

Trees and Rules and Optimization

M. Hakan Akyüz

Açıklanabilir Yapay Öğrenme - II

Outline

- Decision Trees (DTs)
 - CART
- Optimization Models
 - Optimal Classification Trees
 - State-of-the-art
- Rules and Rule Set Generation
 - Column Generation
 - Extensions
- Takeaways

Regression Trees



Prediction: Mean of the response values for the training observations in R_m

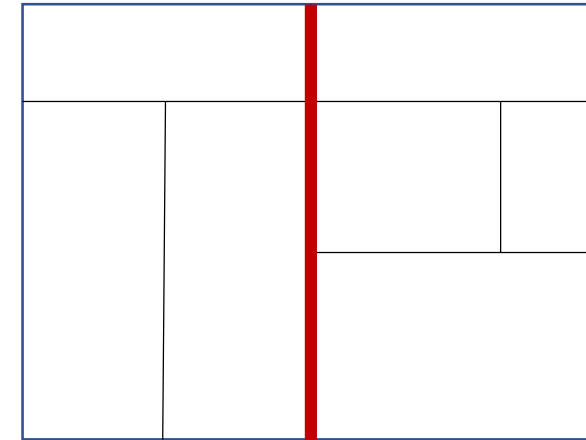
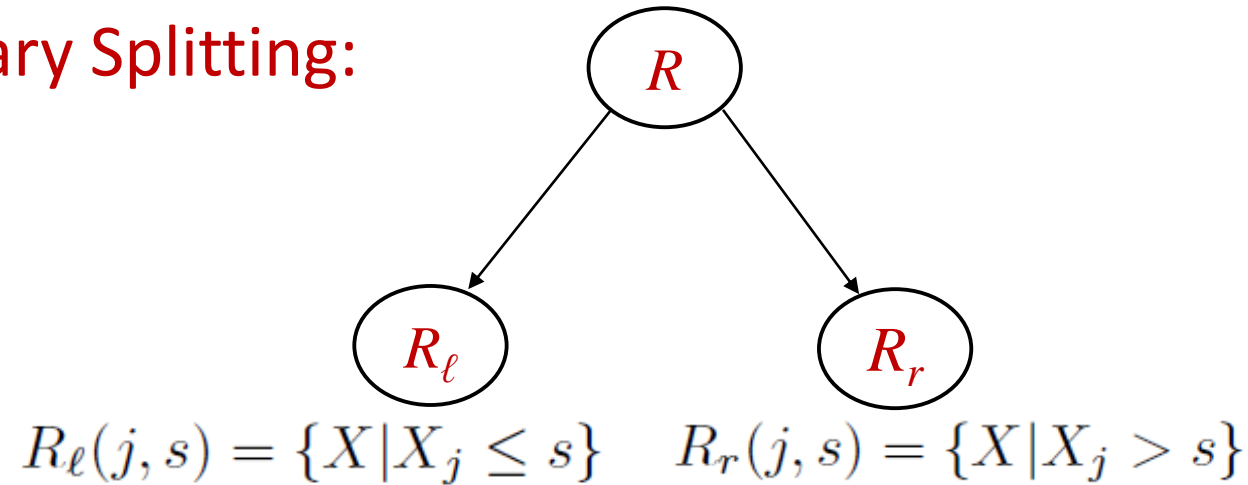
$$\hat{y}_m = \frac{1}{N_m} \sum_{\{x_i \in R_m\}} y_i$$

Goal: Finding the regions such that sum of squares is minimized

$$\sum_{j=1}^m \sum_{\{x_i \in R_j\}} (y_i - \hat{y}_j)^2$$

How to grow a Regression Tree?

Recursive Binary Splitting:



Find j and s minimizing the sum of squared error

$$\sum_{\{x_i \in R_\ell(j, s)\}} (y_i - \hat{y}_\ell)^2 + \sum_{\{x_i \in R_r(j, s)\}} (y_i - \hat{y}_r)^2$$

Until when?

Goal: Avoiding overfitting with a fully grown or large tree
(Selecting a subtree that leads to a lowest test error rate)

CART:

- 1) Grow a large tree until each leaf has a minimum sample size, e.g. 5.
- 2) Pruning large tree based on **cost-complexity**, $C_\alpha(T)$.

$$C_\alpha(T) = \sum_{m=1}^{|T|} \sum_{\{x_i \in R_m\}} (y_i - \hat{y}_m)^2 + \alpha |T|$$

Pruning: Successively collapse the internal node that produces the smallest per-node increase in sum of squared errors.
- Use k-fold cross validation to tune α .

Classification Trees

Similar to a regression tree but uses different error measures based on *impurity* of a region.

Prediction: Majority voting of classes among the training observations in R_m

$$\hat{p}_{mk} = \frac{1}{N_m} \times \sum_{x_i \in R_m} I(y_i = k)$$

$$\hat{y}_m = \arg \max_k \hat{p}_{mk}$$

Error Measures

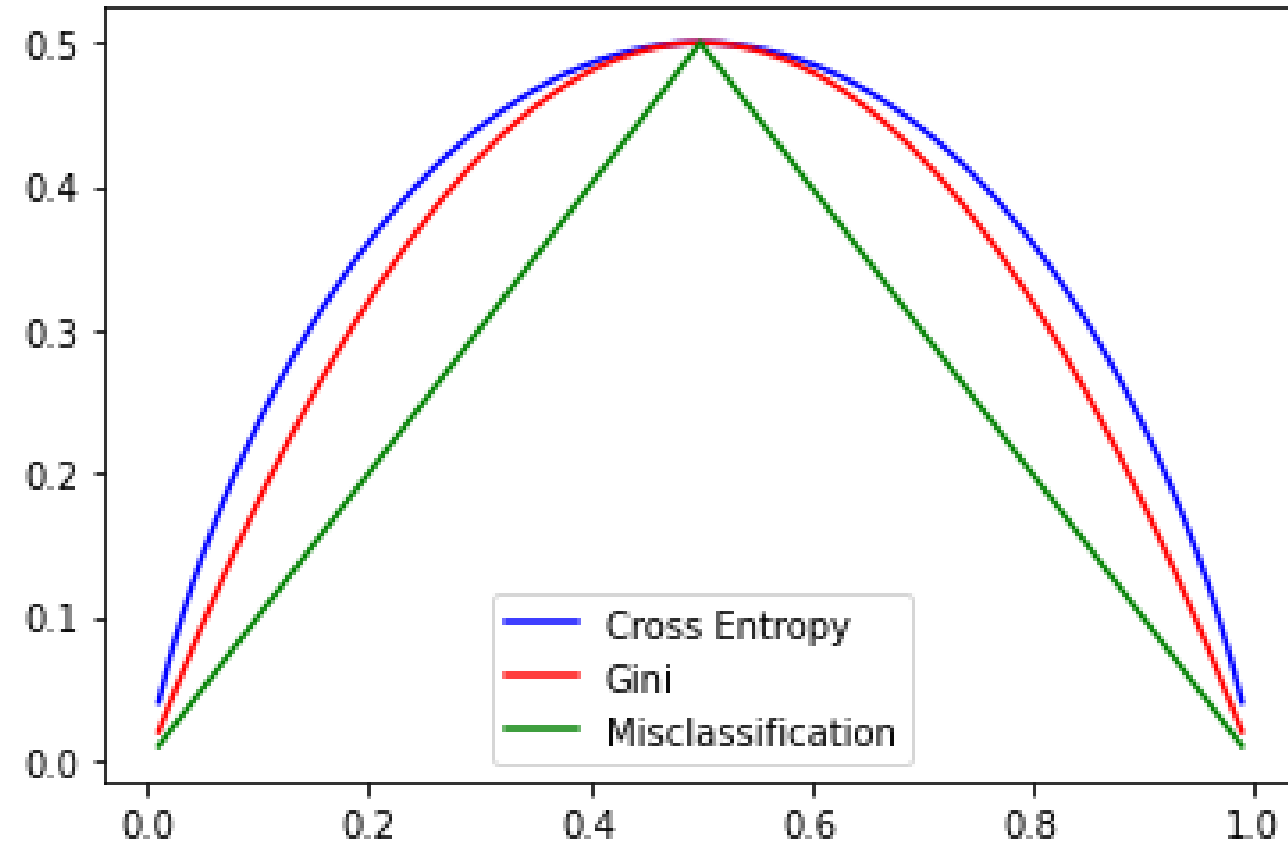
$$\frac{1}{N_m} \times \sum_{x_i \in R_m} I(y_i \neq \hat{y}_m) = 1 - \hat{p}_{m\hat{y}_m} \longrightarrow \text{Misclassification Error}$$

$$\sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}) \longrightarrow \text{Gini Index}$$

$$-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk} \longrightarrow \text{Cross Entropy}$$

Which error measure?

- Use Gini or Cross Entropy to grow the tree.
- For pruning use misclassification error.
Exercise: Why not the other way?
- Weight the node impurity measures by the number $N_{m\ell}$ and N_{mr} of observations in the two child nodes created by splitting node m .



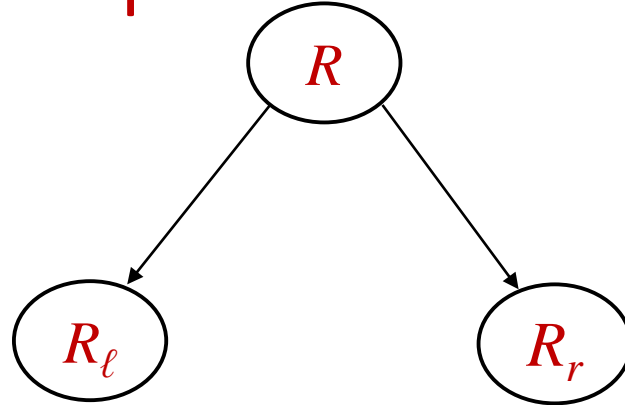
* Two class case. Cross Entropy is scaled to pass through (0.5,0.5).

$$\frac{1}{N_m} \times \sum_{x_i \in R_m} I(y_i \neq \hat{y}_m) = 1 - \hat{p}_m \hat{y}_m \longrightarrow \text{Misclassification Error}$$

$$\sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}) \longrightarrow \text{Gini Index}$$

$$-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk} \longrightarrow \text{Cross Entropy}$$

How to select a split?



$$R_\ell(j, s) = \{X | X_j \leq s\} \quad R_r(j, s) = \{X | X_j > s\}$$

Find j and s minimizing the error (or maximizing the gain of split)

$$\sum_{\{x_i \in R_\ell(j, s)\}} (y_i - \hat{y}_\ell)^2 + \sum_{\{x_i \in R_r(j, s)\}} (y_i - \hat{y}_r)^2$$

Classification Trees

$$\hat{p}_{mk} = \frac{1}{N_m} \times \sum_{x_i \in R_m} I(y_i = k)$$
$$\hat{y}_m = \arg \max_k \hat{p}_{mk}$$

Gini Index:

$$\sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}) \quad \Leftrightarrow$$

$$1 - \sum_{k=1}^K \hat{p}_{mk}^2$$

What is the **gain** of a split "**s**"?

