

Trustworthy AI for Business and Society

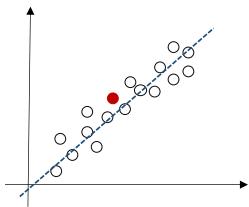
Ilker Birbil

Glass Box Methods



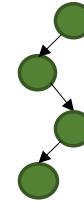
Why Glass Box?

“Additive” Models



$$\hat{y} = 1.0 + 2.0x_1 - 3.0x_2$$

Rules

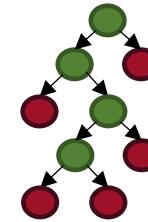


R₁: If Marginal Adhesion ≥ 8 then Malignant

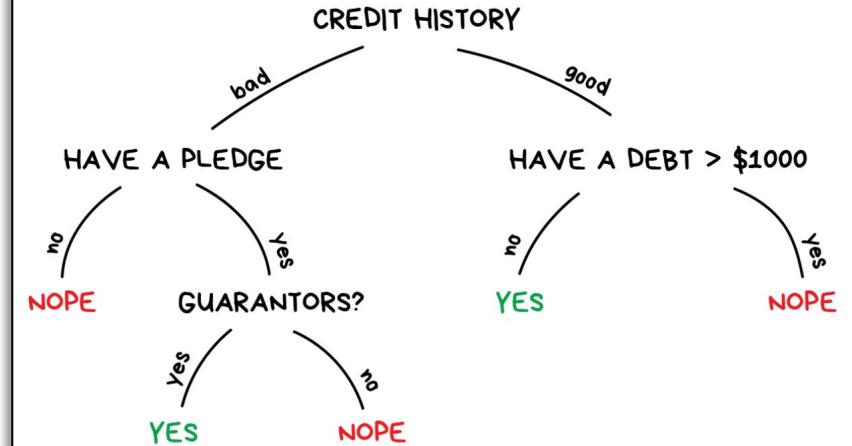
R₂: If Bare Nuclei ≤ 2 and Uniformity of Cell Shape ≤ 3 then Benign

R₃: If Marginal Adhesion ≥ 2 and Uniformity of Cell Size ≥ 5 then Malignant

Trees



GIVE A LOAN?

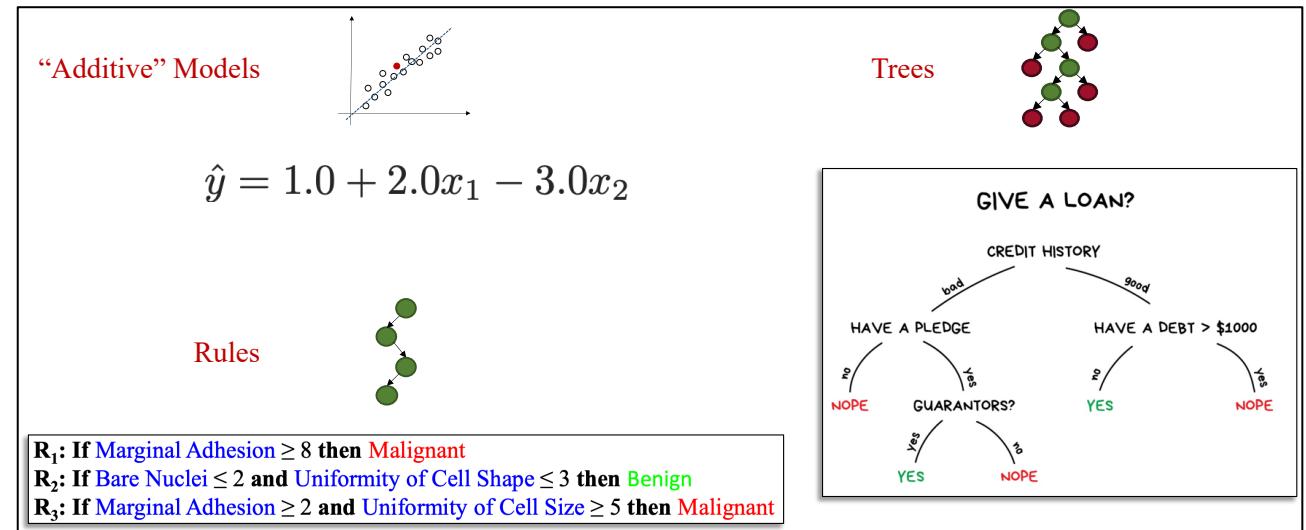


([link](#))



Big Picture

- Linear Models
 - Generalized Additive Models
 - Explainable Boosting Machines
 - Decision Trees
 - Rule-based Approaches



Linear Models

$$\{(x_i, y_i) : i = 1, \dots, n\}$$

$$x_i \in \mathbb{R}^p, y_i \in \mathbb{R}$$

$$Y \approx \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$$y_i \approx \underbrace{\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \dots + \hat{\beta}_p x_{ip}}_{\hat{y}_i}, \quad i = 1, \dots, n$$

prediction

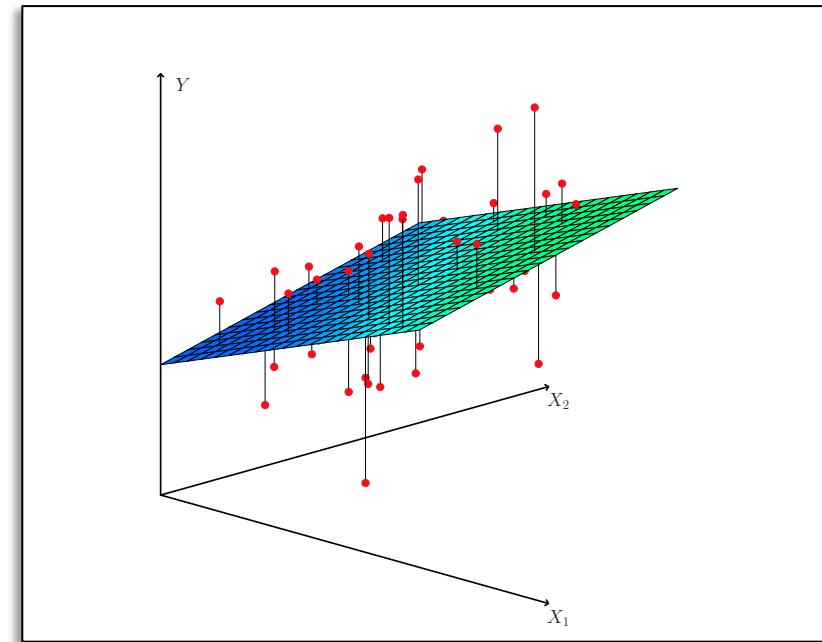
Least
Squares
Method

$$\beta = (\beta_1, \dots, \beta_p)^\top$$

$$\min_{\beta_0, \beta} \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - x_i^\top \beta)^2$$



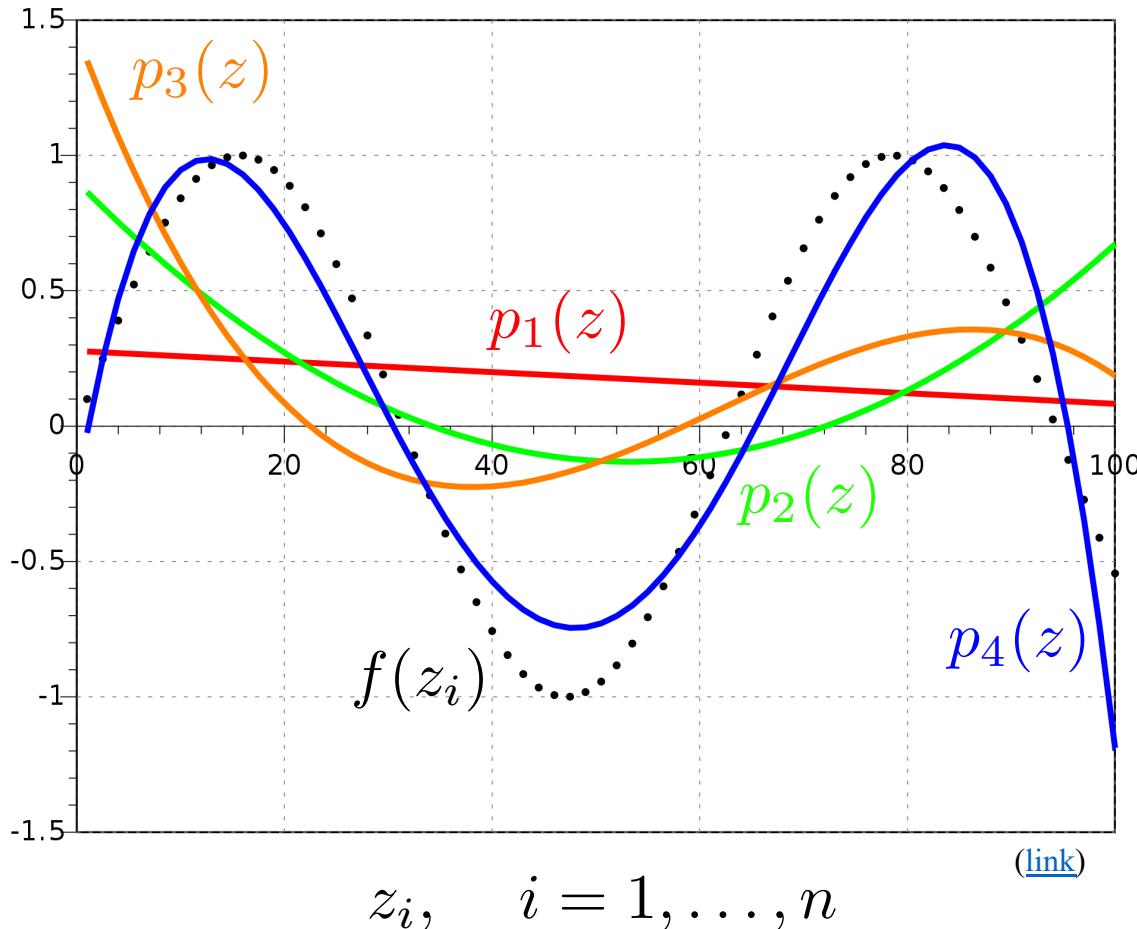
$$\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$$



([link](#))



Linear Models



$$Y \approx \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p$$

$$p_4(z_i) = \hat{\beta}_0 + \hat{\beta}_1 z_i + \hat{\beta}_2 z_i^2 + \hat{\beta}_3 z_i^3 + \hat{\beta}_4 z_i^4$$

$$x_i = (z_i, z_i^2, z_i^3, z_i^4)^\top \in \mathbb{R}^4, \quad i = 1, \dots, n$$

$$\hat{y}_i = p_4(z_i)$$

Training Data $\{(x_i, \underbrace{f(z_i)}_{y_i}) : i = 1, \dots, n\}$



Linear Regression and Regularization

Linear

$$\min_{\beta_0, \beta} \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - x_i^\top \beta)^2$$

Ridge

$\lambda \geq 0$

$$\min_{\beta_0, \beta} \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - x_i^\top \beta)^2 + \frac{\lambda}{2} \|\beta\|_2^2$$

Lasso

$\lambda \geq 0$

$$\min_{\beta_0, \beta} \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - x_i^\top \beta)^2 + \lambda \|\beta\|_1$$

