# TREE MODELS

$$\mathbb{1}\left[x \geq 1.1\right]$$

$$\mathbb{1}\left[x_2 \geq 0\right]$$

$$\mathbb{1}\left[x_1 \geq 0.5\right]$$

$$\mathbb{1}\left[x_2 \geq 1\right]$$

$$\mathbb{1}\left[x_1 \leq 0\right]$$

$$\mathbb{1}\left[x_2 \geq 0\right]$$

F / \ T    F / \ T

1      0    1      0

d  NUMBER OF FEATURES

n  NUMBER OF EXAMPLES

$$\mathcal{X} = \bigcup_{i=1}^{n} R_i$$

$$R_j \cap R_i = \varnothing \quad (i \neq j)$$

$$R_i \cup R_j = R_p \quad \text{PARENT} \left( \underline{\text{ex}} \ \mathbb{R}^d \right)$$

WE CAN DEFINE A REGION IN THE FOLLOWING WAY

$$R_1 = \left\{ X \mid X_j < t, \ X \in R_p \right\}$$

$$R_2 = \left\{ X \mid X_j \geq t, \ X \in R_p \right\} \longrightarrow$$

AND CONTINUE TO SPLIT THE SPACES IN ORDER TO CLASSIFY, IN THIS CASE

WE CAN DEFINE A LOSS FUNCTION AS A SET FUNCTION ON A REGION $\mathbb{R}$ GIVEN A SPLIT OF A $\mathbb{R}_p$ INTO $\mathbb{R}_1, \mathbb{R}_2$ WE CAN COMPUTE $\mathcal{L}(\mathbb{R}_p)$ AS WELL AS THE CARDINALITY WEIGHTED-LOSS OF THE CHILDREN,

WE SELECT THE LEAF REGION, FEATURE, THRESHOLD THAT MINIMIZE THE LOSS

$$J(\mathbb{R}_p) - \frac{|\mathbb{R}_1| J(\mathbb{R}_1) + |\mathbb{R}_2| J(\mathbb{R}_2)}{\mathbb{R}_1 + \mathbb{R}_2}$$

FOR A CLASSIFICATION PROBLEM, WE FOCUS ON $\mathcal{L}_{\text{MISSCLASS}}$. FOR A REGION $\mathbb{R}$ LET $\hat{p}_c$ BE THE PROPORTION OF EXAMPLES IN $\mathbb{R}$ THAT ARE OF CLASS $c$. MISSCLASSIFICATION LOSS ON $\mathbb{R}$ CAN BE EXPRESSED AS

$$J_{\text{MISSCLASS}}(\mathbb{R}) = 1 - \max_c (\hat{p}_c)$$

THE N OF EXAMPLES THAT WOULD BE MISSCLASSIFIED IF WE PREDICTED THE MAJORITY CLASS IN $\mathbb{R}$

- The first split is isolating out more of the positives, but we note that:

$$L(R_p) = \frac{|R_1|\,L(R_1) + |R_2|\,L(R_2)}{|R_1| + |R_2|}$$

$$= \frac{|R_1'|\,L(R_1') + |R_2'|\,L(R_2')}{|R_1'| + |R_2'|} = 100$$

- Thus, not only can we not only are the losses of the two splits identical, but neither of the splits decrease the loss over that of the parent.

WE NEED TO DEFINE A MORE SENSITIVE LOSS. LET'S FOCUS ON CROSS-ENTROPY

$$J_{CE}(R) = - \sum_c \hat{P}_c \log_2 \hat{P}_c$$

$\hat{P} \log_2 \hat{P} = 0$ IF $\hat{P} = 0$, N OF BITS NEEDED TO SPECIFY THE OUTCOME OR CLASS GIVEN A KNOWN DISTRIBUTION

$J_{REDUCTION}$ FROM $P \rightarrow C$ = INFORMATION GAIN

# ALGORITHM

TREE FITTING

$$\begin{cases} \text{FIND } j \in \{1, \ldots, d\}, \theta \in \mathbb{R} \\ \underset{j, \theta}{\text{argmax}} \; G(j, \theta) \quad (G : \text{"gain"}) \end{cases}$$

$S \rightarrow$ left $\qquad$ righ

$$\mathbb{1}\left(x_j^{(i)} \leq \theta\right) \qquad \mathbb{1}\left(x_j^{(i)} > \theta\right)$$

RETURN $S$

# PROS

- WELL- INTERPRETABLE
- ROBUST TO OUTWIRES
- HANDLES MIX OF DISCRETE AND CONTINUOUS FEATURES
- ROBUST TO MONOTONE TRANSFORMATION
- CAN FIT QUICKLY