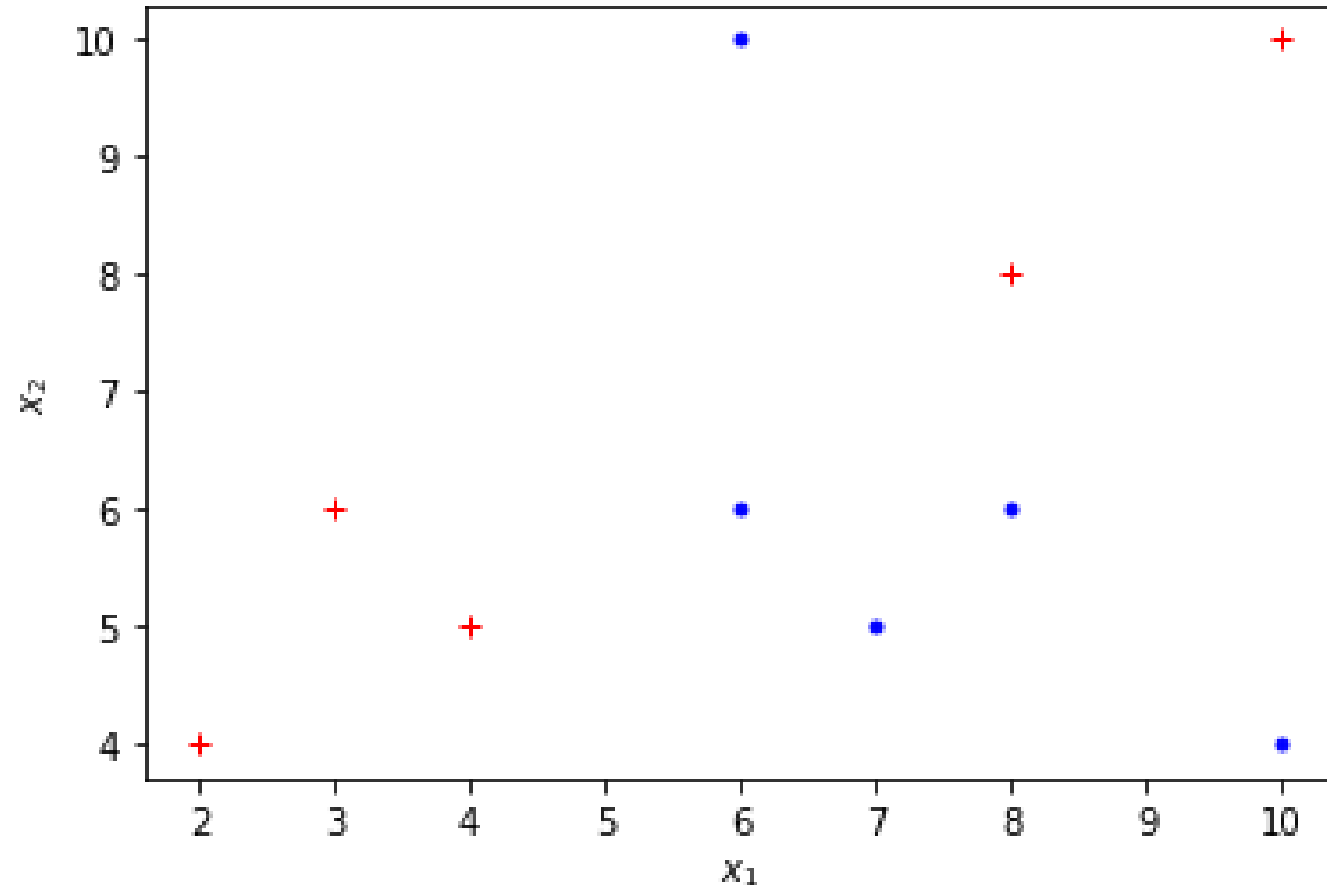
Gain = (Error before the split) – (Error after the split)

$$G(s) = \boxed{E(R)} - \boxed{E(s)} \rightarrow \boxed{E(s) = \frac{N_{ml}}{N_m} E(R_l) + \frac{N_{mr}}{N_m} E(R_r)}$$

Diagram illustrating the components of the Gain formula:

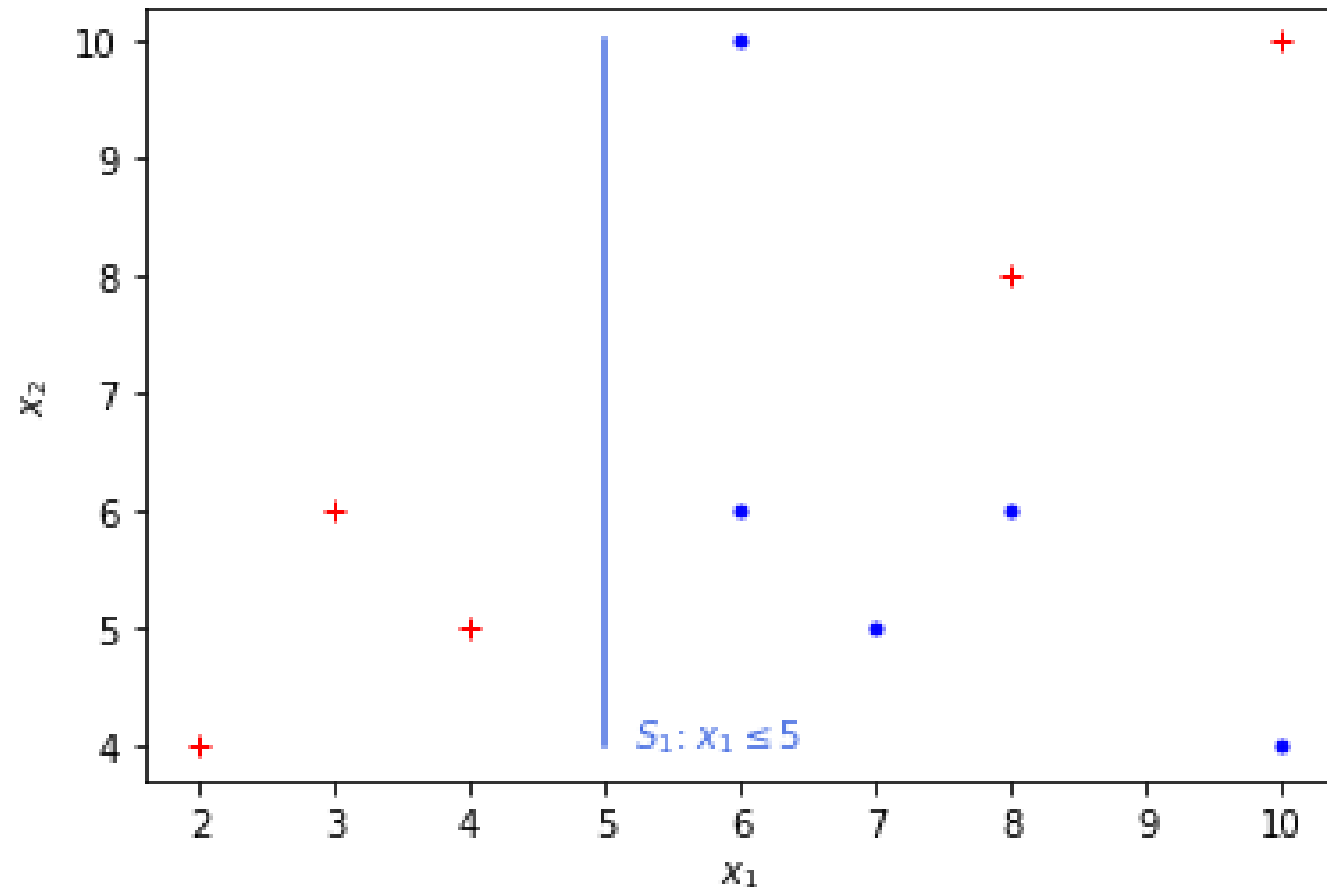
- $\boxed{E(R)}$ points to **Error at node *R***
- $\boxed{E(s)}$ points to **Error of split *s***