Elastic net

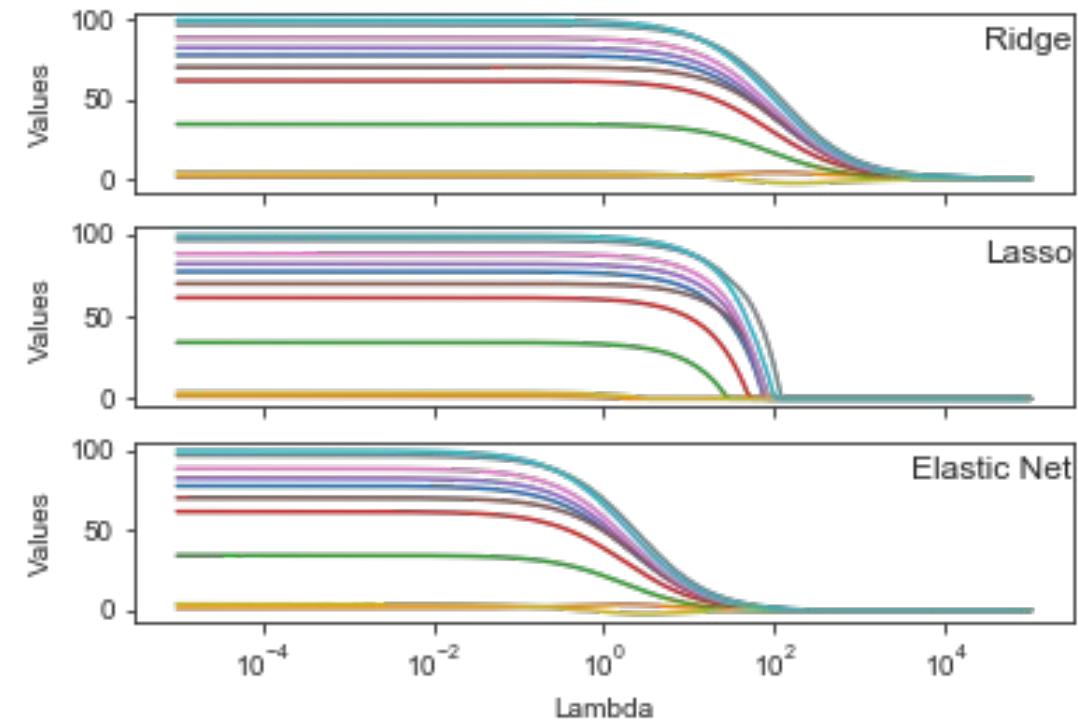
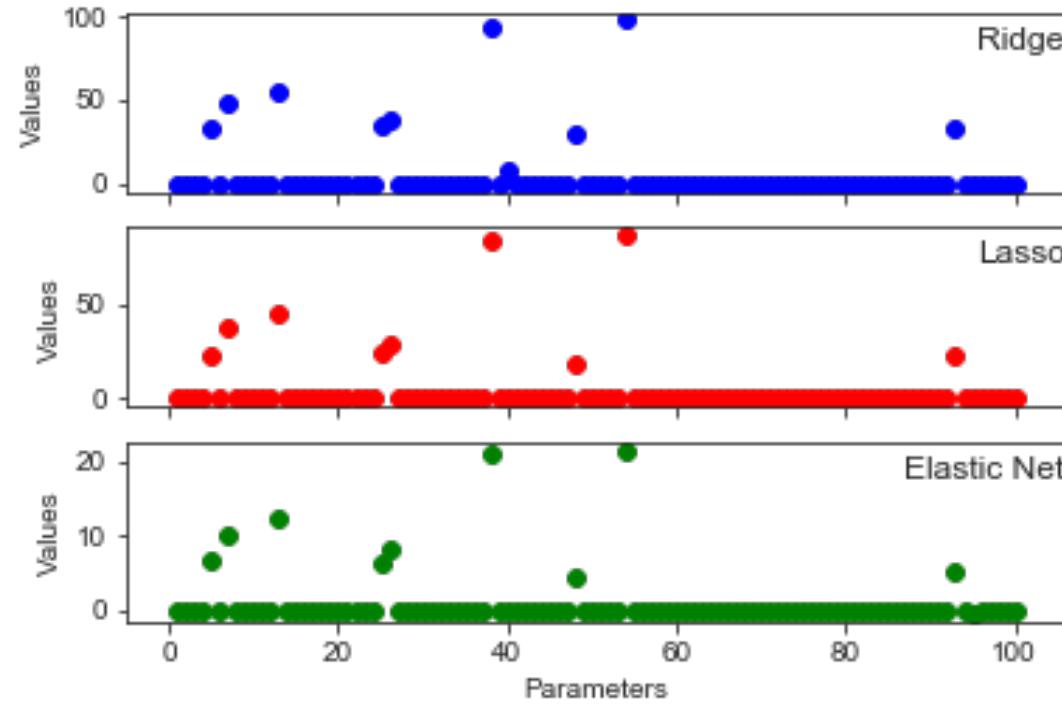
$\lambda \geq 0, \alpha \in [0, 1]$

$$\min_{\beta_0, \beta} \frac{1}{2n} \sum_{i=1}^n (y_i - \beta_0 - x_i^\top \beta)^2 + \lambda \left(\frac{1-\alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$

λ, α hyperparameters



Regularization



Linear Regression and Regularization

Linear

$$\min_{\beta} \frac{1}{2n} \|y - X\beta\|_2^2$$

Ridge

$$\min_{\beta} \frac{1}{2n} \|y - X\beta\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2$$

Lasso

$$\min_{\beta} \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$$

Elastic net

$$\min_{\beta} \frac{1}{2n} \|y - X\beta\|_2^2 + \lambda \left(\frac{1-\alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$

β_0 dropped
(output centered, input standardized)

$$X = \begin{bmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ x_n^\top \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$y = (y_1, \dots, y_n)^\top$$

$$\beta = (\beta_1, \dots, \beta_p)^\top$$



Regularization - Maximum a Posteriori Estimation

$$y_i \sim \mathcal{N}(\hat{y}_i, \sigma^2), i = 1, \dots, n \text{ and i.i.d}$$

$$\mathbb{P}(y|\beta) = \prod_{i=1}^n \mathbb{P}(y_i|\beta)$$

$$\mathbb{P}(\beta|y) \propto \mathbb{P}(y|\beta) \mathbb{P}(\beta)$$

β_0 dropped
(output centered, input standardized)

$$X = \begin{bmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ x_n^\top \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$y = (y_1, \dots, y_n)^\top$$

$$\beta = (\beta_1, \dots, \beta_p)^\top$$

$$\hat{\beta} = \arg \max_{\beta} \mathbb{P}(y|\beta) \mathbb{P}(\beta)$$

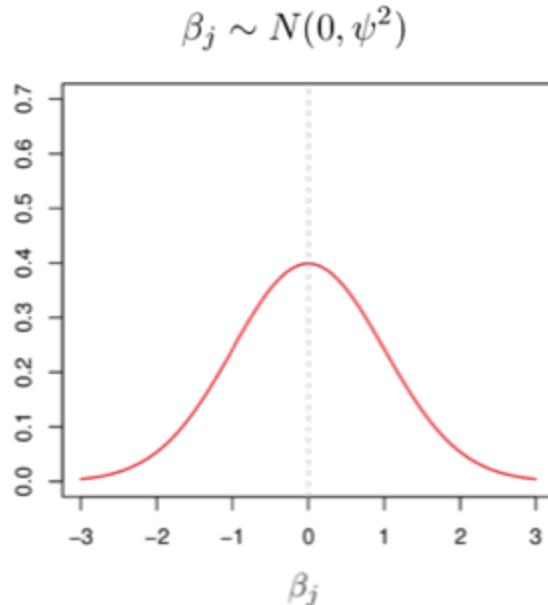
$$= \arg \max_{\beta} \log \mathbb{P}(y|\beta) + \log \mathbb{P}(\beta)$$

$$= \arg \max_{\beta} \sum_{i=1}^n \log \mathbb{P}(y_i|\beta) + \log \mathbb{P}(\beta)$$



Regularization - Maximum a Posteriori Estimation

prior, $\mathbb{P}(\beta)$



$$\hat{\beta} = \arg \max_{\beta} \sum_{i=1}^n \log \mathbb{P}(y_i | \beta) + \log \mathbb{P}(\beta)$$

constant prior

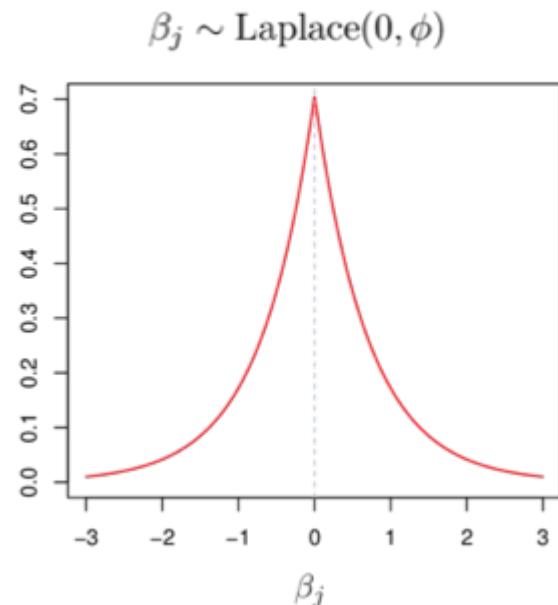
Linear

$\beta_j \sim \mathcal{N}(0, \psi^2), j = 1, \dots, p$ and i.i.d*

Ridge

$\beta_j \sim \text{Laplace}(0, \phi), j = 1, \dots, p$ and i.i.d*

Lasso



*
$$\mathbb{P}(\beta) = \prod_{j=1}^p \mathbb{P}(\beta_j)$$



Classification and Regularization

$$\{(x_i, y_i) : i = 1, \dots, n\}$$

$$x_i \in \mathbb{R}^p, \quad y_i \in \{0, 1\}$$

$$\mathbb{P}(Y = 1 | X = x_i) = \frac{1}{1 + e^{-(\beta_0 + x_i^\top \beta)}} = p(x_i)$$

$$\mathbb{P}(Y = 0 | X = x_i) = \frac{1}{1 + e^{(\beta_0 + x_i^\top \beta)}} = 1 - p(x_i)$$

mean of binary response

$$\mu(X) = \mathbb{P}(Y = 1 | X)$$

$$\text{Log Odds} \quad \log \left(\frac{\mu(X)}{1 - \mu(X)} \right) = \beta_0 + X^\top \beta$$

Maximizing
Log Likelihood

$$\max_{\beta_0, \beta} \frac{1}{n} \sum_{i=1}^n \underbrace{(y_i \log p(x_i) + (1 - y_i) \log(1 - p(x_i)))}_{y_i(\beta_0 + x_i^\top \beta) - \log(1 + e^{(\beta_0 + x_i^\top \beta)})}$$



Classification and Regularization

Logistic Regression

$$\min_{\beta_0, \beta} \frac{1}{n} \sum_{i=1}^n \left(\log(1 + e^{(\beta_0 + x_i^\top \beta)}) - y_i(\beta_0 + x_i^\top \beta) \right)$$

Logistic Regression with Regularization

$$\min_{\beta_0, \beta} \frac{1}{n} \sum_{i=1}^n \left(\log(1 + e^{(\beta_0 + x_i^\top \beta)}) - y_i(\beta_0 + x_i^\top \beta) \right) + \lambda \left(\frac{1-\alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$

$\lambda \geq 0, \alpha \in [0, 1]$ hyperparameters



Interpretation

$$p = 1$$

Features

balance
income
student (1) or not (0)

	Coefficient	Std. Error	Z-statistics	p-value
Intercept	-3.5041	0.0707	-49.55	< 0.0001
Student (1)	0.4049	0.1150	3.52	0.0004

Output

default (1) or not (0)

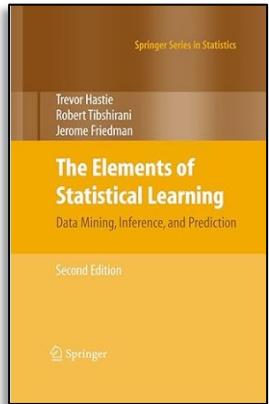
$$p > 1$$

	Coefficient	Std. Error	Z-statistics	p-value
Intercept	-10.8690	0.4923	-22.08	< 0.0001
Balance	0.0057	0.0002	24.74	< 0.0001
Income	0.0030	0.0082	0.37	0.7115
Student (1)	-0.6468	0.2362	-2.74	0.0062

Confounding Effect



Interpretation



[\(link\)](#)

TABLE 4.2. Results from a logistic regression fit to the South African heart disease data.

