

COL744
Machine learning
Nov 29, 2020

Last Class:-

Decision Trees.



$$\text{Entropy} = H(Y) = - \sum_y P(Y=y) \log P(Y=y)$$

↳ Measure to decide which node to split on.

$$MI(Y, X_j) = H(Y) - H(Y|X_j)$$

$\underbrace{\quad}_{\text{in } X_j}$ which minimizes $MI(Y, X_j)$

Gini Index

$$E_{Gini}(Y) = 1 - \sum_y [P(Y=y)]^2$$

$$= 1 - \sum_y P(Y=y) \cdot \underbrace{P(Y=y)}$$

Choose the attribute X_j which results in max reduction in $Gini(Y|X_j)$

$$\arg \max_{X_j} [Gini(Y) - Gini(Y|X_j)]$$

Note:- $Y \in \{0, 1\}$ Let $P(Y=1) = P_1$ $P(Y=0) = P_0$

$$Gini(Y) = 1 - [P_1^2 + P_0^2]$$

$$= 1 - \{ P_1^2 + (1-P_1)^2 \}$$

$$= 1 - [P_1^2 + 1 + P_0^2 - 2P_1]$$

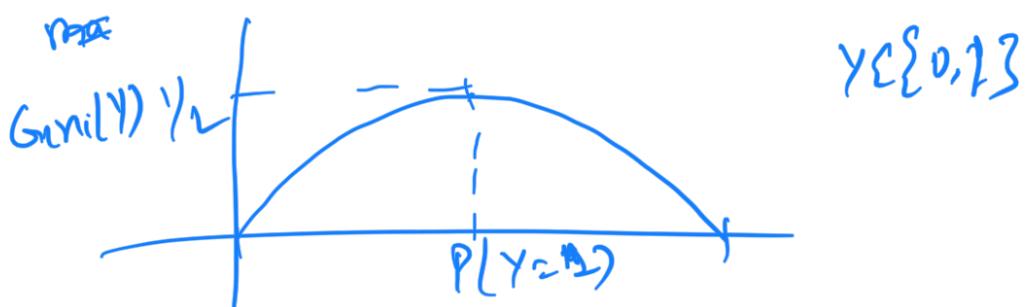
$$= 2P_1 - 2P_1^2 = 2P_1(1-P_1)$$

↳ Which bin is minimized when

$$\text{at } P_1 = 1/2$$

At extremes $Gini(Y) \rightarrow 0$

$P_0 \rightarrow 0$
or $P_1 \rightarrow 0$



Regression Problem:- $y \in \mathbb{R}$ vs $Gini(Y)$ (Regression Tree) & Regression Tree

What is right measure
to split on attribute?
↳ Entropy may not be the
right measure.

For each node (after splitting)
 $\hat{y}_e = \frac{1}{|D_e|} \sum_{i=1}^m y_e^{(i)}$
 Avg. prediction
 $J_{x_d} = \frac{1}{|D|} \sum_{i=1}^m \sum_{j=1}^{|D|} (y_j - \hat{y}_e)^2$ $\hat{y} = \frac{1}{|D|} \sum_{i=1}^m y_e^{(i)}$
 $J_{x_d=e} = \frac{1}{|D_e|} \sum_{i=1}^{|D_e|} (y_e^{(i)} - \hat{y}_e)^2$

$$J_{x_d \text{ split}} = \sum_e P(x_d=e) \cdot [J_{x_d=e}]$$

Look for that attribute which leads to
minimization in error

$$J_{X_j} - J_{X_j, \text{split}}$$

choose

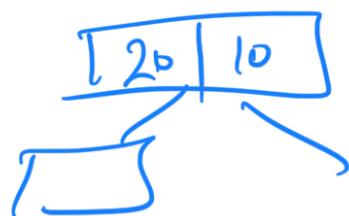
stopping criterion: when do you stop?

↳ stop when $|D_t| \leq \text{min_thresh}$

Finally: splitting on cont. Attribute.

X_j : continuous value (Humidity)

$$X_j \in [0, 100]$$



which ever split leads

to max. reduction

in entropy (variance)

↳ choose that split

Consider a split
at $X_j = 20$
value for which
we have a training
data point

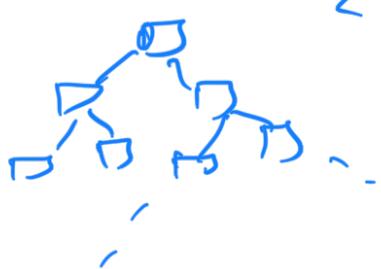
Random Forests:-

Collection of Decision Trees

$$\{\mathbf{x}_i^{(t)}, y_i^{(t)}\}_{i=1}^m$$



Grow Multiple trees
on small variations
of the data



"Same" \rightarrow

$$\begin{cases} D^{(1)}, D^{(2)}, \dots, D^{(T)} \\ |D^{(t)}| = m \\ M^{(1)}, M^{(2)}, \dots, M^{(T)} \end{cases}$$