Classification – $\{ +, \bullet \}$



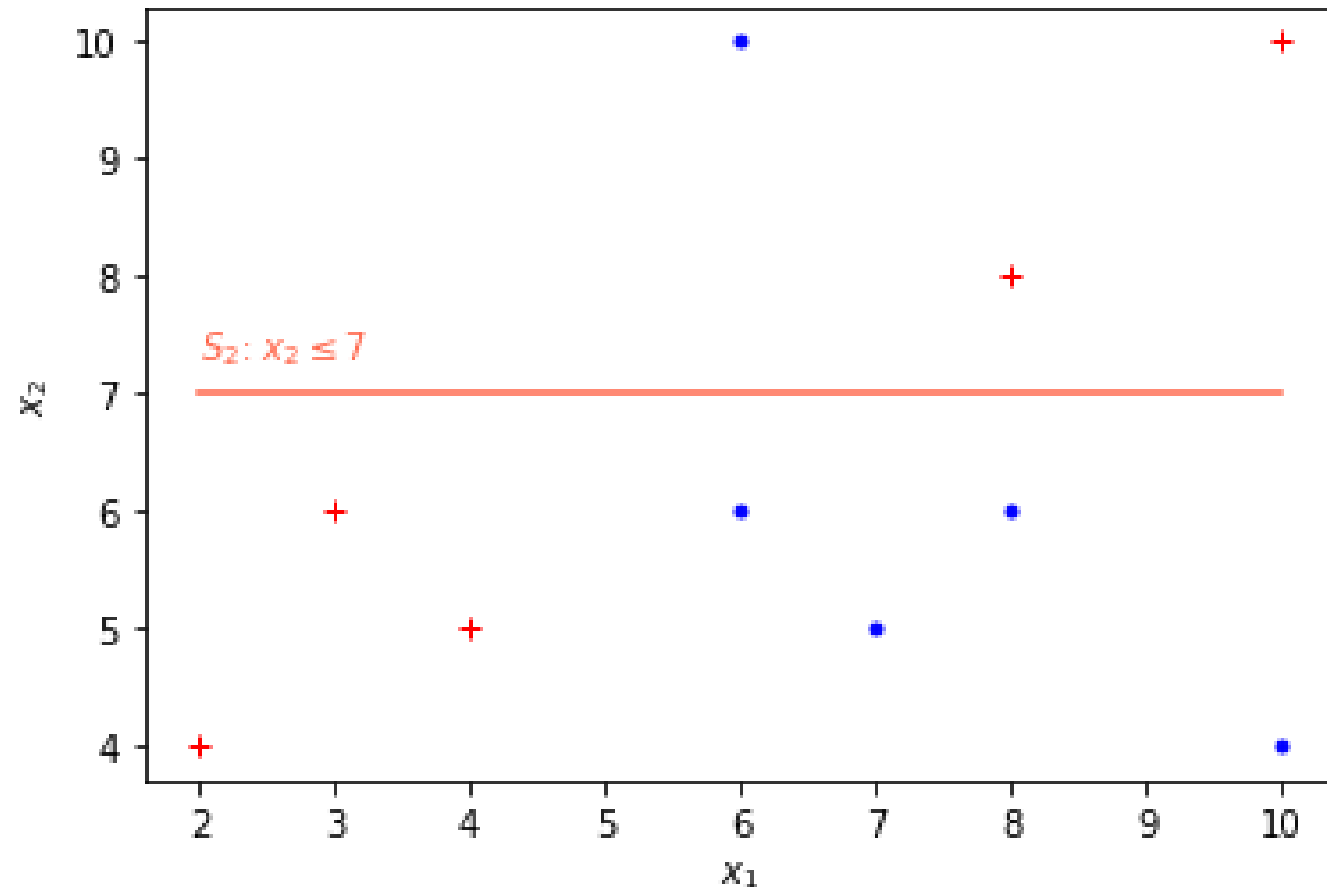
Classification – $\{+, \bullet\}$

Split #1: $x_1 \leq 5$



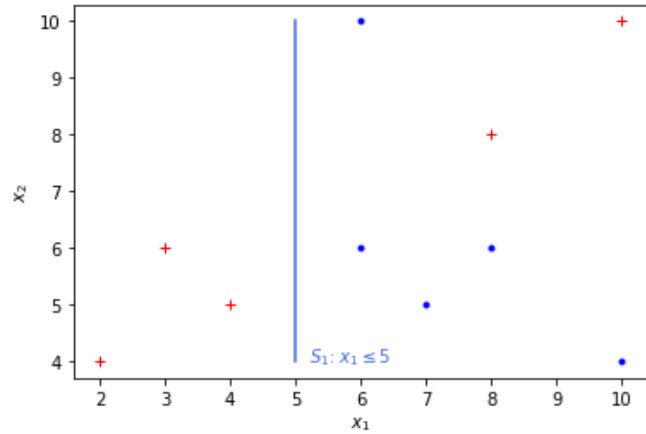
Classification – $\{ +, \bullet \}$

Split #2: $x_2 \leq 7$



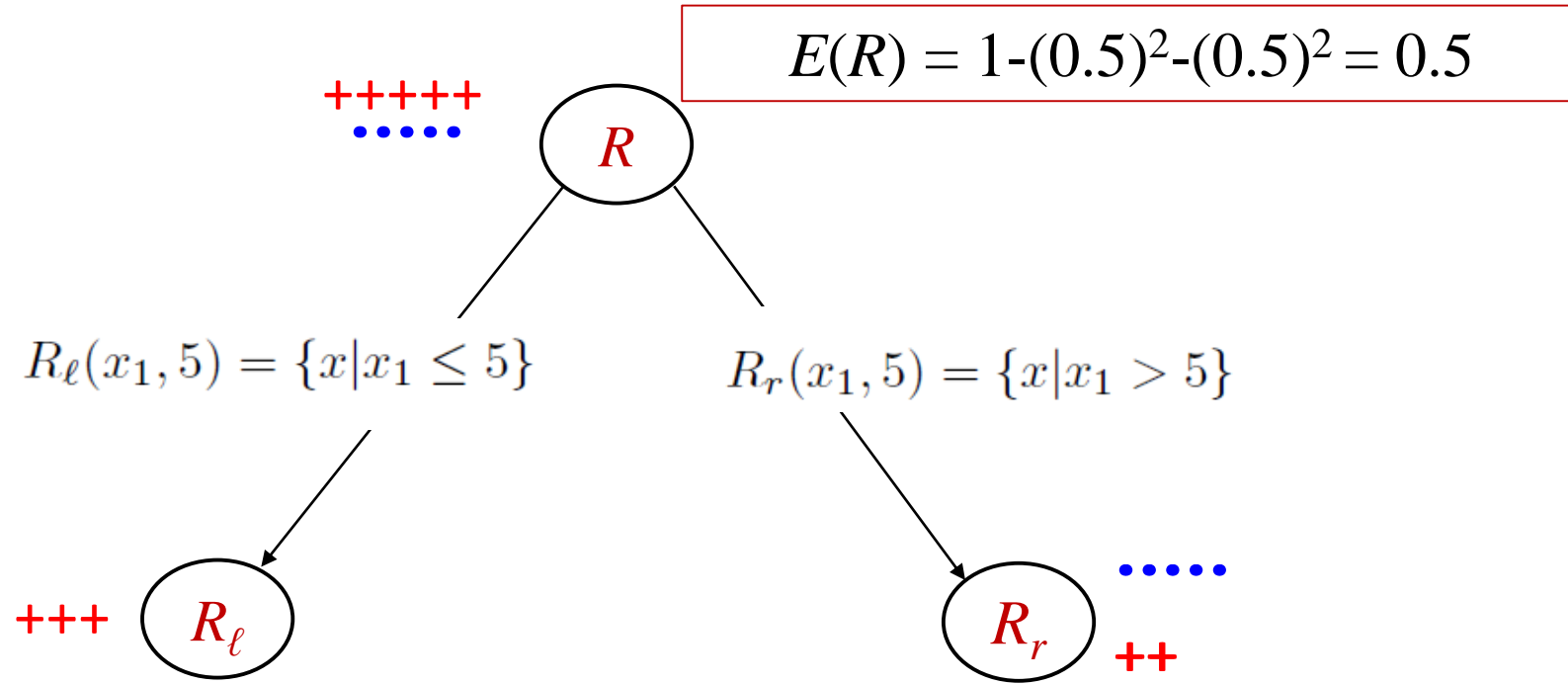
Classification – $\{+, \cdot\}$

Gini Index: $1 - \sum_{k=1}^K \hat{p}_{mk}^2$



Split #1: $x_1 \leq 5$

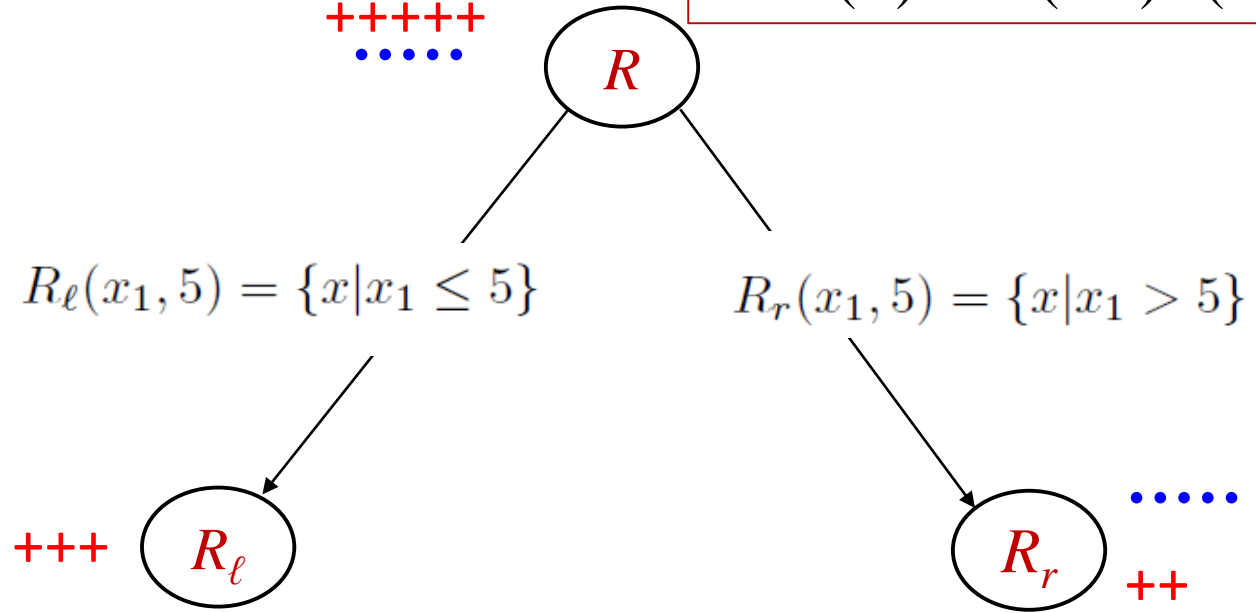
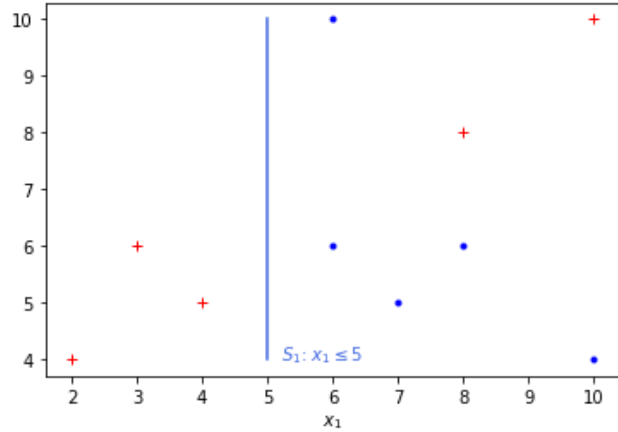
S_1



Classification – { +, • }

Gini Index: $1 - \sum_{k=1}^K \hat{p}_{mk}^2$

Split #1: $x_1 \leq 5$



$$E(R) = 1 - (0.5)^2 - (0.5)^2 = 0.5$$

$$E(R_\ell) = 1 - (1.0)^2 - (0.0)^2 = 0.0$$

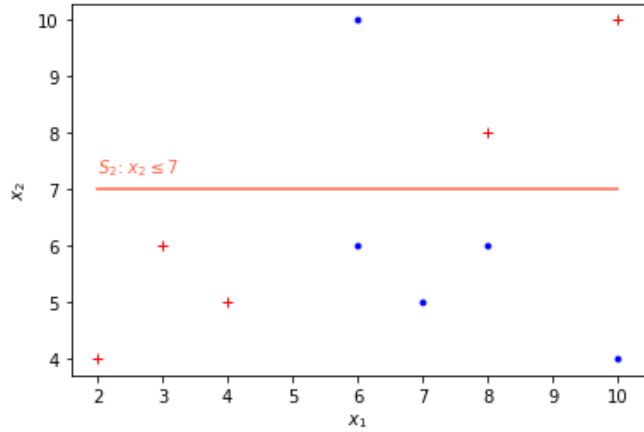
$$E(R_r) = 1 - (2/7)^2 - (5/7)^2 = 0.41$$

$$\begin{aligned} E(S_1) &= (3/10) \cdot E(R_\ell) + (7/10) \cdot E(R_r) \\ &= 0.3 \cdot 0.0 + 0.7 \cdot 0.41 \\ &= 0.29 \end{aligned}$$

$$\begin{aligned} G(S_1) &= E(R) - E(S_1) \\ &= 0.5 - 0.29 \\ &= \mathbf{0.21} \end{aligned}$$