	Coefficient	Std. Error	Z Score
(Intercept)	-4.130	0.964	-4.285
sbp	0.006	0.006	1.023
tobacco	0.080	0.026	3.034
ldl	0.185	0.057	3.219
famhist	0.939	0.225	4.178
obesity	-0.035	0.029	-1.187
alcohol	0.001	0.004	0.136
age	0.043	0.010	4.184

TABLE 4.3. Results from stepwise logistic regression fit to South African heart disease data.

	Coefficient	Std. Error	Z score
(Intercept)	-4.204	0.498	-8.45
tobacco	0.081	0.026	3.16
ldl	0.168	0.054	3.09
famhist	0.924	0.223	4.14
age	0.044	0.010	4.52

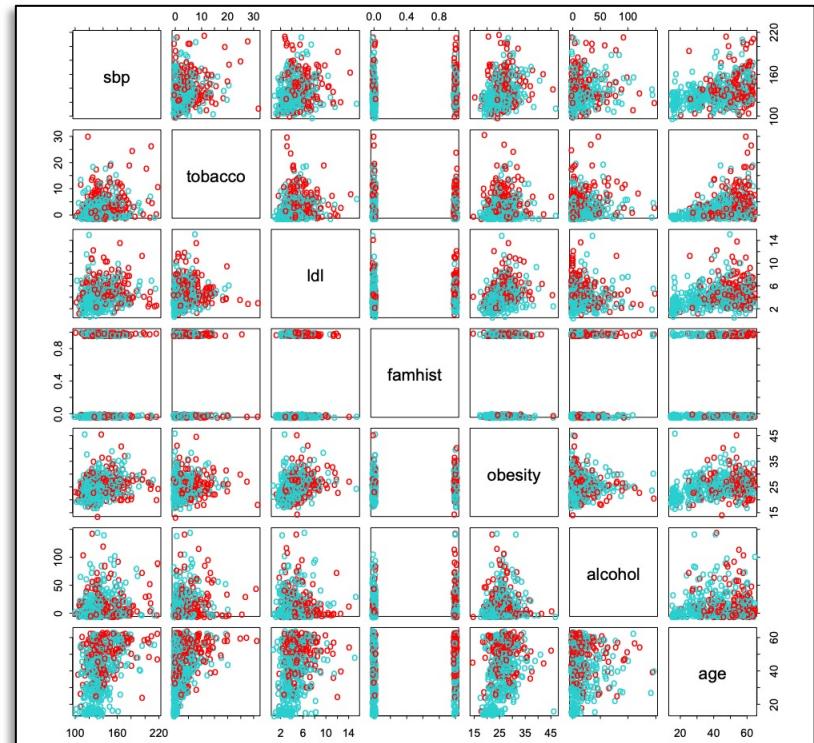
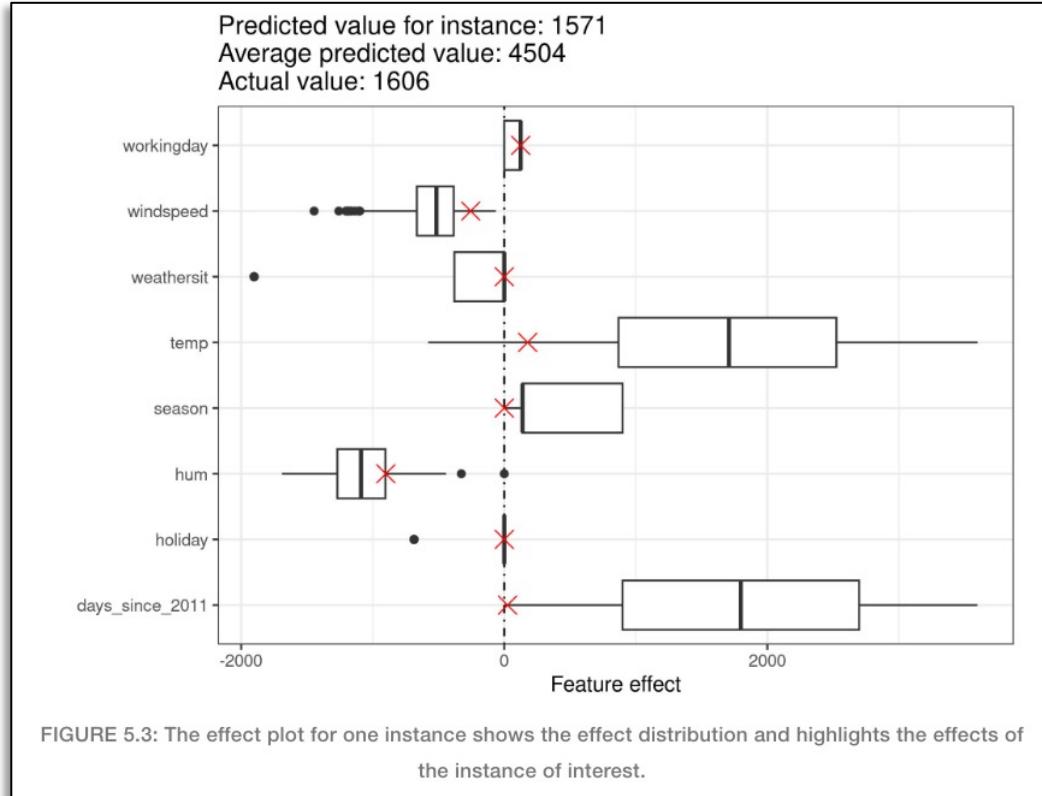


FIGURE 4.12. A scatterplot matrix of the South African heart disease data. Each plot shows a pair of risk factors, and the cases and controls are color coded (red is a case). The variable family history of heart disease (*famhist*) is binary (yes or no).



Interpretation



[\(link\)](#)

©2022 Christoph Molnar
Anything missing? Please send feedback to christoph.molnar.ai@gmail.com
version 1.0

Logistic Regression – Interpretation Cheat Sheet

Summary
Logistic regression is a statistical model used to predict the probability of an outcome. *Example:* Predict if a customer will default on their credit card loan ($Y=1$) or not ($Y=0$) based on income, remaining balance, and student status.

Example: loan default

covariate	β	$\exp(\beta)$	95% CI of β	p-value
Intercept	-11	0.0	[-14.4; -8.4]	0.00
student	-1.8	0.16	[−3.4; −0.3]	0.02
balance	0.62	1.66	[0.48; 0.79]	0.00
income	-0.25	0.78	[-5.0; 4.6]	0.92

p-value
The p-value for a coefficient β_j is the probability to observe this value or a more extreme one under the Null Hypothesis of $\beta_j = 0$ (meaning that the covariate doesn't affect the odds of the outcome). If $p < \alpha$, a pre-defined threshold (often $\alpha = 0.05$), the coefficient is significantly different from 0. *Example:* The effect of student status is significantly different from 0 ($p = 0.02$).

Probability

$$\pi = P(Y = 1|X) = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p))}$$

π is the probability of an outcome ($Y = 1$) given covariates X_1, \dots, X_p and coefficients $\beta_0, \beta_1, \dots, \beta_p$. *Example:* Probability of loan default as a function of income, credit balance, and student status. Unfortunately, coefficients can't be interpreted on the level of probabilities, but we have to use (log) odds.

Intercept

$\exp(\beta_0)$ is the odds of $Y = 1$ given all continuous covariates are 0 and all categorical covariates are set to the reference category. *Example:* The odds of loan default is $\exp(-11) \approx 0.000017$ given student=no, balance=0 and income=0.

Odds

$$\text{Odds} = \frac{\pi}{1 - \pi} = \exp(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p)$$

Expected no. of $Y = 1$ per $Y = 0$. *Example:* Odds of 1.5 means 1.5 loan defaults per non-defaults.

Continuous Coefficient

Template: A one-unit increase in X_j changes the odds for $Y = 1$ multiplicatively by $\exp(\beta_j)$, holding all other covariates constant. *Example:* A \$100 increase in credit balance increases the odds for loan default 1.86-fold, holding student status and income constant.

- Positive $\beta_j \rightarrow$ increasing X_j increases odds.
- Negative $\beta_j \rightarrow$ increasing X_j decreases odds.

Log Odds

$$\ln\left(\frac{\pi}{1 - \pi}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

Logarithm of odds. Also called logits. All following interpretation templates work for log odds when replacing $\exp(\beta_j)$ with β_j "odds" with "log odds", and removing "multiplicatively".

Categorical Coefficient

Template: Changing X_j from [reference] to [other category] changes the odds for $Y = 1$ multiplicatively by $\exp(\beta_j)$, holding all other covariates constant. *Example:* Changing student status from "no" to "yes" decreases the odds of loan default by a factor of 0.16, holding all other covariates constant.

(Log) Odds Ratio

$$OR = \frac{\text{Odds}_1}{\text{Odds}_2} \quad \text{and} \quad \log OR = \ln\left(\frac{\text{Odds}_1}{\text{Odds}_2}\right)$$

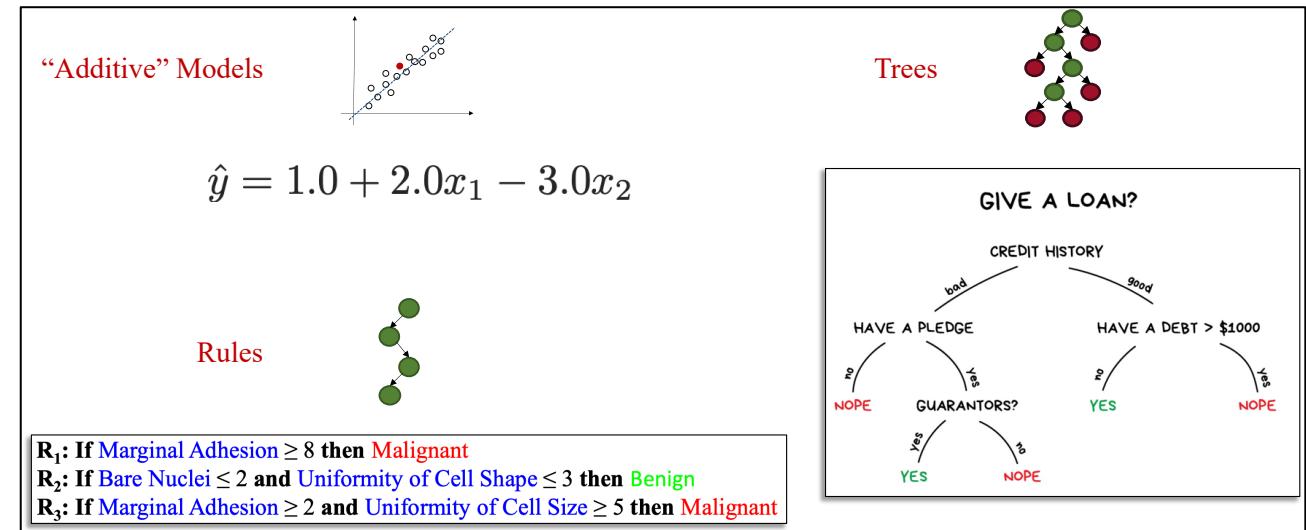
(Log) odds ratios compare (log) odds of two data points. Coefficient interpretations are based on (log) odds ratios, as they compare data points with $X_j + 1$ and with X_j : $OR = \text{Odds}(X_j + 1) / \text{Odds}(X_j) = \exp(\beta_j)$. For categorical covariates the ratio reflects the change from reference to another category. I recommend just using the provided templates and not getting a headache from the (log) odds ratios.