# CONS

- GENERALIZE POORLY
- HIGHLY UNSTABLE

# ENSEMBLE LEARNING

$$F_1, \ldots, F_m$$

$$f(x) = \sum_{i=1}^{n} \beta_i F_i(x) \qquad \beta_i \in \mathbb{R}, \quad \beta_j = \frac{1}{m}$$

THIS IS USEFUL TO RESOLVE UNSTABILITY PROBLEMS OF CERTAIN SUPERVISED ALGORITHMS

THE TRAIN IS DETERMINISTIC, IF I FEED THE SAME INPUTS AND AVG, WE GET THE SAME RESULT

## BAGGING (BOOTSTRAP AGGREGATION)

$$\{S\} = n \qquad \underline{ex} \{1, 3, 5, 7, 7, 8\}$$

$$\uparrow$$
$$x^{(i)}, y^{(i)}$$

$$P(X = 1) = \frac{1}{6}$$

$$P(X = 7) = \frac{1}{3}$$

A BOOTSTRAP IS A SUBSET $S'$ WE TRAIN AND MAKE DIFFERENT PREDICTIONS. THEN LINEARLY COMBINE THEM
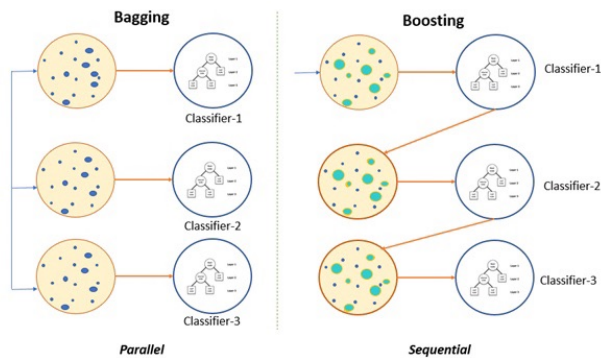
$Y$ : # OF UNIQUE DATAPOINTS FROM $S$

$X_1, \ldots, X_n$ : WHETHER THE $i$-TH DATAPOINT IS IN A BOOTSTRAPPED SET

$$\mathbb{E}[Y] = \sum_{i=1}^{n} [X_i] = \sum_{i=1}^{n} P(X_i = 1)$$

$$P(X_i = 0) = \left(1 - \frac{1}{n}\right)^n$$

$$P(X_i = 1) = 1 - \left(1 - \frac{1}{n}\right)^n = \underbrace{\qquad}_{\approx 0,632} c^{-1}$$

$$\lim_{n \to +\infty} \left(1 + \frac{c}{n}\right)^n = e^c$$



Bagging      Boosting

Classifier-1

Classifier-2

Classifier-3

Classifier-1
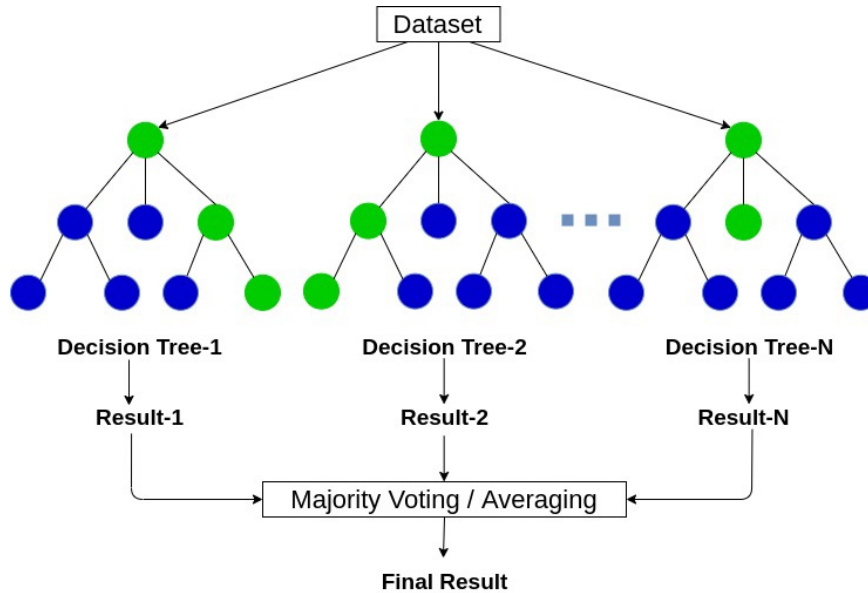
Classifier-2

Classifier-3

Parallel      Sequential

# RANDOM FOREST

→ BAGGING

→ SAMPLE A SUBSET OF FEATURES FOR SPLITS

WE HAVE MANY TREES

# BOOSTING

$$f_{i-1} = \sum_{j=1}^{i-1} \beta_j F_j \rightarrow f_i = f_{i-1} + F_j = \sum_{j=1}^{i} \beta_j F_j$$

WE ARE BASICALLY COMBINING BAD CLASSIFIERS INTO ONE GOOD, BY AMALGAMATE THEM.