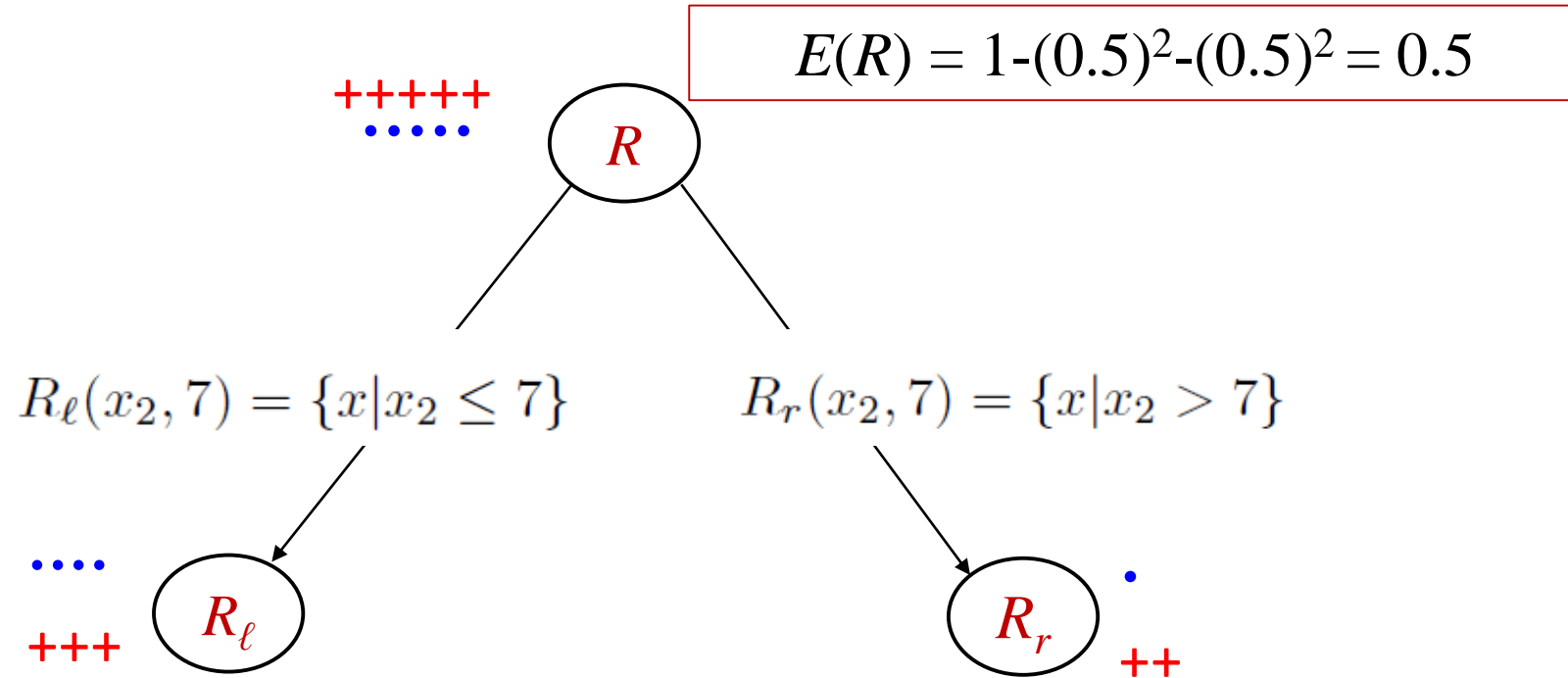
Classification – $\{+, \bullet\}$

Gini Index: $1 - \sum_{k=1}^K \hat{p}_{mk}^2$



Split #2: $x_2 \leq 7$

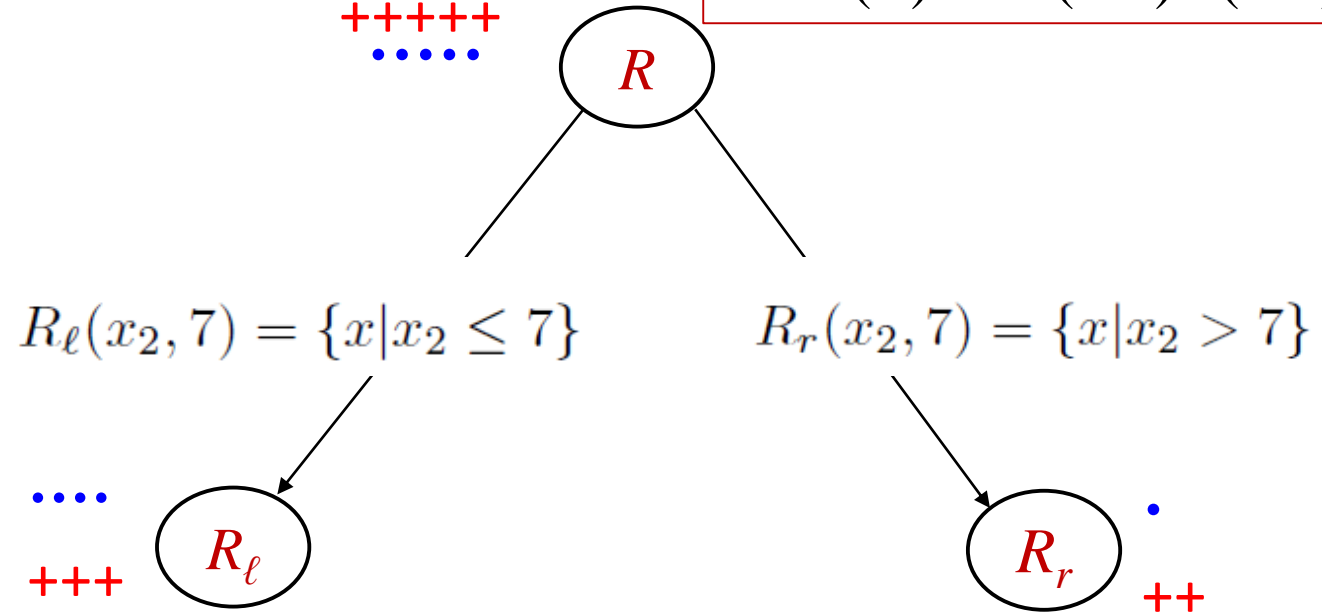
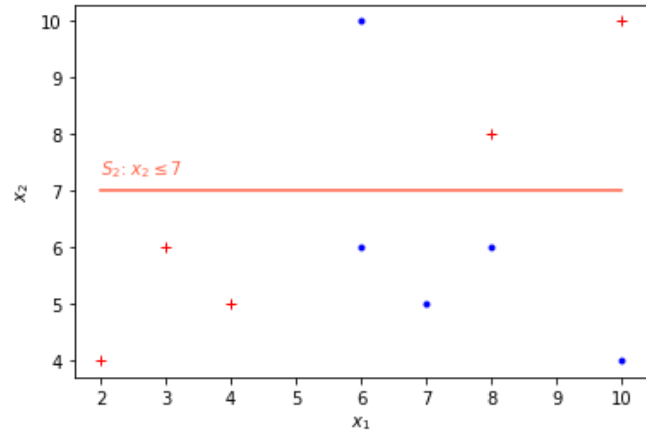
S_2



Classification – { +, • }

Gini Index: $1 - \sum_{k=1}^K \hat{p}_{mk}^2$

Split #2: $x_2 \leq 7$



$$E(R) = 1 - (0.5)^2 - (0.5)^2 = 0.5$$

$$R_\ell(x_2, 7) = \{x | x_2 \leq 7\}$$

$$R_r(x_2, 7) = \{x | x_2 > 7\}$$

$$E(R_\ell) = 1 - (4/7)^2 - (3/7)^2 = 0.49$$

$$E(R_r) = 1 - (1/3)^2 - (2/3)^2 = 0.44$$

$$\begin{aligned} E(S_2) &= (7/10) \cdot E(R_\ell) + (3/10) \cdot E(R_r) \\ &= 0.7 \cdot 0.49 + 0.3 \cdot 0.44 \\ &= 0.343 + 0.132 \\ &= 0.48 \end{aligned}$$

$$\begin{aligned} G(S_2) &= E(R) - E(S_2) \\ &= 0.5 - 0.48 \\ &= \mathbf{0.02} \end{aligned}$$

Optimization Models?

Generic Model

$$\begin{array}{ll}\text{minimize} & f(\boldsymbol{x}) \\ \text{subject to} & g(\boldsymbol{x}) \geq \mathbf{0}.\end{array}$$

Integer Programming Model

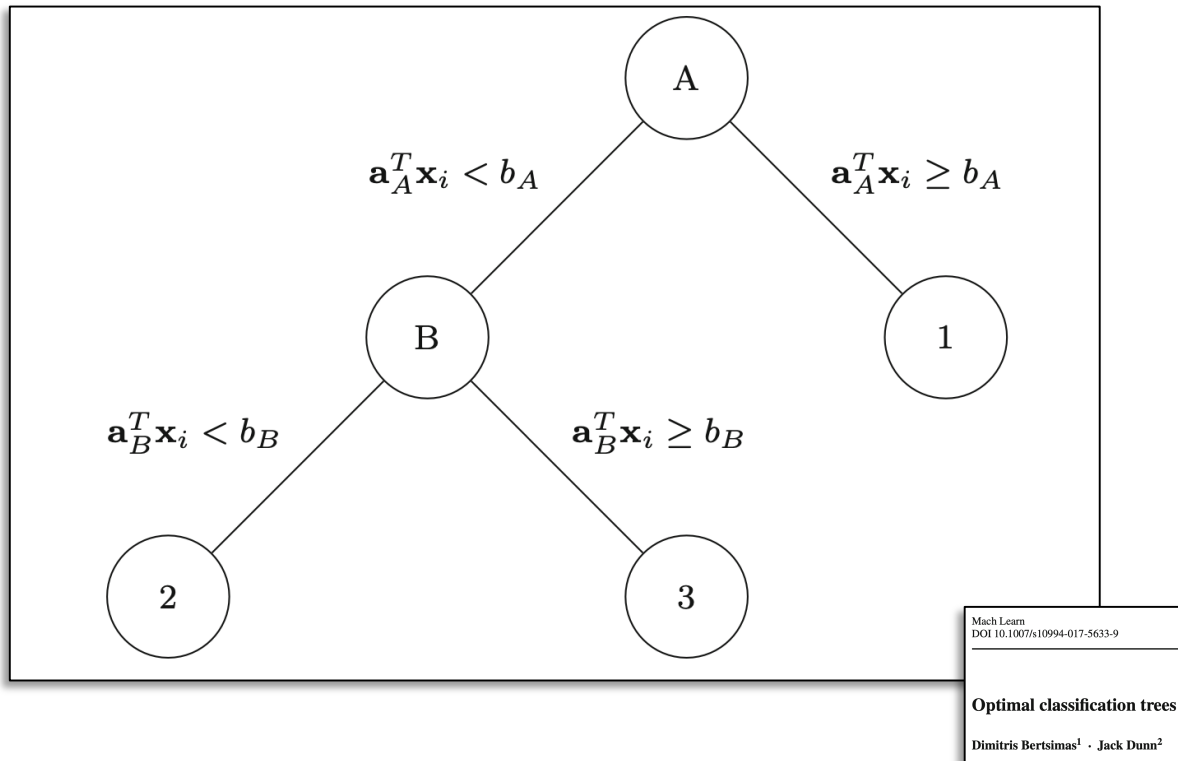
$$\begin{array}{ll}\text{minimize} & \boldsymbol{c}^\top \boldsymbol{x} \\ \text{subject to} & \boldsymbol{A}\boldsymbol{x} \geq \boldsymbol{b}, \\ & \boldsymbol{x} \in \{\mathbf{0}, \mathbf{1}\}.\end{array}$$

Linear Programming Model

$$\begin{array}{ll}\text{minimize} & \boldsymbol{c}^\top \boldsymbol{x} \\ \text{subject to} & \boldsymbol{A}\boldsymbol{x} \geq \boldsymbol{b}, \\ & \boldsymbol{x} \geq \mathbf{0}.\end{array}$$

Optimal Classification Trees (OCT)

$$\mathbf{x}_i \in \mathbb{R}^p, y_i \in \{1, \dots, K\}, i = 1, \dots, n$$



classification error

number of branch nodes
in T

minimize

$$R_{xy}(T) + \alpha |T|$$

tree

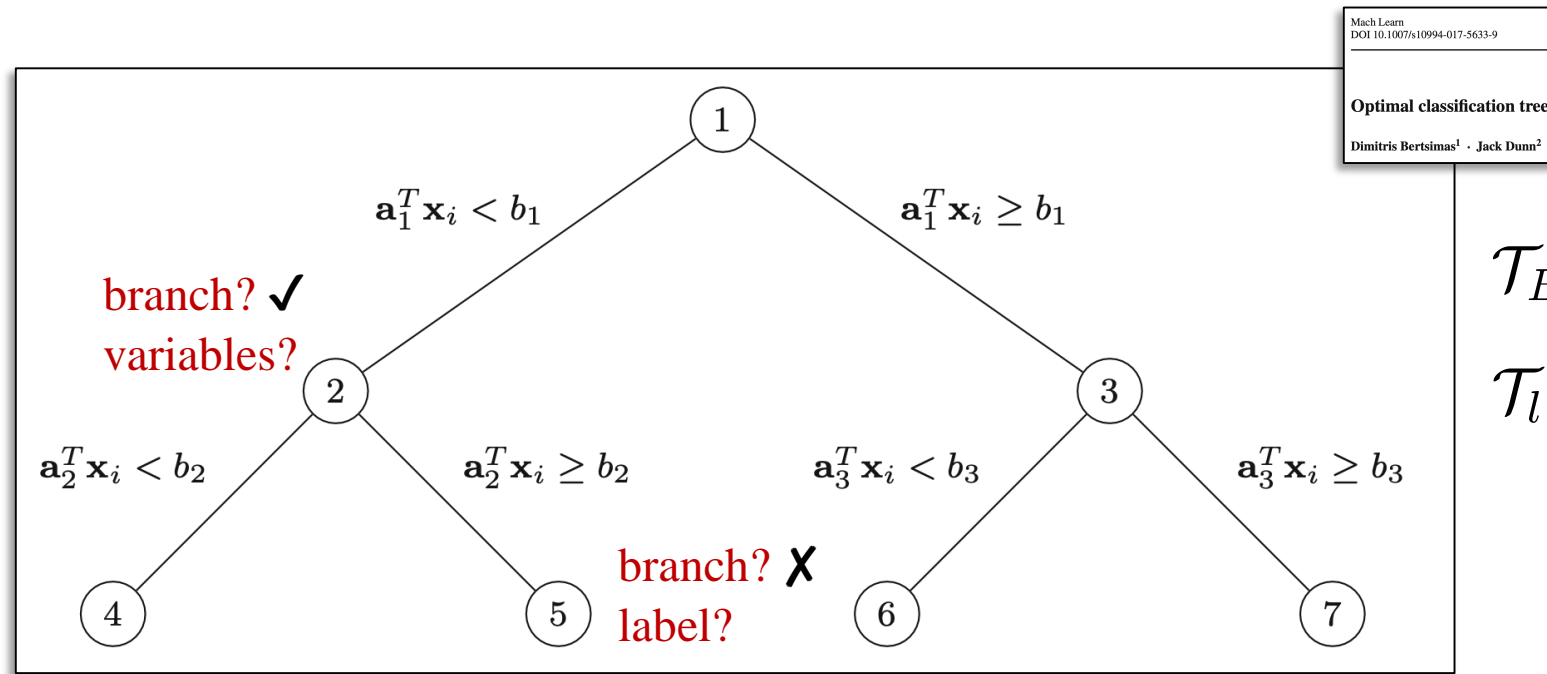
subject to

$$N_x(l) \geq N_{\min}, \quad l \in \text{leaves of } T$$

number points in leaf l

minimum number of
samples in any leaf

OCT - Setup



\mathcal{T}_B : branch nodes

\mathcal{T}_l : leaf nodes

$p(t)$: parent node of t

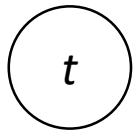
$A(t)$: set of ancestors of node t

$A_L(t)$: “left” ancestors of node t

$A_R(t)$: “right” ancestors of node t

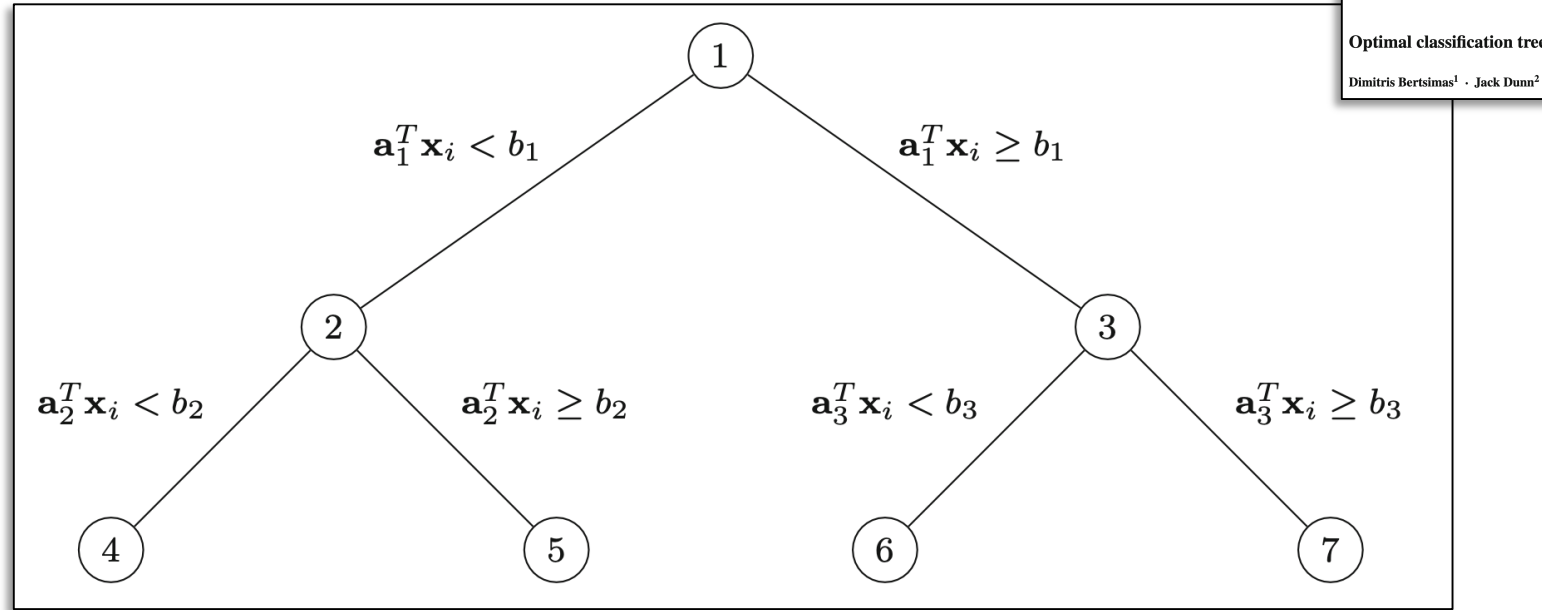
$$A_L(5) = \{1\}, \quad A_R(5) = \{2\}$$