[\(link\)](#)



Big Picture

- Linear Models
- Generalized Additive Models
- Explainable Boosting Machines
- Decision Trees
- Rule-based Approaches



Generalized Additive Models (GAMs)

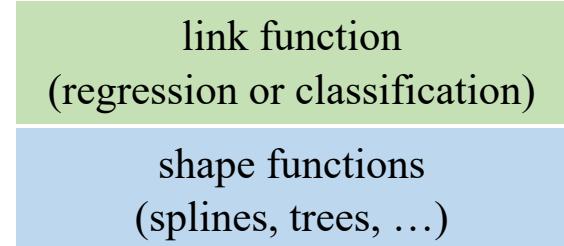
Model	Form	Intelligibility	Accuracy
Linear Model	$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$	+++	+
Generalized Linear Model	$g(y) = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$	+++	+
Additive Model	$y = f_1(x_1) + \dots + f_n(x_n)$	++	++
Generalized Additive Model	$g(y) = f_1(x_1) + \dots + f_n(x_n)$	++	++
Full Complexity Model	$y = f(x_1, \dots, x_n)$	+	+++

Table 1: From Linear to Additive Models.

([link](#))

$$g(\mathbb{E}[Y]) = f_0 + \sum_{j=1}^p f_j(X_j)$$

contribution of each feature



Backfitting: The shape functions are trained in-place* sequentially by fitting each one to the residual up to that iteration.

When (gradient) boosting is used to learn the shape functions, at each iteration of boosting the algorithm cycles through the features.

* previously trained version is discarded



Generalized Additive Models (GAMs)

regression
(*identity*)

binary classification
(*logit*)

count data
(*log*)

$$g(\mathbb{E}[Y]) = f_0 + \sum_{j=1}^p f_j(X_j)$$

$$g(\mathbb{E}[Y]) = Y$$

$$g(\mathbb{E}[Y]) = \log\left(\frac{\mathbb{E}[Y]}{1 - \mathbb{E}[Y]}\right)$$

$$g(\mathbb{E}[Y]) = \log(\mathbb{E}[Y])$$

link function
(regression or classification)

$$\log\left(\frac{\mu(X)}{1 - \mu(X)}\right) = \beta_0 + X^\top \beta$$

Link function:

We specify this using: `GAM(link='...')`

Link functions take the distribution mean to the linear prediction.

- `'identity'`
- `'logit'`
- `'inverse'`
- `'log'`
- `'inverse-squared'`

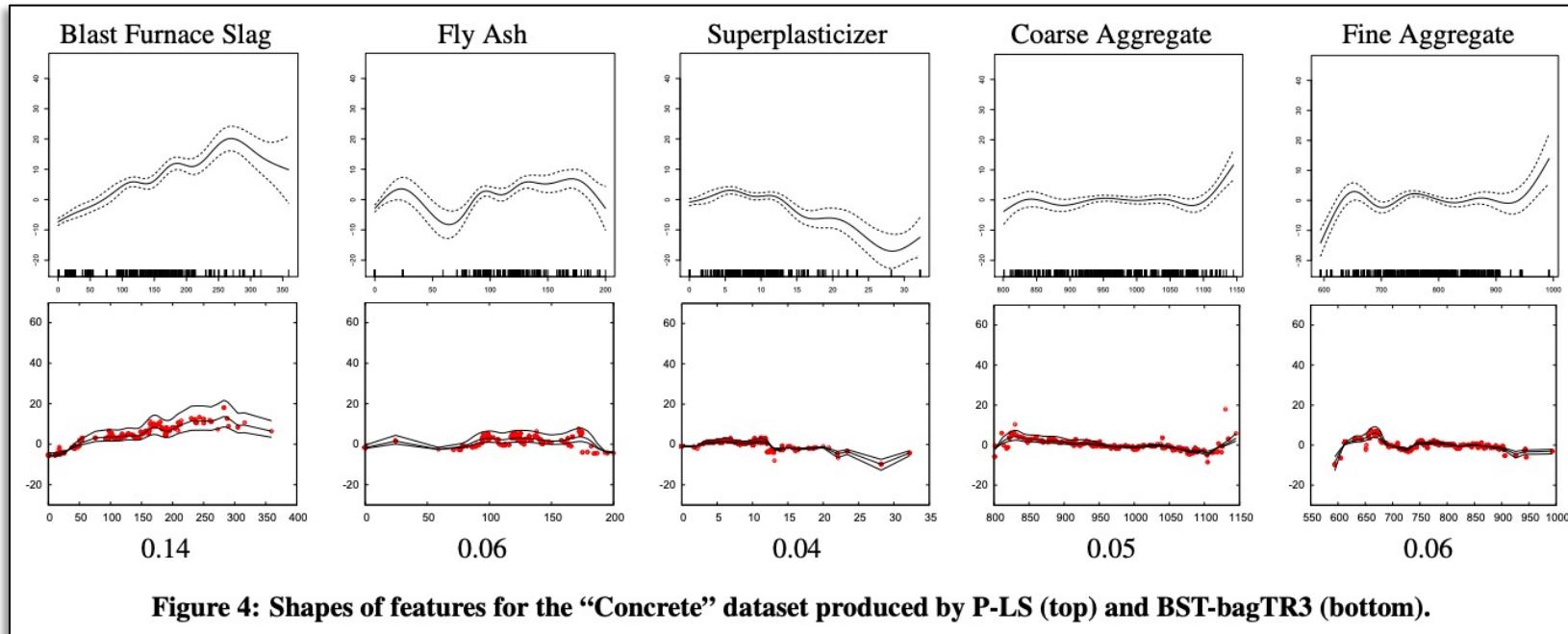


common models:

- `LinearGAM` identity link and normal distribution
- `LogisticGAM` logit link and binomial distribution
- `PoissonGAM` log link and Poisson distribution
- `GammaGAM` log link and gamma distribution
- `InvGauss` log link and inv_gauss distribution

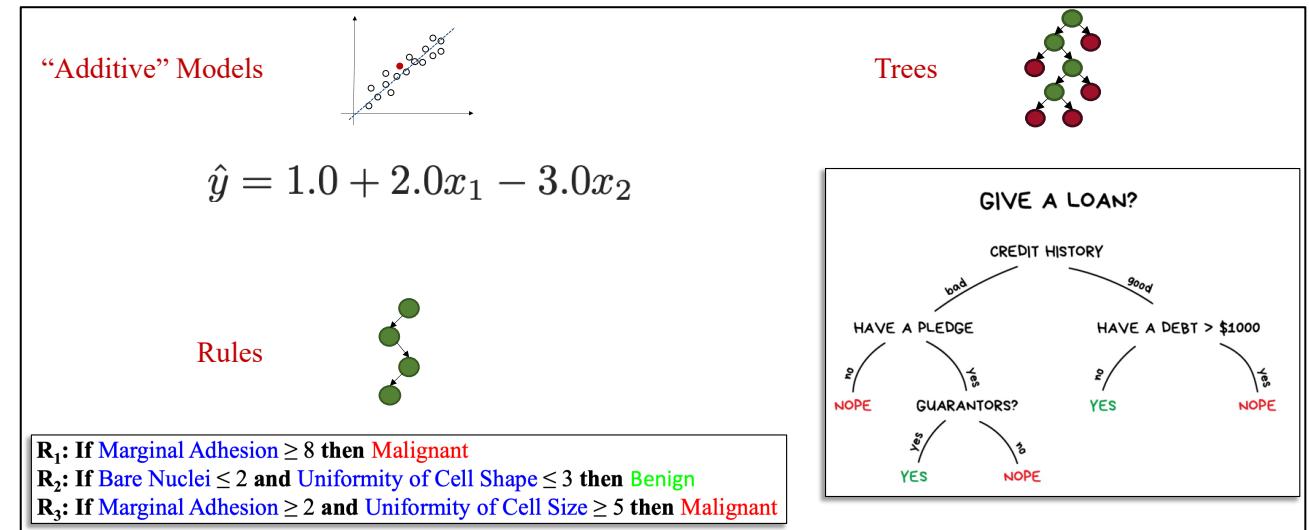


Interpretation



Big Picture

- Linear Models
- Generalized Additive Models
- Explainable Boosting Machines
- Decision Trees
- Rule-based Approaches



Explainable Boosting Machines (EBMs)

Our first contribution in this paper is to build models that are more powerful than GAMs, but are still intelligible. We observe that two-dimensional interactions can still be rendered as heatmaps of $f_{ij}(x_i, x_j)$ on the two-dimensional x_i, x_j -plane, and thus a model that includes only one- and two-dimensional components is still intelligible.

([link](#))

$$g(\mathbb{E}[Y]) = f_0 + \sum_{j=1}^p f_j(X_j) + \sum_{j=1}^p \sum_{k \neq j}^p f_{jk}(X_j, X_k)$$

one- and two-dimensional
shape functions

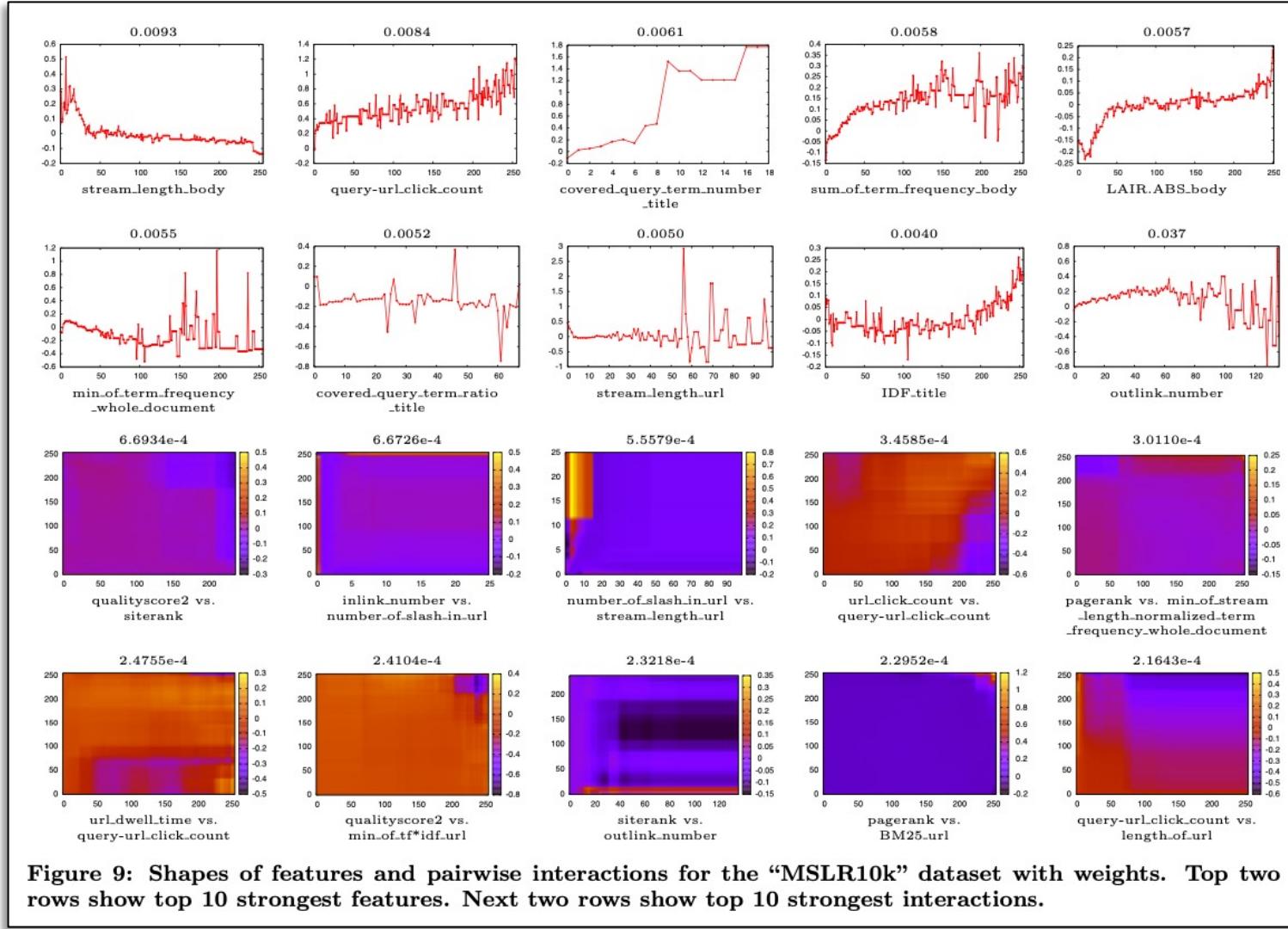
cross-interaction between
two features

When the number of feature is high, considering all interactions is very costly.
Thus, only **strong** interactions are used after ranking the interactions.