EX IN CLASSIFIERS

$$\delta, \; \varepsilon, \; \gamma \; \in \; (0, 1)$$

TOL OF UNCERTAINTY    TOL OF ERROR    IMPROVEMENT IN ACCURACY    MATTER OF TIME UNTIL GETS TO PERF. CLASS.

$\gamma$-WEAR LEARNERS    VS    STRONG LEARNERS

LOSS $P(Acc(f) > \frac{1}{2} + \gamma) \geq 1 - \delta$    $P(J(f) < \varepsilon) \geq 1 - \delta$ if $> m(\delta, \varepsilon)$

MORE WEAR INTO STRONG

# FORWARD STAGEWISE ADDITIVE MODELING (ALG)

AT ITER $i$, LOOK FOR $F_i$ (PARAMETRIZED BY $\theta$)

$$\beta_i, \theta_i = \arg\min_{\beta, \theta} \mathcal{L}(\underbrace{f_{i-1}}_{\text{CONSTANT}} + \beta F_i)$$

$$= \arg\min \sum_{i=1}^{n} \ell\left(y^{(i)}, f_{i-1}(x^{(i)}) + \beta F_i(x^{(i)}; \theta)\right)$$

$$F_i = F_i(\cdot \; ; \theta_i)$$

$$f_i(x) = f_{i-1}(x) + \beta_i F_i(x)$$

$$F_i = \arg\min_{F} \mathcal{L}(f_{i-1} + F) = \arg\min_{F} \sum_{i=1}^{n} \ell\left(y^{(i)}, f_{i-1}(x^{(i)}) + F(x^{(i)})\right)$$

WE ARE OPTIMIZING W.R.T FUNCTION, WE IMAGINE FUNCTIONS AS CONTINUOUS DIMENSIONAL VECTORS

LOSS FOR GB

$$\mathcal{L}(f) = \sum_{i=1}^{n} l\left(y^{(i)}, f(x^{(:)})\right) = \sum_{j=1}^{K} \sum_{x^{(i)} \in R_j} l\left(y^{(i)}, w_j\right)$$

X GB

$$\mathcal{L}(f) = \sum_{i=1}^{n} l\left(y^{(i)}, f(x^{(:)})\right) = \sum_{j=1}^{K} \sum_{x^{(i)} \in R_j} l\left(y^{(i)}, w_j\right) + \gamma \bar{J} + \frac{1}{2} \lambda \sum_{j=1}^{J} w_j^2$$

N OF LAYERS ↑

PREDICTION ↑

WE "DISCRETIZE" AND OBTAIN

$$f = \begin{bmatrix} f_{i-1}(x^{(i)}) \\ \vdots \\ f_{i-1}(x^{(n)}) \end{bmatrix} \xrightarrow{\text{AND OBTAIN}} g = \begin{bmatrix} \dfrac{\partial \ell(y^{(i)}, f_{i-1}(x^{(1)}))}{\partial f_{i-1}(x^{(1)})} \\ \vdots \\ \dfrac{\partial \ell(y^{(n)}, f_{i-1}(x^{(n)}))}{\partial f_{i-1}(x^{(n)})} \end{bmatrix}$$

<span style="color:blue">DERIVATIVE OF LOSS W.R.T. PREDICTION</span>

WE ATTEMPT TO RE-CREATE GRADIENT DESCENT U.R. IN A CONTINUOUS SPACE

$$f = f - \alpha g$$

WE WANT TO FIT A MODEL FOR $g$

$$F_i \approx \arg\min_F \sum_{j=1}^{n} (-g_j - F(x^{(1)}))^2$$

DECISION - TREES FIT THE MODEL VERY NICELY, SINCE THEY ARE THE QUINTESSENTIAL OF WEAK LEARNERS ( WE NEED QUICK LEARNERS )

YOU LEARN THE GRADIENT AND ADD TREE BY TREE


VARIANTS:

→ EXTREME GRADIENT BOOST
→ NG BOOST