$$A(t) = A_L(t) \cup A_R(t)$$



OCT – Splitting Constraints

Univariate Trees



\mathbf{a}_t : unit vector

$d_t \in \{0, 1\}$: split or not

$$\sum_{j=1}^p a_{jt} = d_t, \quad t \in \mathcal{T}_B$$

$$0 \leq b_t \leq d_t, \quad t \in \mathcal{T}_B$$

$$a_{jt} \in \{0, 1\}, \quad j = 1, \dots, p; \quad t \in \mathcal{T}_B$$

$$d_t \leq d_{p(t)}, \quad t \in \mathcal{T}_B \setminus \{1\}$$

OCT – Sample Tracking Constraints

Univariate Trees

z_{it} : \mathbf{x}_i is in node t or not

l_t : leaf t contains points or not

$$z_{it} \leq l_t, \quad t \in \mathcal{T}_L$$

$$\sum_{i=1}^n z_{it} \geq N_{\min} l_t, \quad t \in \mathcal{T}_L$$

$$\sum_{t \in \mathcal{T}_L} z_{it} = 1, \quad i = 1, \dots, n$$

$$\mathbf{a}_m^\top \mathbf{x}_i + \varepsilon \leq b_m + M_1(1 - z_{it}) \quad i = 1, \dots, n; t \in \mathcal{T}_L; m \in A_L(t)$$

$$\mathbf{a}_m^\top \mathbf{x}_i \geq b_m - M_2(1 - z_{it}) \quad i = 1, \dots, n; t \in \mathcal{T}_L; m \in A_R(t)$$

OCT – Objective Function

Univariate Trees

$$Y_{ik} = \begin{cases} +1, & \text{if } y_i = k; \\ -1, & \text{otherwise} \end{cases}$$

$$N_{kt} = \frac{1}{2} \sum_{i=1}^n (1 + Y_{ik}) z_{it}, \quad k = 1, \dots, K; \quad t \in \mathcal{T}_L$$

$$c_t = \arg \max_{k=1, \dots, K} \{N_{kt}\} \quad \text{and} \quad c_{kt} = \begin{cases} 1, & \text{if } c_t = k; \\ 0, & \text{otherwise} \end{cases}$$

$$N_t = \sum_{i=1}^n z_{it}, \quad t \in \mathcal{T}_L$$

$$L_t = N_t - \max_{k=1, \dots, K} \{N_{kt}\} = \min_{k=1, \dots, K} \{N_t - N_{kt}\}$$

$$L_t \geq N_t - N_{kt} - M(1 - c_{kt}), \quad k = 1, \dots, K; t \in \mathcal{T}_L,$$

$$L_t \leq N_t - N_{kt} + M c_{kt}, \quad k = 1, \dots, K; t \in \mathcal{T}_L,$$

$$L_t \geq 0, \quad t \in \mathcal{T}_L$$

$$(M = n)$$

OCT – Overall Model

Univariate Trees

$$\begin{aligned}
 &\text{minimize} && \frac{1}{\hat{L}} \sum_{t \in \mathcal{T}_L} L_t + \alpha \sum_{t \in \mathcal{T}_B} d_t \\
 &\text{subject to} && L_t \geq N_t - N_{kt} - M(1 - c_{kt}), \quad k = 1, \dots, K; t \in \mathcal{T}_L, \\
 & && L_t \leq N_t - N_{kt} + M c_{kt}, \quad k = 1, \dots, K; t \in \mathcal{T}_L, \\
 & && L_t \geq 0, \quad t \in \mathcal{T}_L \\
 & && N_{kt} = \frac{1}{2} \sum_{i=1}^n (1 + Y_{ik}) z_{it}, \quad k = 1, \dots, K; t \in \mathcal{T}_L \\
 & && N_t = \sum_{i=1}^n z_{it}, \quad t \in \mathcal{T}_L \\
 & && \mathbf{a}_m^\top \mathbf{x}_i + \varepsilon \leq b_m + M_1(1 - z_{it}) \quad i = 1, \dots, n; t \in \mathcal{T}_B; m \in A_L(t) \\
 & && \mathbf{a}_m^\top \mathbf{x}_i \geq b_m - M_2(1 - z_{it}) \quad i = 1, \dots, n; t \in \mathcal{T}_B; m \in A_R(t) \\
 & && z_{it} \leq l_t, \quad t \in \mathcal{T}_L \\
 & && \sum_{i=1}^n z_{it} \geq N_{\min} l_t, \quad t \in \mathcal{T}_L \\
 & && \sum_{i=1}^n z_{it} = 1, \quad i = 1, \dots, n \\
 & && \sum_{j=1}^p a_{jt} = d_t, \quad t \in \mathcal{T}_B \\
 & && 0 \leq b_t \leq d_t, \quad t \in \mathcal{T}_B \\
 & && d_t \leq d_{p(t)}, \quad t \in \mathcal{T}_B \setminus \{1\} \\
 & && z_{it}, l_t \in \{0, 1\}, i = 1, \dots, n; t \in \mathcal{T}_L, \\
 & && a_{jt}, d_t \in \{0, 1\}, i = 1, \dots, n; t \in \mathcal{T}_B
 \end{aligned}$$

$$\begin{aligned}
 &\text{minimize} && R_{xy}(T) + \alpha|T| \\
 &\text{subject to} && N_x(l) \geq N_{\min}, \quad l \in \text{leaves of } T
 \end{aligned}$$

Number of Binary Variables

$$\mathcal{O}(n2^D)$$

tree depth

OCT – Take Away Messages

- A large-scale mixed integer optimization model
- Relies heavily on the speed of state-of-the-art solvers (*e.g.*, Gurobi)
- Can be easily extended to splits with hyperplanes (**multivariate trees**)
- Empirically shown that overfitting does not occur
- The solution time increases drastically with the number of samples and the tree depth
- Random Forest still performs better when it comes to accuracy

Mach Learn
DOI 10.1007/s10994-017-5633-9

Optimal classification trees

Dimitris Bertsimas¹ · Jack Dunn²

State-of-the-art

Mach Learn
DOI 10.1007/s10994-017-5633-9

Optimal classification trees

Dimitris Bertsimas¹ · Jack Dunn²

$n = 5456$; $p = 2$; $K = 4$; (2017)

Strong Optimal Classification Trees

Sina Aghaei

Center for Artificial Intelligence in Society, University of Southern California, Los Angeles, CA 90089, USA, saghaei@usc.edu

Andrés Gómez

Department of Industrial and Systems Engineering, Los Angeles, CA 90089, USA, gomezand@usc.edu

Phebe Vayanos

Center for Artificial Intelligence in Society, University of Southern California, Los Angeles, CA 90089, USA,
phebe.vayanos@usc.edu

$n = 3196$; $p = 38$; $K = 2$; (2021)

The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-19)

Learning Optimal Classification Trees Using a Binary Linear Program Formulation

Sicco Verwer

Yingqian Zhang

$n = 4601$; $p = 57$; $K = 2$; (2019)

arXiv version 1 (2020) 1-48

Submitted 4/00; Published 10/00

MurTree: Optimal Classification Trees via Dynamic Programming and Search

Emir Demirović

*Delft University of Technology
Delft, The Netherlands*

E.DEMIROVIC@TUDELFT.NL

Anna Lukina

*Institute of Science and Technology Austria
Klosterneuburg, Austria*

ANNA.LUKINA@IST.AC.AT

Emmanuel Hebrard

*LAAS CNRS
Toulouse, France*

HEBRARD@LAAS.FR

Jeffrey Chan

*RMIT University
Melbourne, Australia*

JEFFREY.CHAN@RMIT.EDU.AU

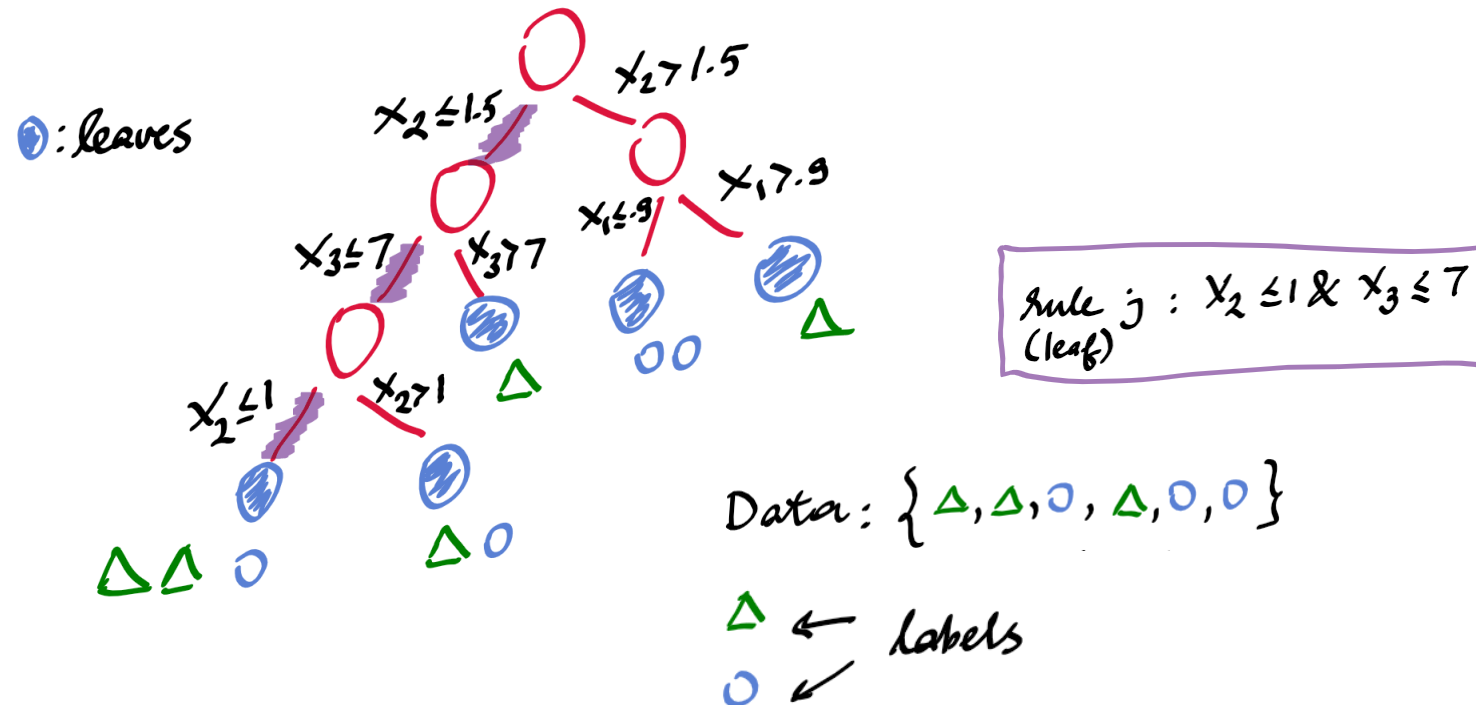
$n = 43500$; $p = 181$; $K = 7$; (2020)

Rules and Rule Set Generation?