Testing strong interactions: ANOVA, partial dependence, statistical tests,
greedy forward selection strategy.

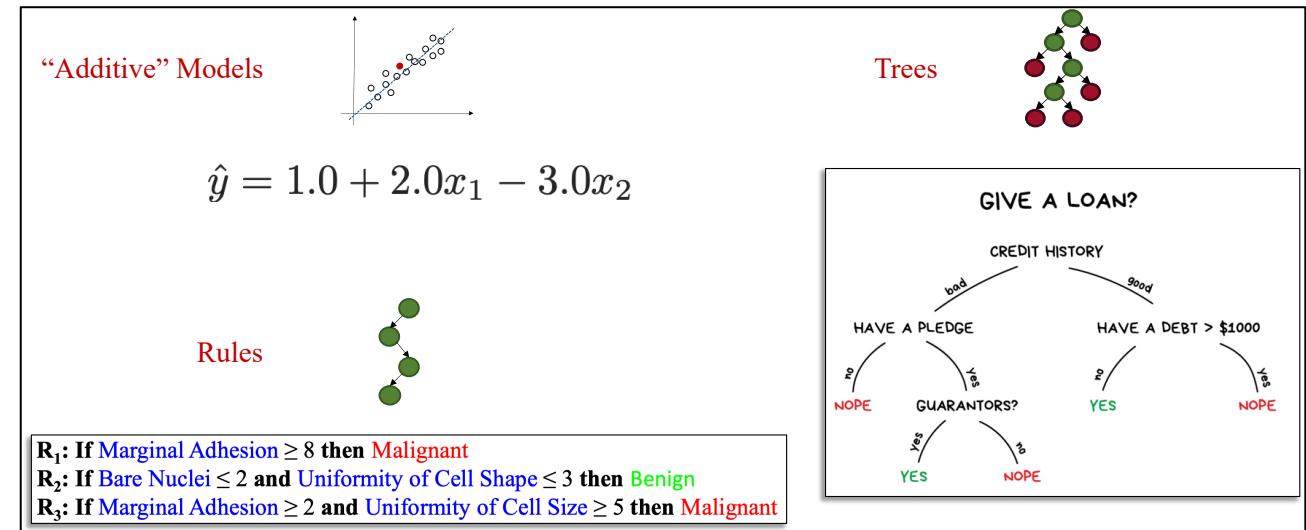


Interpretation

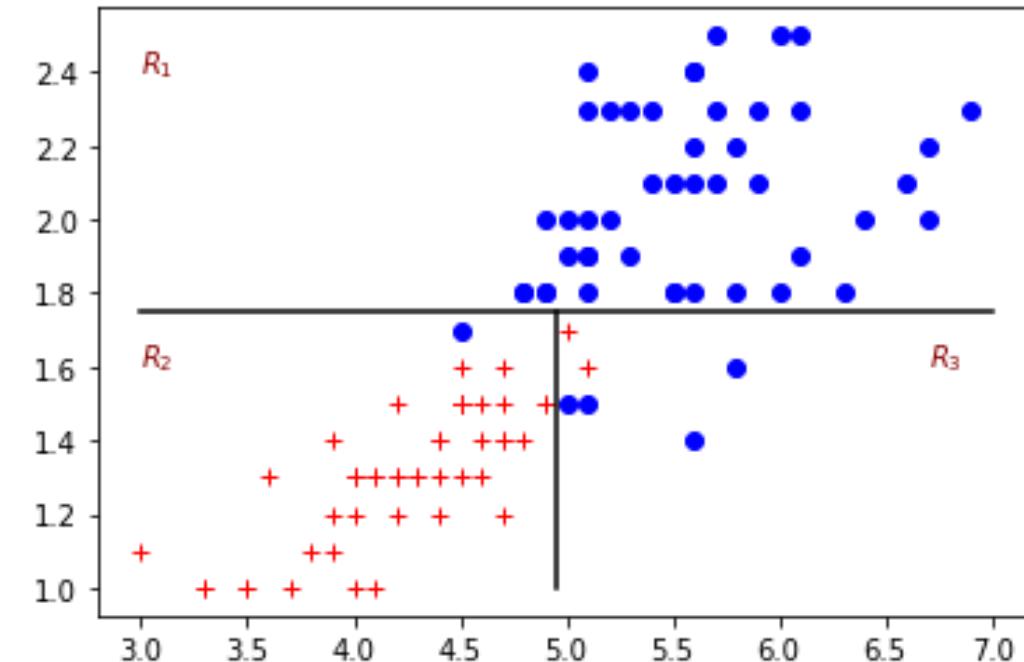
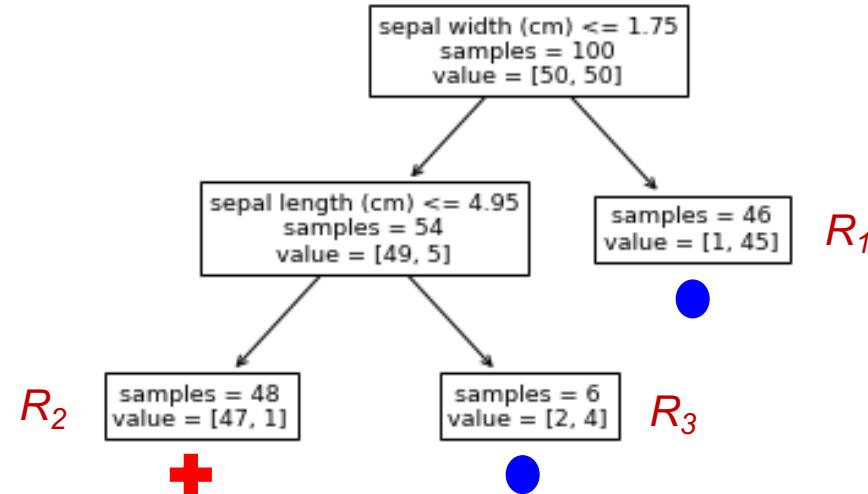


Big Picture

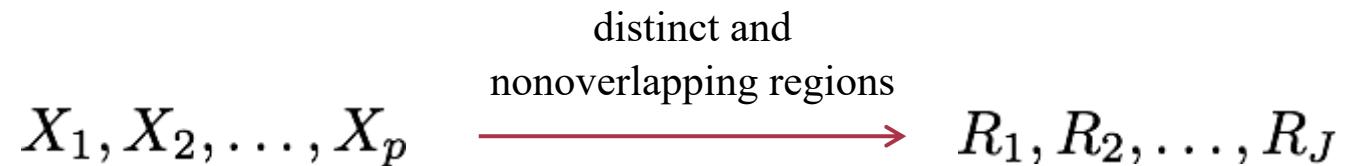
- Linear Models
- Generalized Additive Models
- Explainable Boosting Machines
- Decision Trees
- Rule-based Approaches



Example Tree



Regression Trees



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

R_1, R_2, \dots, R_J ?

Goal: Finding the regions such that RSS is minimized

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

\hat{y}_{R_j} : mean response within R_j



Regression Trees

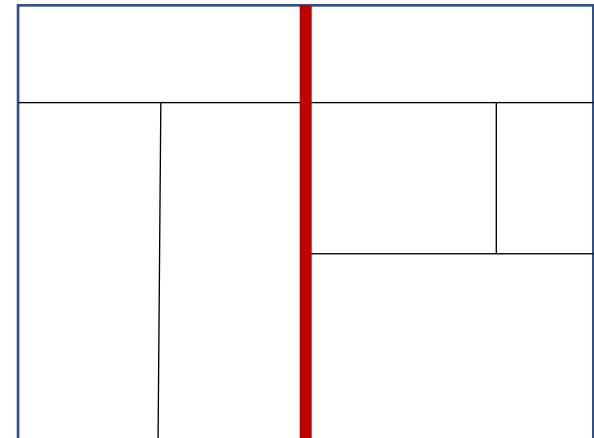
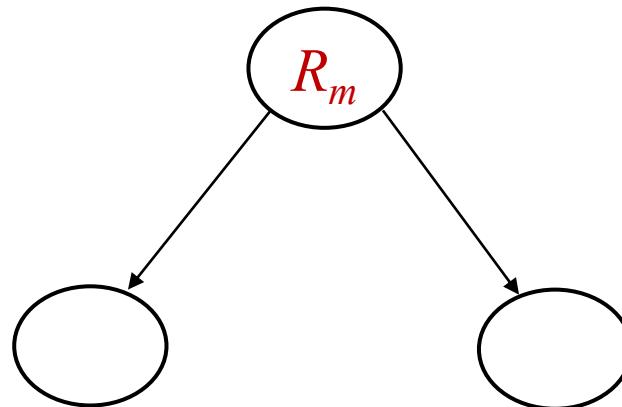
Recursive Binary Splitting

n_{R_1} : number of samples in region $R_1(j, s)$

n_{R_2} : number of samples in region $R_2(j, s)$

$$n_R = n_{R_1} + n_{R_2}$$

$$E(R_m) = \sum_{i:x_i \in R_m} (y_i - \hat{y}_{R_m})^2$$



$$R_1(j, s) = \{X | X_j < s\} \quad R_2(j, s) = \{X | X_j \geq s\}$$

Find j and s that minimizes the weighted error

$$\frac{n_{R_1}}{n_R} E(R_1(j, s)) + \frac{n_{R_2}}{n_R} E(R_2(j, s))$$

`sklearn.tree.DecisionTreeRegressor`

Parameters: `criterion : {"squared_error", "friedman_mse", "absolute_error", "poisson"}, default="squared_error"`



Regression Trees

Tree Pruning

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

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

$T \subset T_0$: subtree

$|T|$: number of terminal nodes in T

R_m : region corresponding to
 the m th terminal node

α : fixed parameter

$\alpha \uparrow$ $|T| \downarrow$

Use k -fold cross validation to choose α



Regression Trees

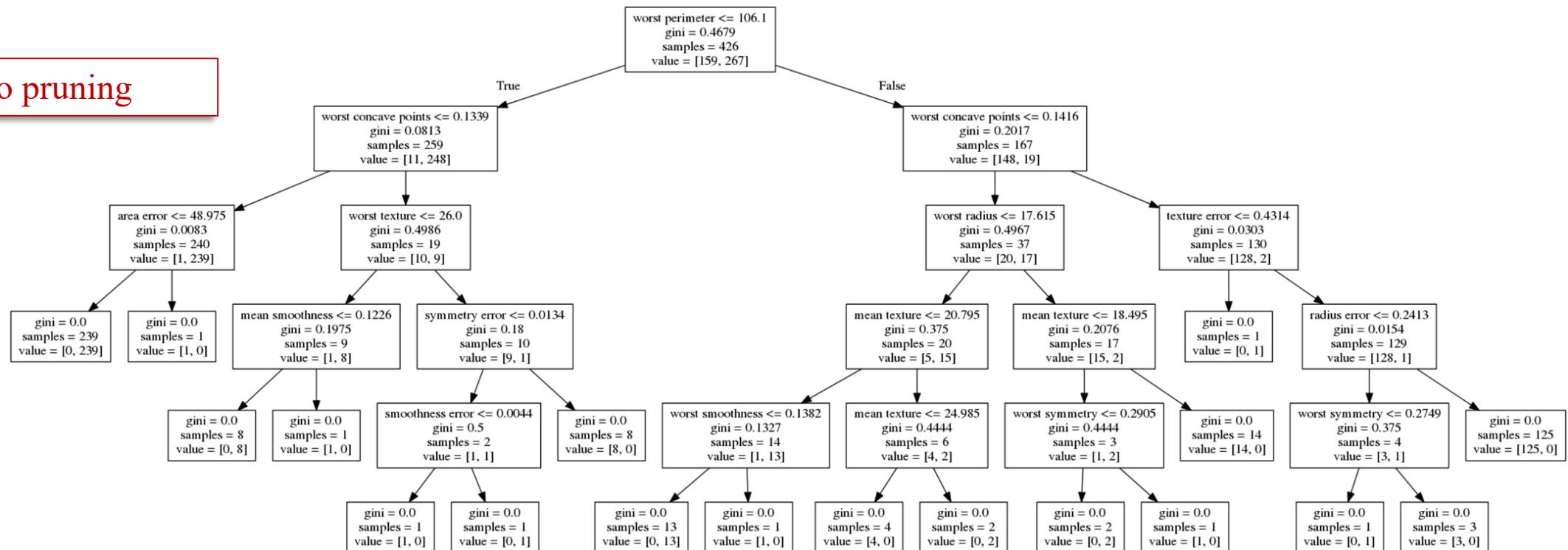
Tree Pruning

`sklearn.tree.DecisionTreeRegressor`

```
class sklearn.tree.DecisionTreeRegressor(*, criterion='squared_error', splitter='best', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None, min_impurity_decrease=0.0, ccp_alpha=0.0) [source]
```

