b=1}^B f_b(X_{\text{test}})$$

• Feature Bagging (ensemble)

Also do bagging on features

Each tree only use a random subset of features

Forced to use all features to learn!

• Boosting

(n learners)

y

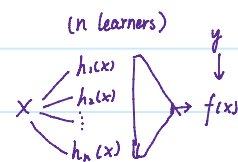
given w use an ensemble w learn:

• Boosting

Combine weak learners to build a strong one

other possible methods:

- Stacked Learner



Modify the sample likelihood of data with wrong prediction

Framework. for $t=1$ to T

construct distribution D_t on $\{1, \dots, m\}$

Find classifier h_t ^(weak), error $\epsilon_t = \Pr_{i \sim D_t} [h_t(x_i) \neq y_i]$

Output final classifier H_{final}

eg. Adaboost. $D_1(i) = \frac{1}{m}$

$$\text{Given } D_t, h_t \Rightarrow D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } y_i = h_t(x_i) \\ e^{\alpha_t} & \text{if } y_i \neq h_t(x_i) \end{cases}$$

$Z_t = \text{normalize const}$

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

$$H_{\text{final}} = \text{sign} \left(\sum_t \alpha_t h_t(x) \right)$$

Theorem write $\epsilon_t = \frac{1}{2} - \gamma_t$

$$\Rightarrow \text{Error}_{\text{train}}(H_{\text{final}}) \leq \prod_t 2\sqrt{\epsilon_t(1-\epsilon_t)} \leq e^{-2\sum_t \gamma_t^2}$$

Proof.

$$\text{Error} = \frac{1}{m} \sum_i \begin{cases} 1 & \text{if } y_i \neq H_{\text{final}}(x_i) \\ 0 & \text{else} \end{cases}$$

$$= \frac{1}{m} \sum_i \begin{cases} 1 & \text{if } y_i f(x_i) \leq 0 \\ 0 & \text{else} \end{cases}$$

$$\leq \frac{1}{m} \sum_i e^{-y_i f(x_i)} = \sum_i D_{\text{final}}(i) \prod_t Z_t = \prod_t Z_t$$

$$D_{\text{final}}(i) = \frac{1}{m} \frac{e^{-y_i \sum_t \alpha_t h_t(x_i)}}{\prod_t Z_t} = \frac{1}{m} \frac{e^{-y_i f(x_i)}}{\prod_t Z_t}$$

$$Z_t = \sum_i D_t(i) e^{-\alpha_t y_i h_t(x_i)}$$

$$= \epsilon_t e^{\alpha_t} + (1 - \epsilon_t) e^{-\alpha_t}$$

$$= 2\sqrt{\epsilon_t(1-\epsilon_t)}$$

$$\Rightarrow \epsilon_t = \frac{1}{2} - \gamma_t \Rightarrow Z_t = \sqrt{1 - 4\gamma_t^2} \leq e^{-2\gamma_t^2}$$

$$\Rightarrow \text{Error} \leq e^{-2\sum_t \gamma_t^2}$$

Generalization