RULE-BASED METHODS USING OPTIMIZATION

IF AGE > 25 AND PASTCREDIT=TRUE THEN ACCEPT

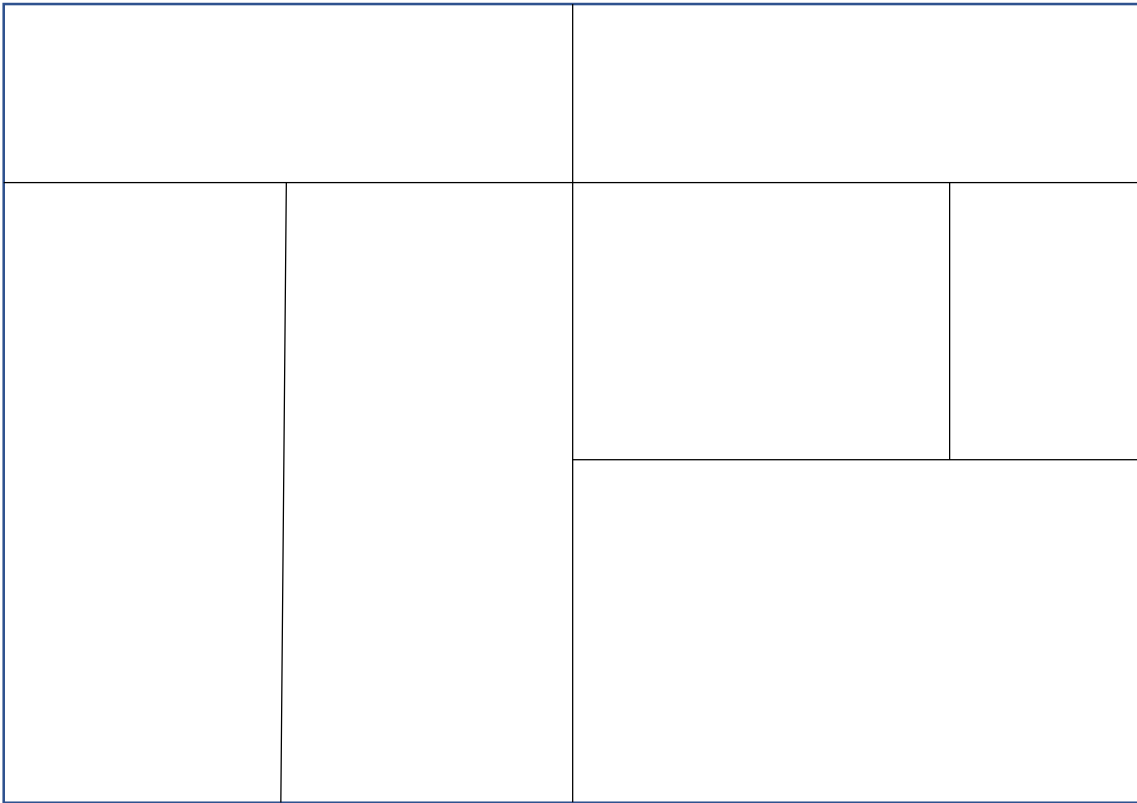
DECISION
TREES



DECISION TREES

vs

RULES



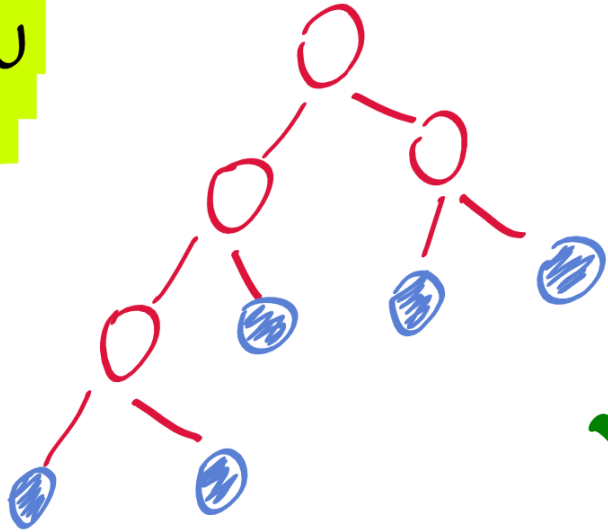
Rules are simple and easy to understand. In some cases they are easier to understand than the decision trees.

What are possible measures for interpretability?

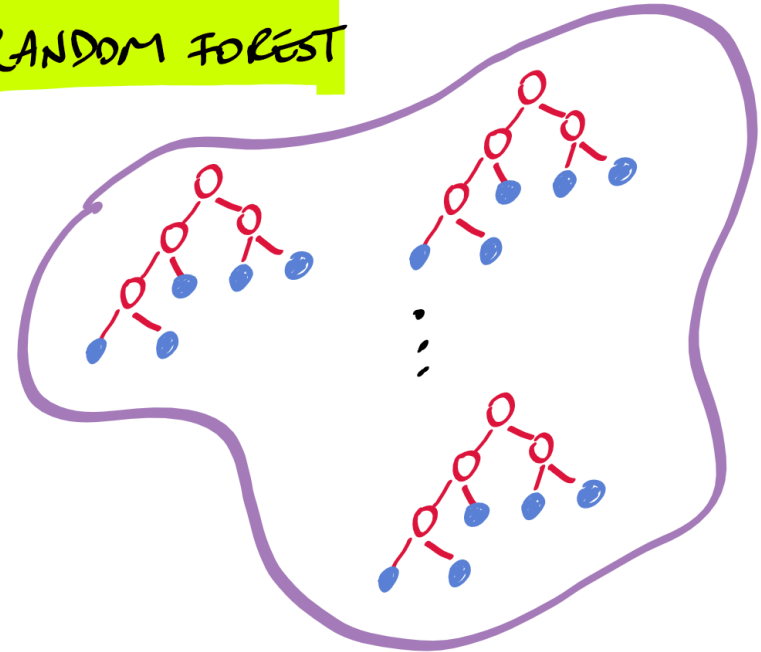
- Set Size (number of rules)
- Rule length (number of clauses)
- Cover (number of samples covered)
- Overlap (between two rules)

Rule Generation

DECISION
TREES



RANDOM FOREST



Variance ↓ Accuracy ↑

Interpretability ↓

**Rule Generation for Classification:
Scalability, Interpretability, and Fairness**

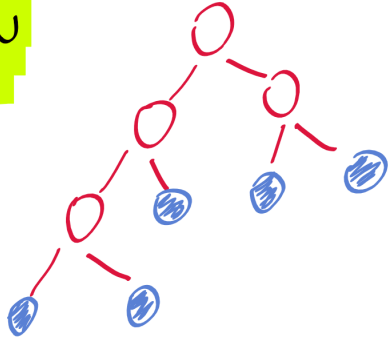
Adia Lumadjeng*, Tabea Röber**, M. Hakan Akyüz*, Ş. İlker Birbil**

*Erasmus University Rotterdam, 3000 DR, Rotterdam, The Netherlands

**University of Amsterdam, 11018 TV, Amsterdam P.O. Box 15953, The Netherlands

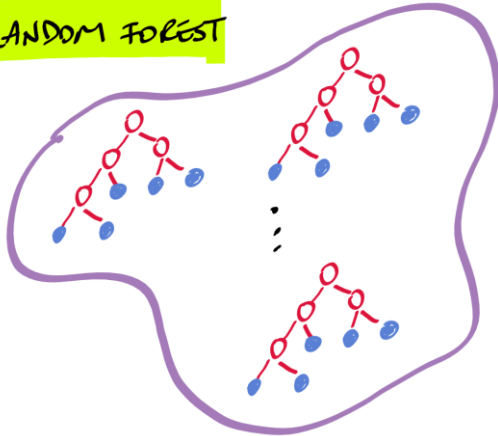
Scalability

DECISION
TREES



Rules
Rules
Rules
...

RANDOM FOREST



BOOSTING



Optimization
with
Constraints

~~minimize $f(x)$
subject to $g(x) \geq 0$.~~

minimize $c^T x$
subject to $Ax \geq b$,
 ~~$x \in \{0, 1\}$.~~

minimize $c^T x$
subject to $Ax \geq b$,
 $x \geq 0$.

LP!

Scalable Rule Generation for Learning

Modeling

$$\mathcal{D} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}\}$$

$$\mathbf{x}_i \in \mathbb{R}^p, y_i \in \{1, \dots, K\}$$

$$y_i = k \implies \mathbf{y}_i = \left(-\frac{1}{K-1}, -\frac{1}{K-1}, \dots, \overset{k^{\text{th}}}{1}, \dots, -\frac{1}{K-1}\right)^\top \in \mathbb{R}^K$$

prediction of rule j
for sample i $\longrightarrow \mathbf{R}_j(\mathbf{x}_i) \in \mathbb{R}^K$

prediction
for sample i $\longrightarrow \hat{\mathbf{y}}_i(\mathbf{w}) = \sum_{j \in J} \boxed{a_{ij}} \mathbf{R}_j(\mathbf{x}_i) \boxed{w_j}$

classification loss
for sample i $\longrightarrow \mathcal{L}(\hat{\mathbf{y}}_i(\mathbf{w}), \mathbf{y}_i) = \max\{1 - \boxed{\kappa} \hat{\mathbf{y}}_i(\mathbf{w})^\top \mathbf{y}_i, 0\}$

Total Loss $\longrightarrow \sum_{i \in \mathcal{I}} \underbrace{\mathcal{L}(\hat{\mathbf{y}}_i(\mathbf{w}), \mathbf{y}_i)}_{\boxed{v_i}}$

set of rules

rule j covers sample i or not

weight of rule j

$$\kappa = (K - 1)/K$$

auxiliary variables

Master Problem

$$v_i \geq \max\{1 - \kappa \hat{\mathbf{y}}_i(\mathbf{w})^\top \mathbf{y}_i, 0\}$$

$$\begin{array}{ll} \text{minimize} & \underbrace{\lambda \sum_{j \in \mathcal{J}} c_j w_j}_{(i)} + \underbrace{\sum_{i \in \mathcal{I}} v_i}_{(ii)} \\ \text{subject to} & \sum_{j \in \mathcal{J}} \hat{a}_{ij} w_j + v_i \geq 1, \quad i \in \mathcal{I}, \\ & v_i \geq 0, \quad i \in \mathcal{I}, \\ & w_j \geq 0, \quad j \in \mathcal{J} \end{array}$$