[source]

no pruning



(link)

Regression Trees

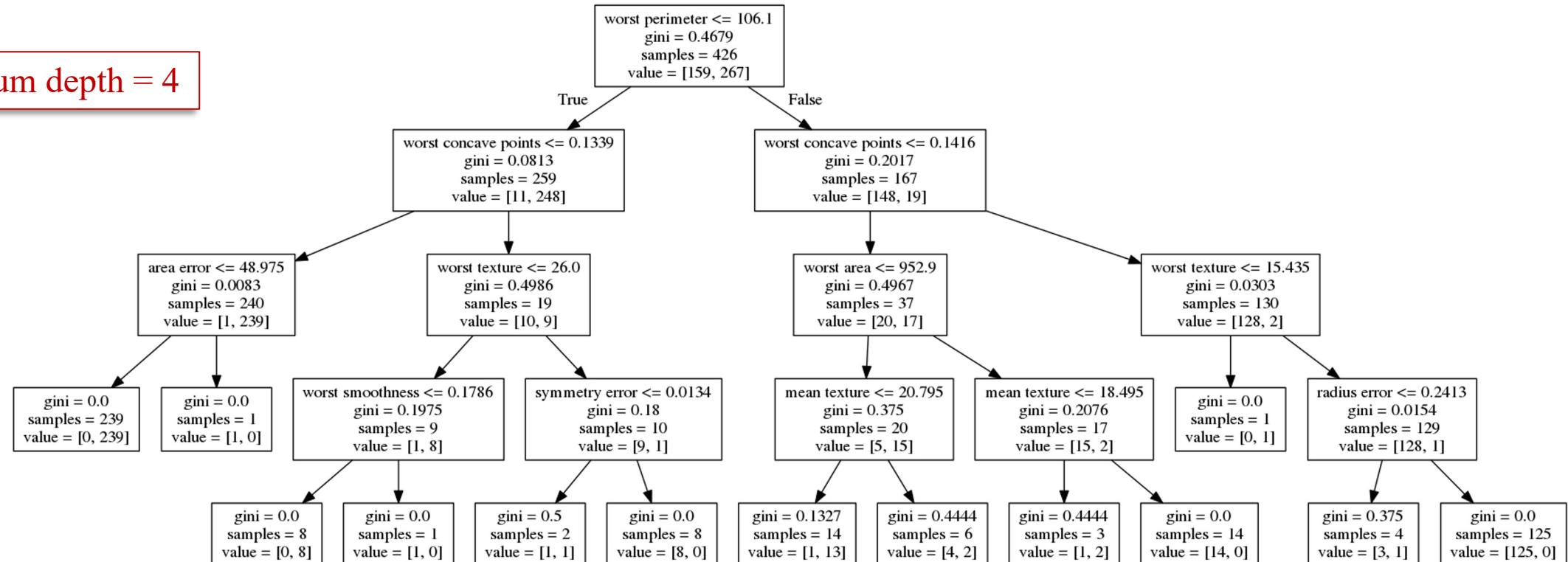
Tree Pruning

sklearn.tree.DecisionTreeRegressor

```
class sklearn.tree.DecisionTreeRegressor(*, criterion='squared_error', splitter='best', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, ccp_alpha=0.0)
```

[\[source\]](#)

maximum depth = 4


[\(link\)](#)


Regression Trees

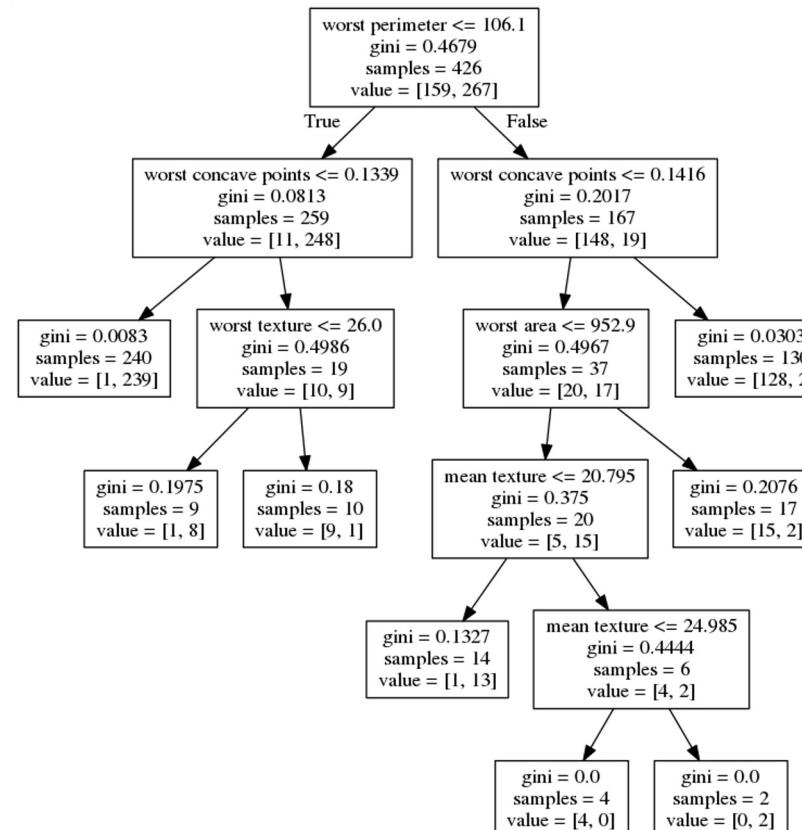
Tree Pruning

sklearn.tree.DecisionTreeRegressor

```
class sklearn.tree.DecisionTreeRegressor(*, criterion='squared_error', splitter='best', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, ccp_alpha=0.0)
```

[\[source\]](#)

maximum leaf nodes = 8


[\(link\)](#)


Regression Trees

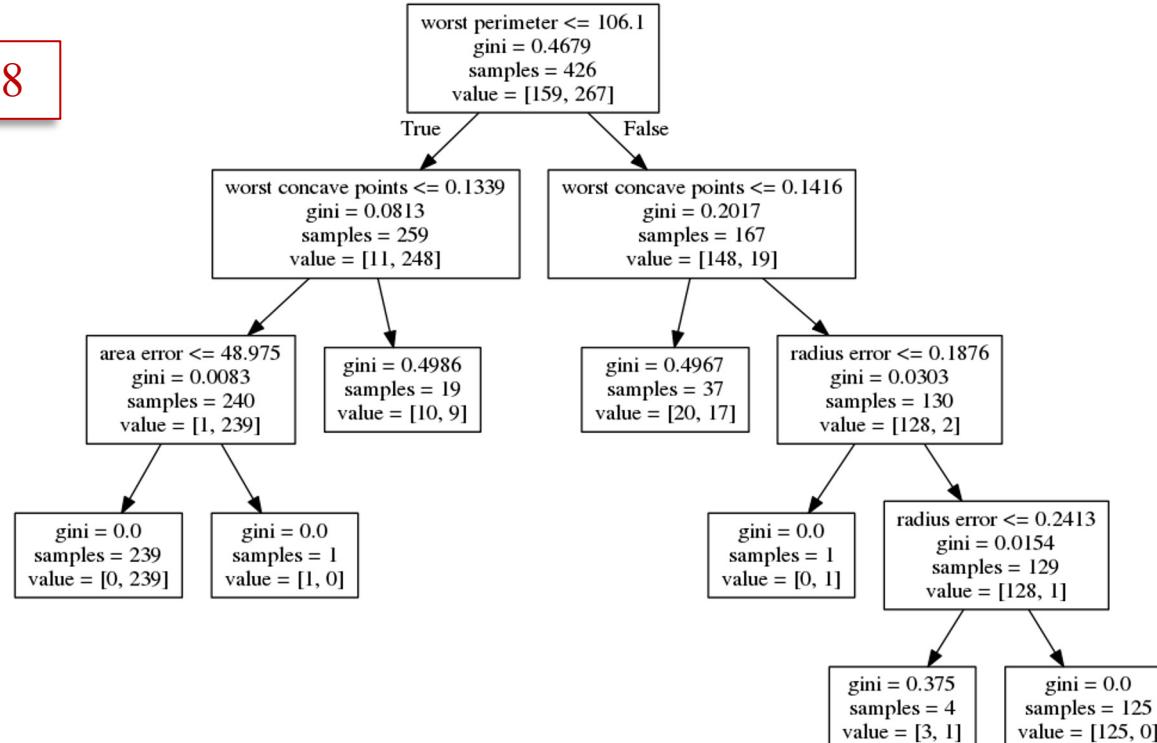
Tree Pruning

sklearn.tree.DecisionTreeRegressor

```
class sklearn.tree.DecisionTreeRegressor(*, criterion='squared_error', splitter='best', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=None, random_state=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, ccp_alpha=0.0)
```

[\[source\]](#)

minimum samples for split = 8


[\(link\)](#)


Classification Trees

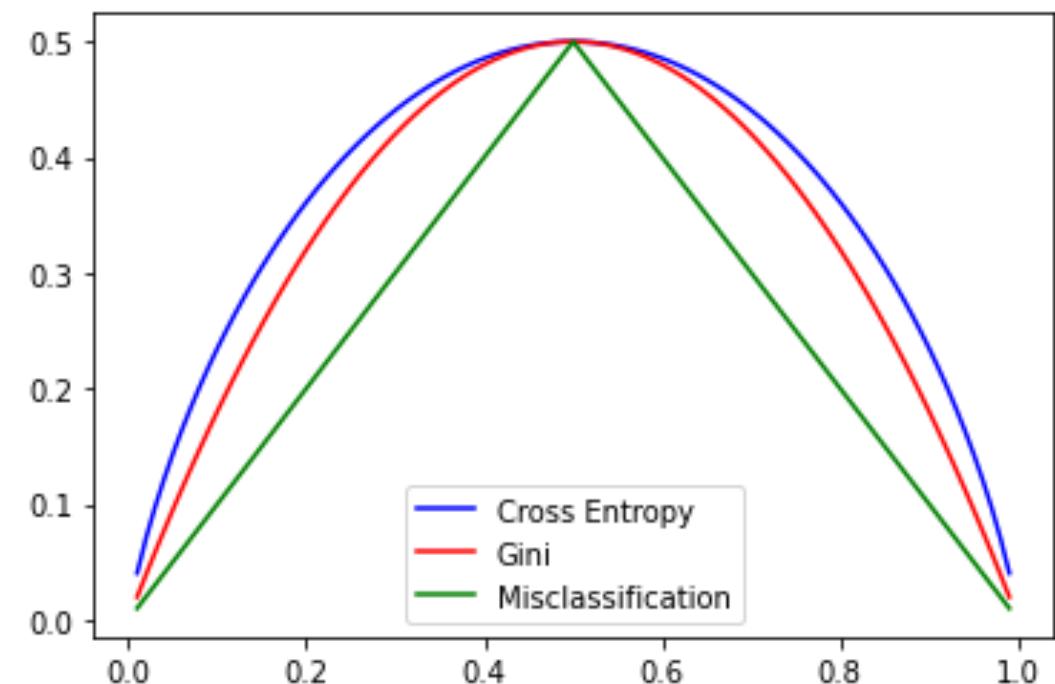
Similar to a regression tree but uses different error measures based on *purity* of a region