Bootstrapped Sample

$$H_t \sum_{D^{(t)} = S^3}^1 \dots T_3 \sum_{1 \leq i \leq m}$$

sample with replacement

$$\Leftrightarrow (\underline{x^{(t)}, y^{(t)}}) \sim \{\underline{x^{(i)}, y^{(i)}}\}_{i=1}^m$$

$$\underline{D^{(t)}} = \underline{S^3} \cup \{\underline{x^{(i)}, y^{(i)}}\}$$

3

$$\{\underline{x^{(i)}, y^{(i)}}\}_{i=1}^m$$

Learning "ensemble" of models on Bootstrapped samples

Decision Trees | Random

Random Forest

(collection of decision trees)

Create a T trees, one for every bootstrapped sample of data

Random Forest

① Create

Bootstrapped sample $D^{(1)} - ST$

② Grow

T trees on $D^{(1)}$ referred to as $M^{(1)} - ST$, \dots , $M^{(T)} - ST$ } y_i Model

③ At

test time:-

$$p^{(t)}(x) = M^{(t)}(x)$$

$$h^{\text{avg}}(x) = \frac{1}{T} \sum_{t=1}^T p^{(t)}(x)$$

Prediction of ensemble to test

" " " "

prediction
corresponding
to t th model

Bagging: "Barn leaves"
 $\{x_i, y_i\}_{i=1}^m$ Useful when the original model is high variance.

20	10
----	----

X_d :- But split on Attr. to
 ↓ Grow further $X_d + X_d'$
 "unstable" learners.

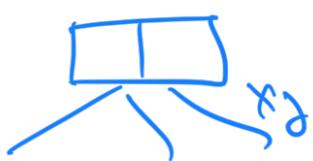
\Rightarrow A high amount of variance in model lead we are going to learn ("overfitting")

"Bagging" :- Reduces this variance by forming an ensemble of models of bootstrapped samples

Specif. Use: Random Forests

Use "Bagging" prediction { ① Learn $M^1 \dots M^T$ modls on Bootstrapped Samples

② While growing each tree M^t split on a subset of randomly chosen attributes.



$$X_d^* = \arg \min_{X_d \subseteq \{x_1, \dots, x_n\}} \sum_{i=1}^m L(Y_i, X_d)$$

$\{x_i, y_i\}_{i=1}^m$

21	9
----	---

X_d :- But attr. to split on
 ↓ Grow further

"Robust" to
noise

choose $\frac{1}{T}$ a subset of attributes

$$x_d^* = \underset{x_F}{\text{argmax}} \underline{\text{MI}}(Y, X_d)$$

F#... # of
different attributes
split m

$$\Leftrightarrow x_F \subset \{x_1, \dots, x_n\}$$

\Rightarrow Helps from "Different" feature.
"Variance" :- Overshooting

③ Out of Bag error :-

$\hookrightarrow M(t)$:- $\hat{x}^{(t)}$:- example
Not chosen as part
of bootstrapped set $D^{(t)}$

validation
set that
for first tree \hookleftarrow $\left[\begin{array}{l} \hat{x}^{(t)} := \{x^{(t)}, y^{(t)}\} \\ \text{s.t. } (x^{(t)}, y^{(t)}) \notin D^{(t)} \end{array} \right]$

$$J^{(t)} = \sum_{i=1}^N \mathbb{1}_{\{f_i^{(t)}(x^{(t)}) \neq y^{(t)}\}}$$

\Downarrow out of bag error

A Total out of bag error

$$= \frac{1}{T} \sum_{t=1}^T J^{(t)}$$

Gradient Boosted Decision Trees :-

"In - 1 min" :- "weak learners" :- slightly

BOOSTING

better than random

$$\frac{M_a^1}{M_a^1 - M_T} \xrightarrow{\text{loss}} M(t)$$

ensemble of learners

$M(t)$ is learned by focusing
on errors made by
previous model

$$M^{(1)} = \text{argmin}_{M^{(1)}} \sum_{i=1}^n \ell(y_i, M^{(1)}(x_i))$$

Search in space y

$M^{(1)}$

Minimizing the
loss over

$$\text{argmin}_{M^{(1)}} \sum_{i=1}^n \ell(y_i, M^{(1)}(x_i))$$

$$\sum_{i=1}^n \ell(y_i, M^{(1)}(x_i))$$

squared loss

$$M^{(1)}(x_i) = \sum_{t=1}^T M_t(x_i)$$

Functional Gradient
Boosting

$$\sum_{i=1}^n \ell(y_i, M^{(1)}(x_i))$$