$$\hat{a}_{ij} = \kappa a_{ij} \mathbf{R}_j(\mathbf{x}_i)^\top \mathbf{y}_i$$

(i) total **weighted** cost of using rules

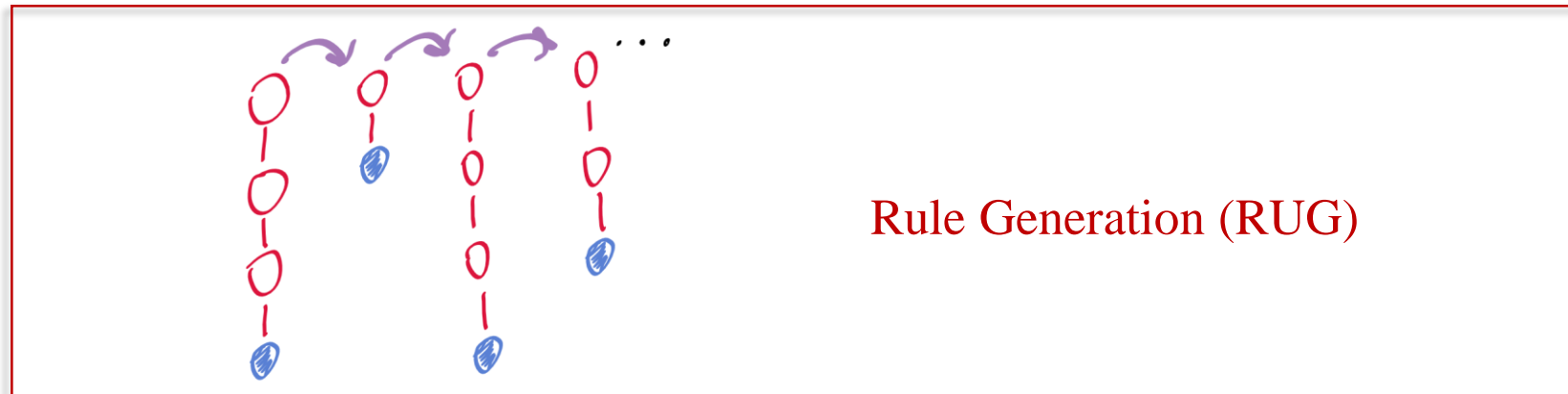
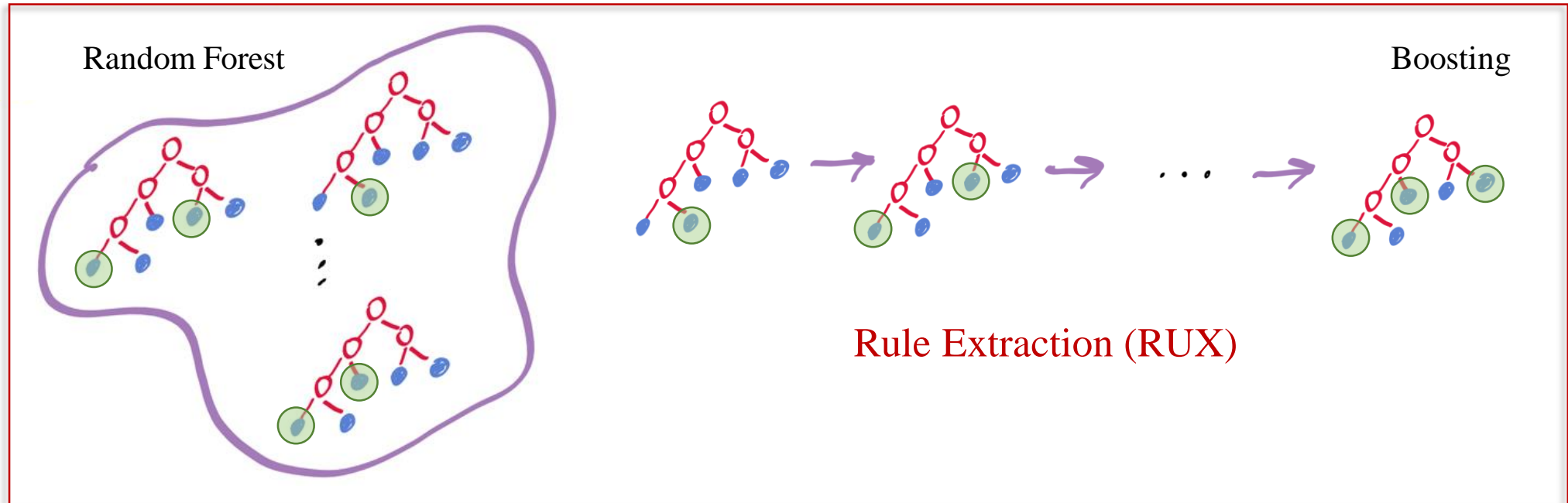
(ii) total loss from misprediction

scaling hyperparameter

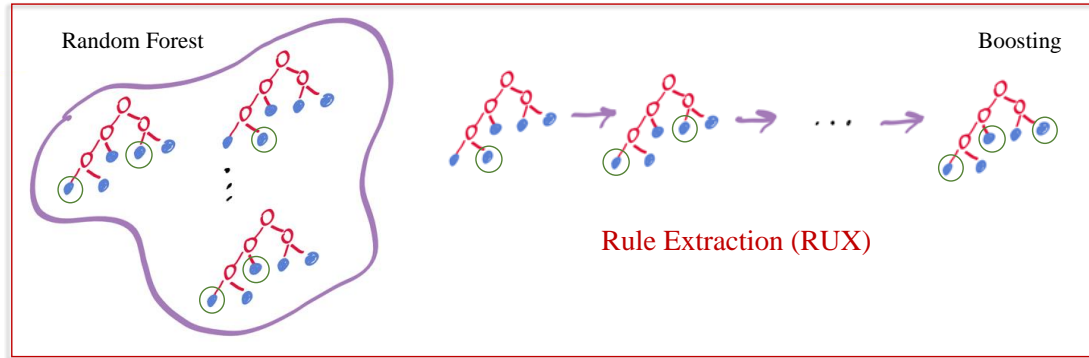
Set of rules (aka columns)?

Rule Discovery

Set of rules (aka columns)?



RUX



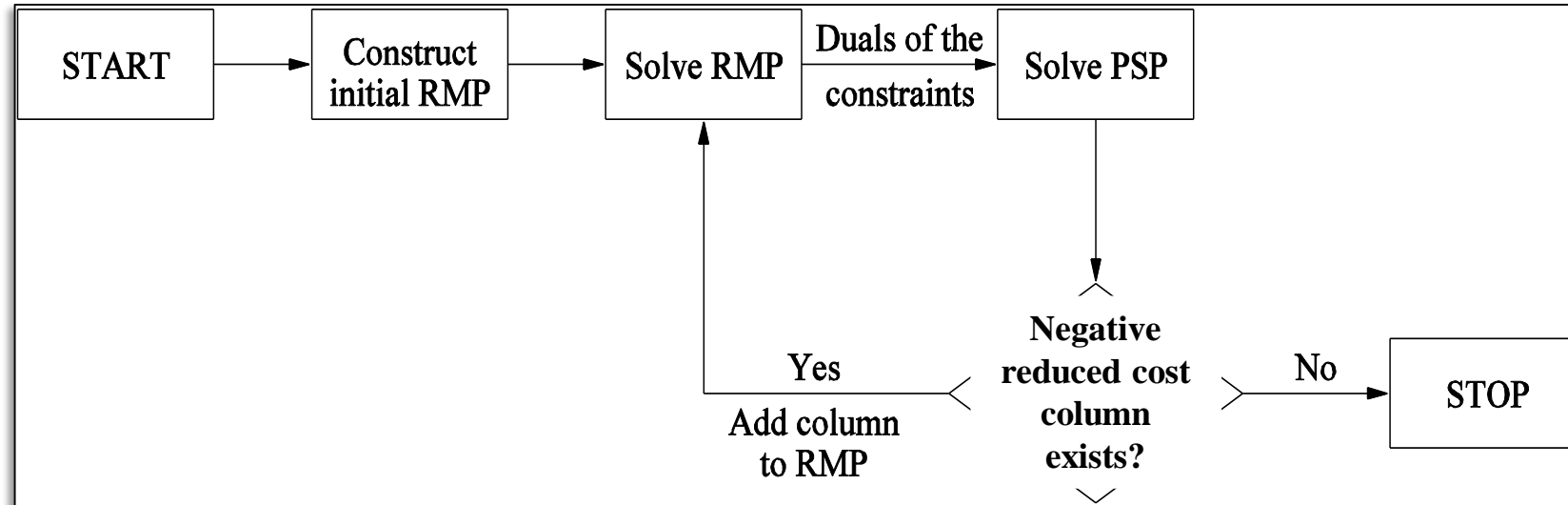
Master Problem

$$\begin{aligned} &\text{minimize} && \lambda \sum_{j \in \mathcal{J}} c_j w_j + \sum_{i \in \mathcal{I}} v_i \\ &\text{subject to} && \sum_{j \in \mathcal{J}} \hat{a}_{ij} w_j + v_i \geq 1, && i \in \mathcal{I}, \\ & && v_i \geq 0, && i \in \mathcal{I}, \\ & && w_j \geq 0, && j \in \mathcal{J} \end{aligned}$$

Algorithm - Rule Extraction

- 1: **Input:** training data, $\mathcal{D} = \{(\mathbf{x}_i, y_i) : i \in \mathcal{I}\}$
 - 2: Train a tree ensemble to construct rule pool \mathcal{J}
 - 3: $\mathcal{J}^* \leftarrow$ Solve master problem using \mathcal{J}
 - 4: **return** \mathcal{J}^* and $w_j, j \in \mathcal{J}^*$
-

Column Generation



RMP*

$$\begin{array}{ll} \text{minimize} & \lambda \sum_{j \in \mathcal{J}_t} c_j w_j + \sum_{i \in \mathcal{I}} v_i \\ \text{subject to} & \sum_{j \in \mathcal{J}_t} \hat{a}_{ij} w_j + v_i \geq 1, \quad i \in \mathcal{I}, \\ & v_i \geq 0, \quad i \in \mathcal{I}, \\ & w_j \geq 0, \quad j \in \mathcal{J}_t \end{array}$$

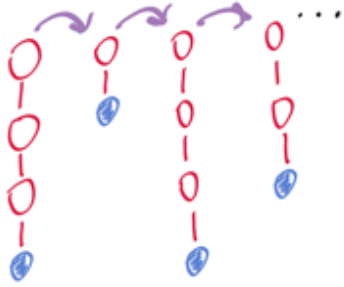
* restricted master problem

RUG

Dual Problem (J_t)

$$\begin{aligned} &\text{maximize} && \sum_{i \in \mathcal{I}} \beta_i \\ &\text{subject to} && \sum_{i \in \mathcal{I}} \hat{a}_{ij} \beta_i \leq \lambda c_j, \quad j \in \mathcal{J}_t, \\ &&& 0 \leq \beta_i \leq 1, \quad i \in \mathcal{I} \end{aligned}$$

Rule Generation (RUG)



PSP*

$$\min_{j \in \mathcal{J} / \mathcal{J}_t} \left\{ \lambda c_j - \sum_{i \in \mathcal{I}} \hat{a}_{ij} \beta_i^{(t)} \right\}$$

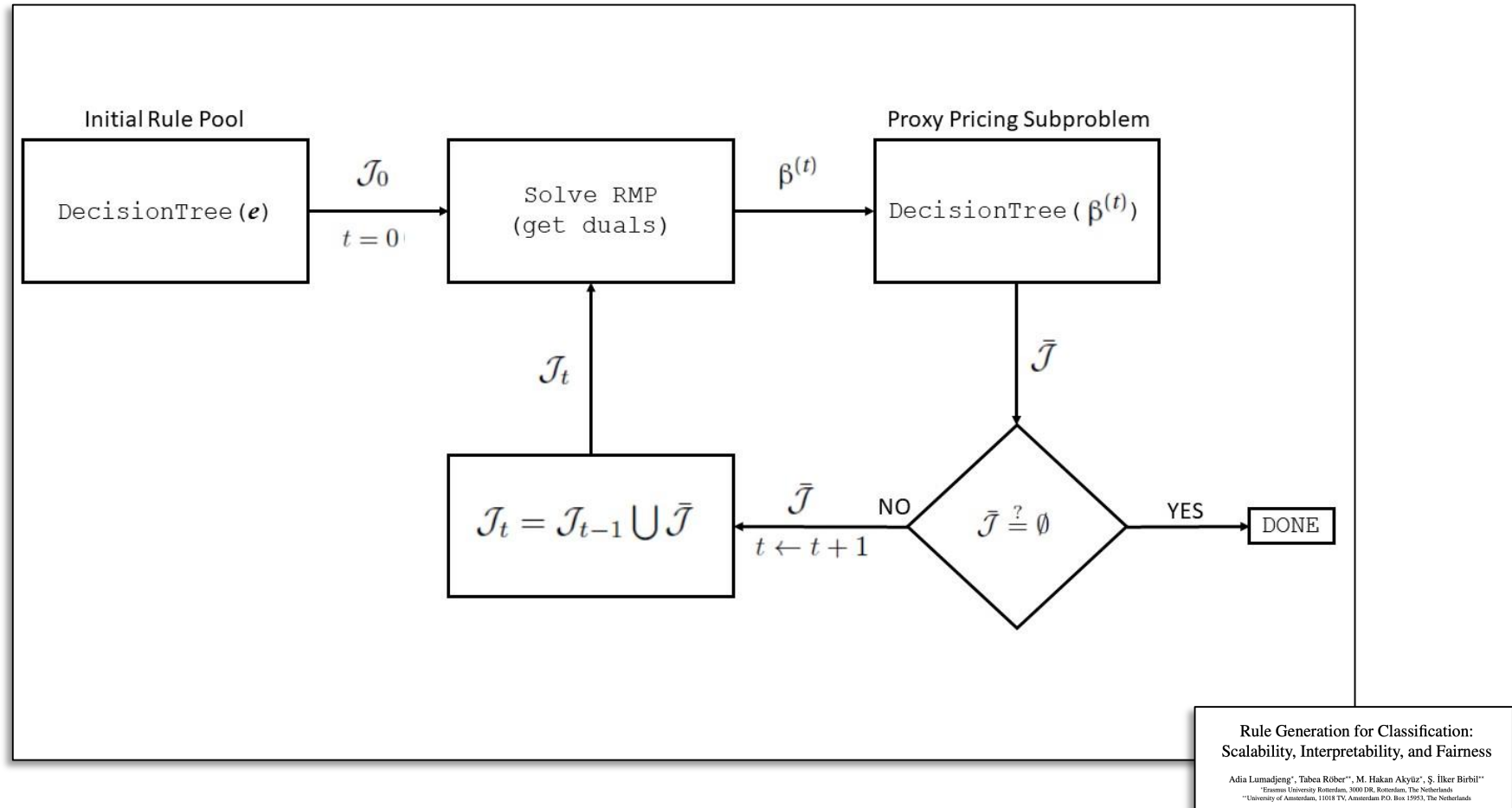
Proposition. *PSP is NP-hard*

* pricing subproblem

Proxy PSP

PSP

$$\min_{j \in \mathcal{J} / \mathcal{J}_t} \left\{ \lambda c_j - \sum_{i \in \mathcal{I}} \hat{a}_{ij} \beta_i^{(t)} \right\}$$