\hat{p}_{mk} : proportion of class- k training observations in the m th region

Misclassification error: $1 - \max_k \{\hat{p}_{mk}\}$

Cross entropy: $-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$

Gini index: $\sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$



Classification Trees

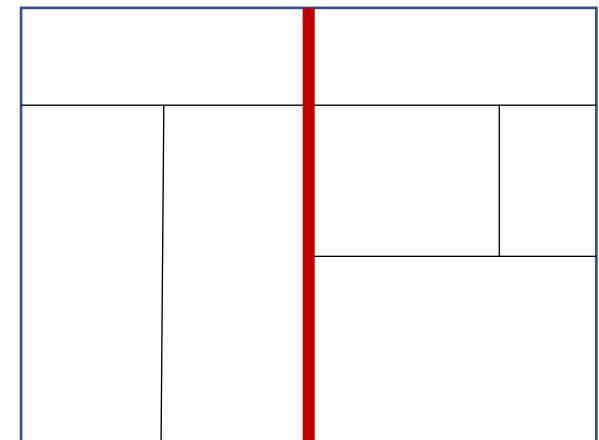
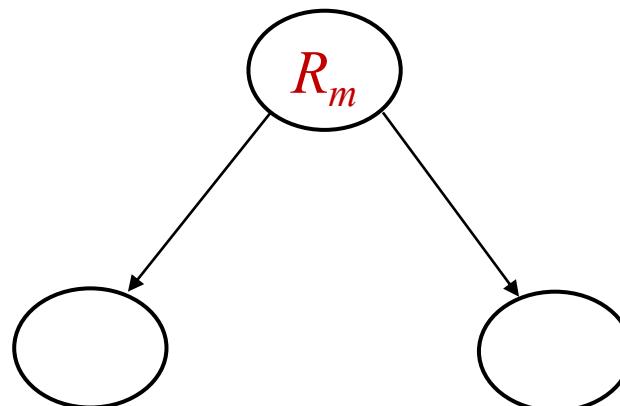
Recursive Binary Splitting

n_{R_1} : number of samples in region $R_1(j, s)$

n_{R_2} : number of samples in region $R_2(j, s)$

$$n_R = n_{R_1} + n_{R_2}$$

$$E(R_m) = 1 - \max_k \{\hat{p}_{mk}\}$$



$$R_1(j, s) = \{X | X_j < s\} \quad R_2(j, s) = \{X | X_j \geq s\}$$

Find j and s that minimizes the weighted error

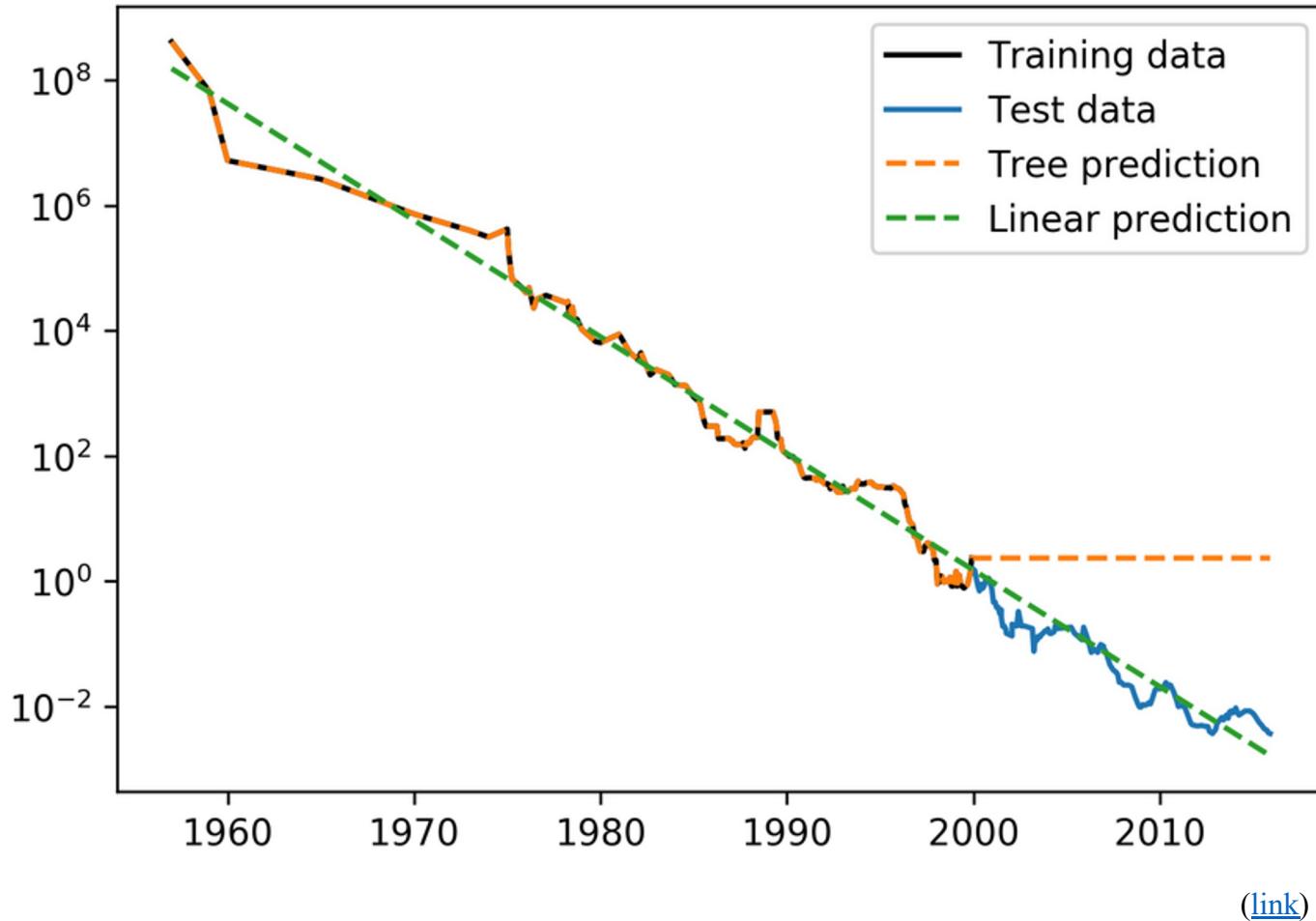
$$\frac{n_{R_1}}{n_R} E(R_1(j, s)) + \frac{n_{R_2}}{n_R} E(R_2(j, s))$$

`sklearn.tree.DecisionTreeClassifier`

Parameters: `criterion : {"gini", "entropy", "log_loss"}, default="gini"`



Lack of Extrapolation

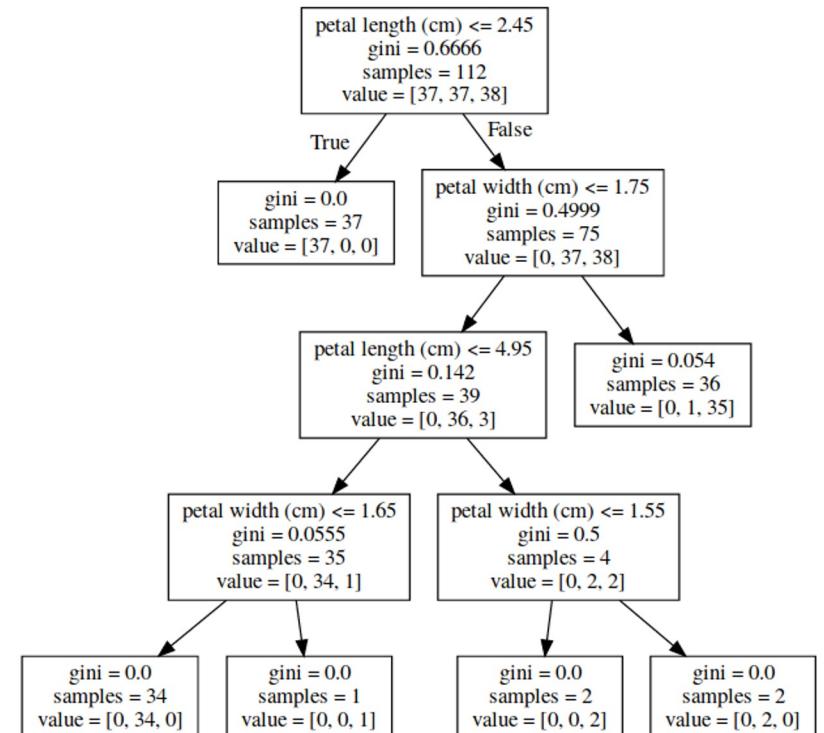


([link](#))



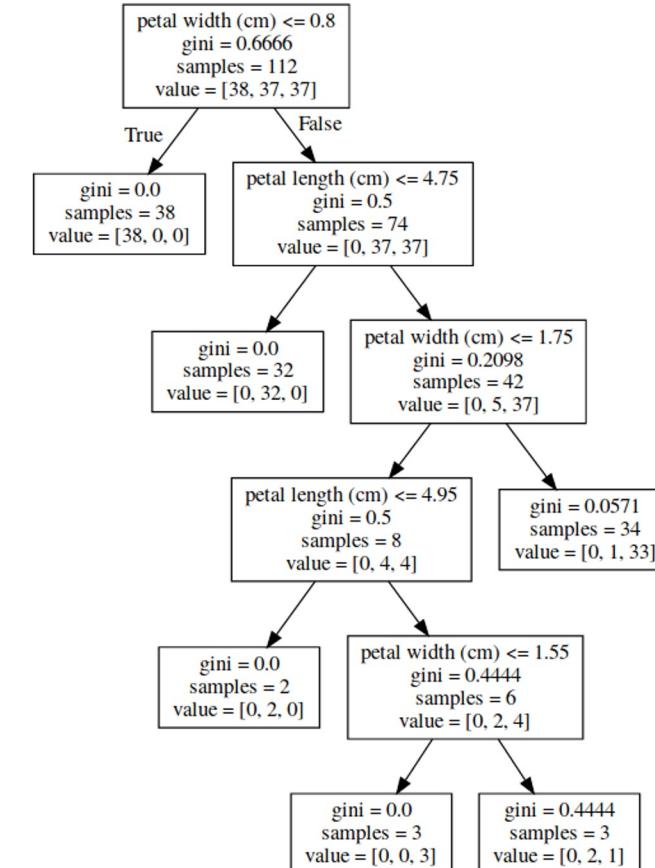
Instability

```
iris = load_iris()
X_train, X_test, y_train, y_test = train_test_split(
    iris.data, iris.target, stratify=iris.target, random_state=0)
tree = DecisionTreeClassifier(max_leaf_nodes=6).fit(X_train, y_train)
tree_dot = export_graphviz(tree, out_file=None, feature_names=iris.feature_names)
graphviz.Source(tree_dot)
```