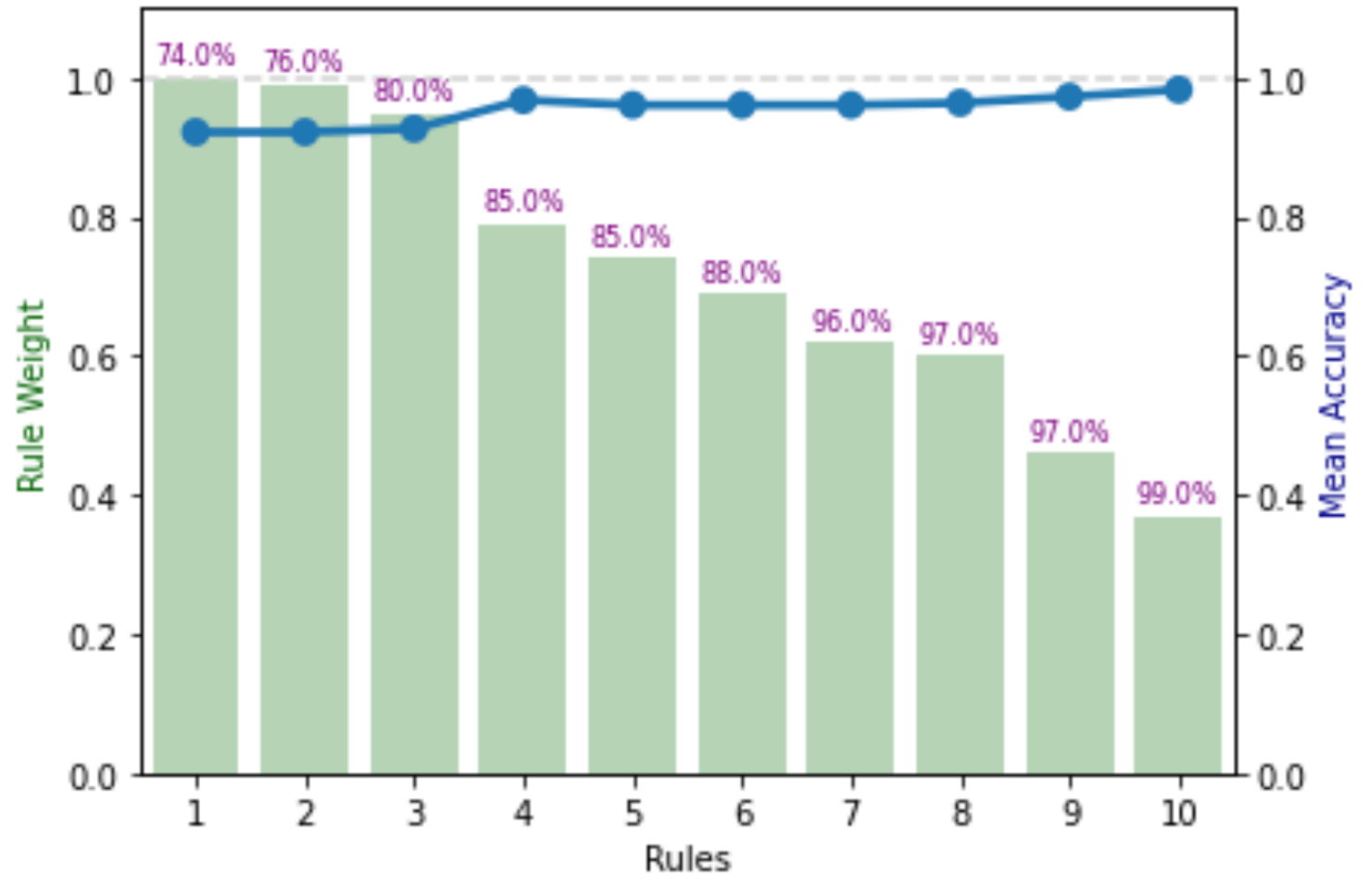
“Banknote” Dataset

Accuracy

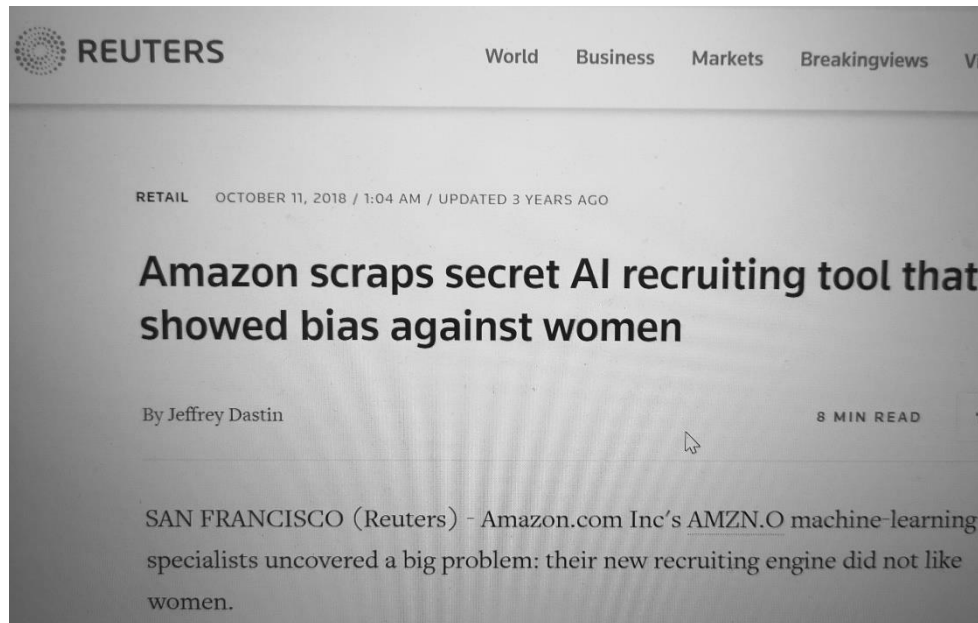
Random Forest: 0.9660194174757282
AdaBoost: 0.9975728155339806
RUXRF: 0.9878640776699029
RUXADA: 0.9975728155339806
RUG: 1.0

Number of Rules

Random Forest: 712
AdaBoost: 776
RUXRF: 32
RUXADA: 39
RUG: 40



Extensions: Fairness



"Everyone wanted this holy grail," one of the people said. "They literally wanted it to be an engine where I'm going to give you 100 resumes, it will spit out the top five, and we'll hire those."



<https://thenextweb.com/news/ai-models-need-to-be-interpretable-rather-than-just-explainable>

Fairness

- Sensitive Attribute (female/male, race, age etc., G)
- Additional constraints are needed.
- There are more than 20 fairness definitions (Verma and Rubin, 2018).

- Group fairness (Statistical parity):

$$P(\hat{Y}=1 | G = m) = P(\hat{Y} = 1 | G = f).$$

- Predictive parity:

$$P(Y = 1 | \hat{Y} = 1, G = m) = P(Y = 1 | \hat{Y} = 1, G = f).$$

- Equal opportunity (False negative rate):

$$P(\hat{Y} = 0 | Y = 1, G = m) = P(\hat{Y} = 0 | Y = 1, G = f).$$

- Equalized Odds (Disparate Mistreatment):

$$P(\hat{Y} = 1 | Y = i, G = m) = P(\hat{Y} = 1 | Y = i, G = f), i \in 0, 1.$$

- Overall accuracy equality:

$$P(\hat{Y} = Y, G = m) = P(\hat{Y} = Y, G = f).$$

Binary
classification

What happens when there are multiple classes?

Equalized Odds:

$$P(\hat{y}_i \neq y_i | y_i = k, G = g) = P(\hat{y}_j \neq y_j | y_j = k, G = g'), \quad \forall i, j \in \mathcal{I}, \forall k \in \mathcal{K} \text{ and } \forall g, g' \in \mathcal{G}$$

$$\begin{aligned} & \text{minimize} && \lambda \sum_{j \in \mathcal{J}} c_j w_j + \sum_{i \in \mathcal{I}} v_i \\ & \text{subject to} && \sum_{j \in \mathcal{J}} \hat{a}_{ij} w_j + v_i \geq 1, \quad i \in \mathcal{I}, \\ & && v_i \geq 0, \quad i \in \mathcal{I}, \\ & && w_j \geq 0, \quad j \in \mathcal{J} \end{aligned}$$

Fairness
Constraints

$$\left\{ \begin{aligned} \sum_{i \in \mathcal{P}_{k,1}} v_i - \sum_{i \in \mathcal{P}_{k,2}} v_i &\leq \epsilon && \forall k \in \mathcal{K} \text{ and } \forall g, g' \in \mathcal{G} \\ \sum_{i \in \mathcal{P}_{k,2}} v_i - \sum_{i \in \mathcal{P}_{k,1}} v_i &\leq \epsilon && \forall k \in \mathcal{K} \text{ and } \forall g, g' \in \mathcal{G} \end{aligned} \right.$$

FairRUX

FairRUG

$$\mathcal{P}_{k,g} = \{i \in \mathcal{I} : y_i = k, G = g\}$$

Equal Overall Mistreatment:

$$\sum_{k \in \mathcal{K}} P(\hat{y}_i \neq y_i = k | G = g) = \sum_{k \in \mathcal{K}} P(\hat{y}_i \neq y_i = k | G = g'), \quad \forall g, g' \in \mathcal{G}$$

$$\begin{aligned} & \text{minimize} && \lambda \sum_{j \in \mathcal{J}} c_j w_j + \sum_{i \in \mathcal{I}} v_i \\ & \text{subject to} && \sum_{j \in \mathcal{J}} \hat{a}_{ij} w_j + v_i \geq 1, \quad i \in \mathcal{I}, \\ & && v_i \geq 0, \quad i \in \mathcal{I}, \\ & && w_j \geq 0, \quad j \in \mathcal{J} \end{aligned}$$

Fairness
Constraints

$$\left\{ \begin{aligned} \sum_{i \in \mathcal{I}_g} v_i - \sum_{i \in \mathcal{I}_{g'}} v_i &\leq \epsilon && \forall g, g' \in \mathcal{G} \\ \sum_{i \in \mathcal{I}_{g'}} v_i - \sum_{i \in \mathcal{I}_g} v_i &\leq \epsilon && \forall g, g' \in \mathcal{G} \end{aligned} \right.$$

FairRUX

FairRUG

$$\mathcal{I}_g = \{i \in \mathcal{I} : G = g\}$$

Rule Generation – Take Aways

- Linear programming brings in scalability
- The generated rules come with their associated weights
- Proxy PSP can be solved very fast
- Fairness constraints can be easily added
- There could still be too many rules leading to a less interpretable model
- Accuracy can be inferior compared against ensemble methods

Rule Generation for Classification: Scalability, Interpretability, and Fairness

Adia Lumadjeng*, Tabea Röber**, M. Hakan Akyüz*, Ş. İlker Birbil**

*Erasmus University Rotterdam, 3000 DR, Rotterdam, The Netherlands

**University of Amsterdam, 11018 TV, Amsterdam P.O. Box 15953, The Netherlands

$n = 245057$; $p = 3$; $K = 2$;

$n = 58509$; $p = 48$; $K = 11$; (2022)

Optional Readings

Interpretable and Fair Boolean Rule Sets via Column Generation

Connor Lawless¹ Sanjeeb Dash² Oktay Günlük¹ Dennis Wei²

Mach Learn

DOI 10.1007/s10994-017-5633-9

Optimal classification trees

Dimitris Bertsimas¹ · Jack Dunn²

Rule Generation for Classification: Scalability, Interpretability, and Fairness

Adia Lumadjeng*, Tabea Röber**, M. Hakan Akyüz*, Ş. İlker Birbil**

*Erasmus University Rotterdam, 3000 DR, Rotterdam, The Netherlands

**University of Amsterdam, 11018 TV, Amsterdam P.O. Box 15953, The Netherlands