[\(link\)](#)

```
X_train, X_test, y_train, y_test = train_test_split(
    iris.data, iris.target, stratify=iris.target, random_state=1)
tree = DecisionTreeClassifier(max_leaf_nodes=6).fit(X_train, y_train)
tree_dot = export_graphviz(tree, out_file=None, feature_names=iris.feature_names)
graphviz.Source(tree_dot)
```



[\(link\)](#)



Computational Complexity

Volume 5, number 1

INFORMATION PROCESSING LETTERS

May 1976

CONSTRUCTING OPTIMAL BINARY DECISION TREES IS NP-COMPLETE*

Laurent HYAFIL

IRIA – Laboria, 78150 Rocquencourt, France

and

Ronald L. RIVEST

Dept. of Electrical Engineering and Computer Science, M.I.T., Cambridge, Massachusetts 02139, USA

Received 7 November 1975, revised version received 26 January 1976

Binary decision trees, computational complexity, NP-complete

([link](#))

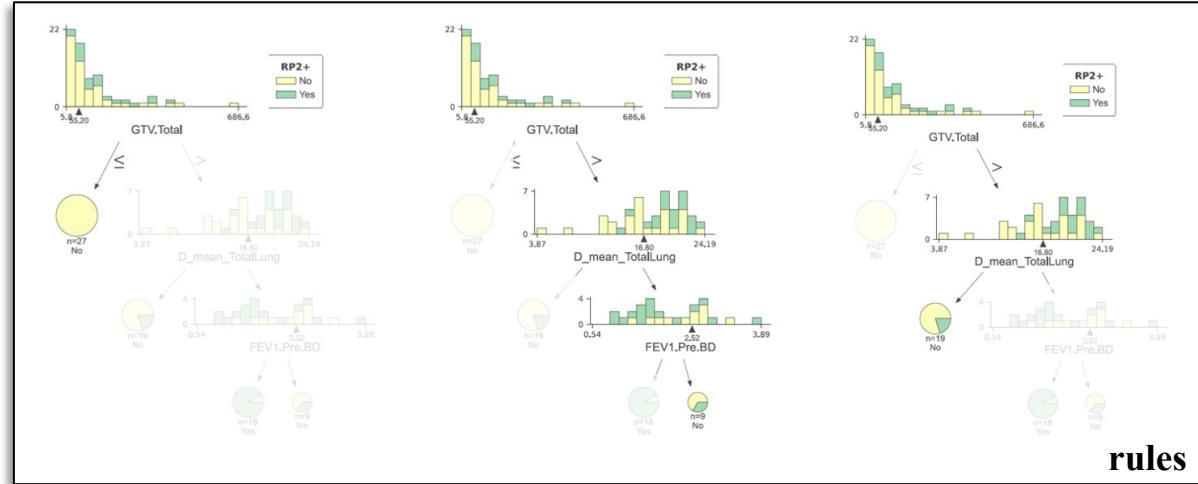
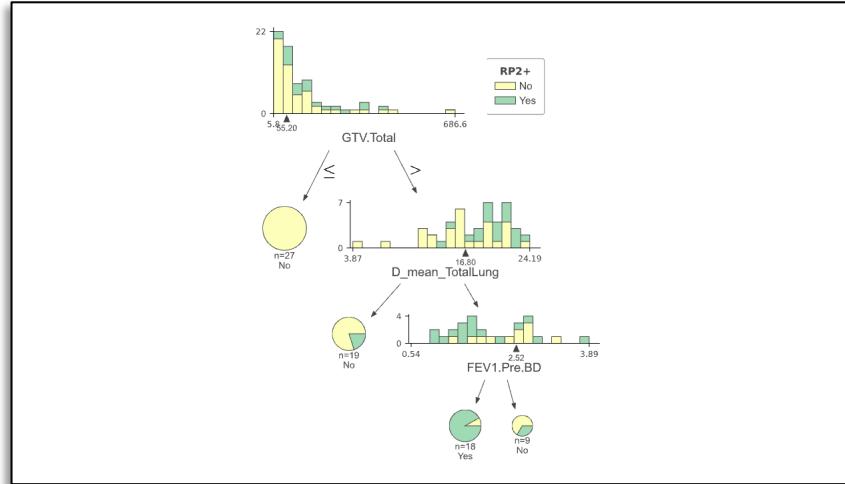
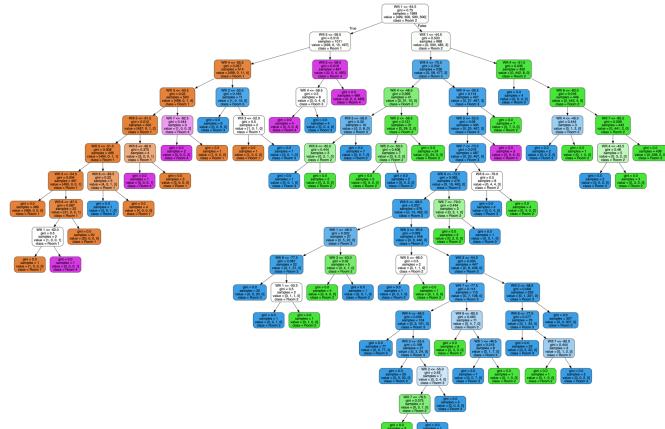
We demonstrate that constructing optimal binary decision trees is an NP-complete problem, where an optimal tree is one which minimizes the expected number of tests required to identify the unknown object.



Optimization for Tree Growing

$$\mathcal{D} = \{(x_i, y_i) : i \in \mathcal{I}\}$$

$$\mathcal{I} = \{1, \dots, n\}, \text{ sample indices}$$

**rules**[\(link\)](#)

$$\min_{T \in \mathcal{T}} \ell_{\mathcal{D}}(T) + \lambda \Omega(T)$$

total loss

model complexity

all possible trees

(Binary) Optimal Classification Trees: (Bin)OCT
 Optimal Sparse Decision Trees: OSDT \rightarrow FSDT

...



Optimization for Tree Growing

$$\min_{T \in \mathcal{T}} \ell_{\mathcal{D}}(T) + \lambda \Omega(T)$$

total loss
model complexity
all possible trees

classification error	number of branches 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 of points in leaf l	minimum number of samples in any leaf

Optimal Classification Trees (OCT)⁺

For a tree $d = (d_{un}, \delta_{un}, d_{\text{split}}, \delta_{\text{split}}, K, H_d)$, we define its objective function as a combination of the misclassification error and a sparsity penalty on the number of leaves:

$$R(d, \mathbf{x}, \mathbf{y}) = \ell(d, \mathbf{x}, \mathbf{y}) + \lambda H_d. \quad (4)$$

$R(d, \mathbf{x}, \mathbf{y})$ is a regularized empirical risk. The loss $\ell(d, \mathbf{x}, \mathbf{y})$ is the misclassification error of d , i.e., the fraction of training data with incorrectly predicted labels. H_d is the number of leaves in the tree d . λH_d is a regularization term that penalizes bigger trees.

Optimal Sparse Decision Trees (OSDT)[♦]

⁺ Bertsimas, D. & Dunn, J. “[Optimal classification trees](#),” Machine Learning, 106(7), 1039–1083, 2017.

[♦] Hu, X., Rudin, C. & Seltzer, M. “[Optimal sparse decision trees](#),” in Advances in Neural Information Processing Systems, 7267–7275, 2019.



Optimization for Tree Growing

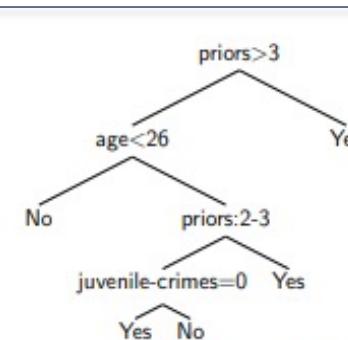
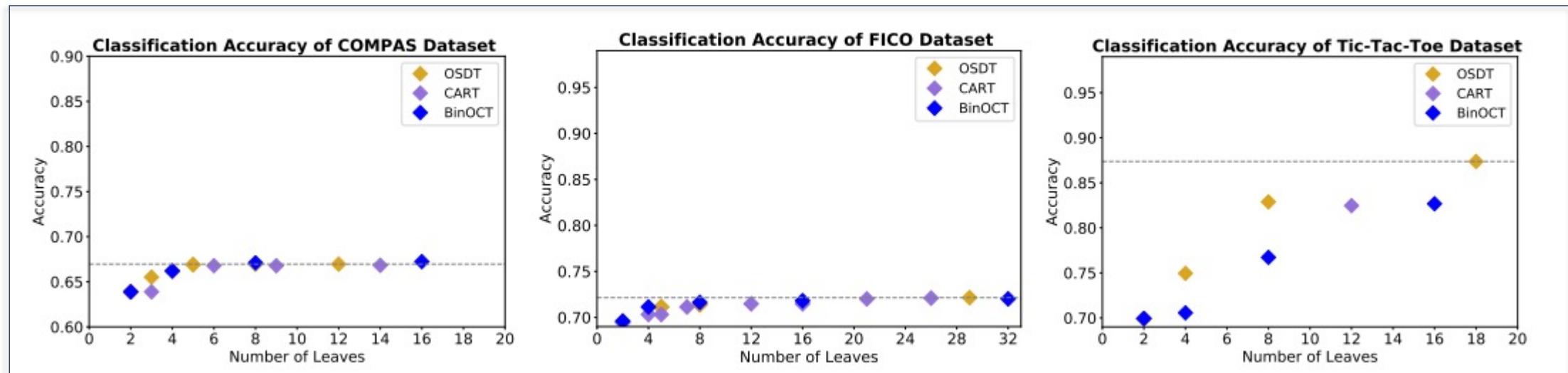


Figure 4: An optimal decision tree generated by OSDT on the COMPAS dataset. ($\lambda = 0.005$, accuracy: 66.90%)



Big Picture

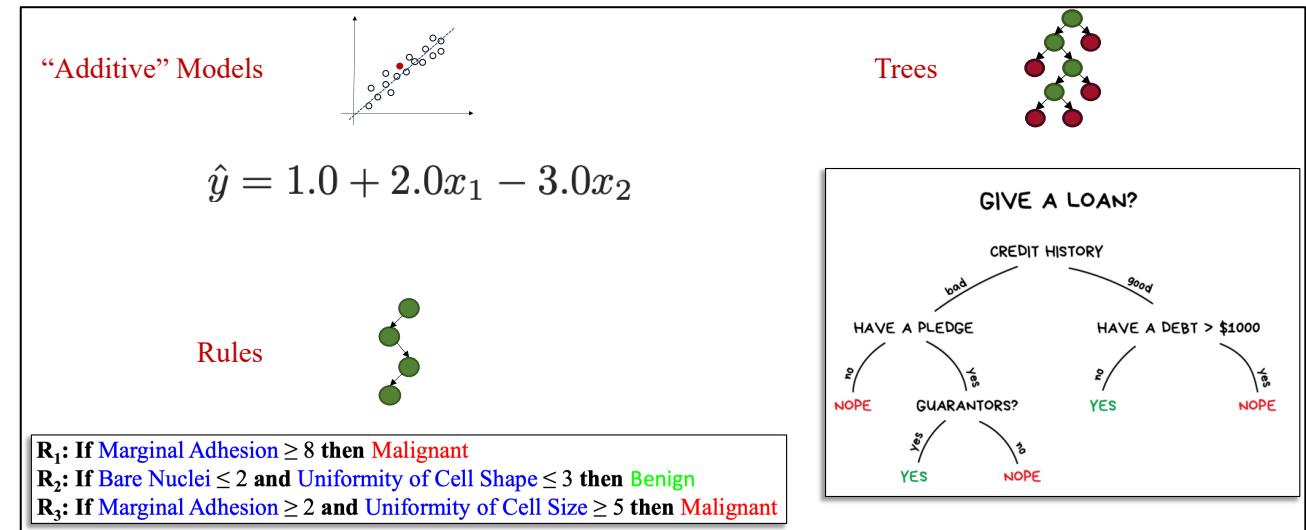
- Linear Models

- Generalized Additive Models

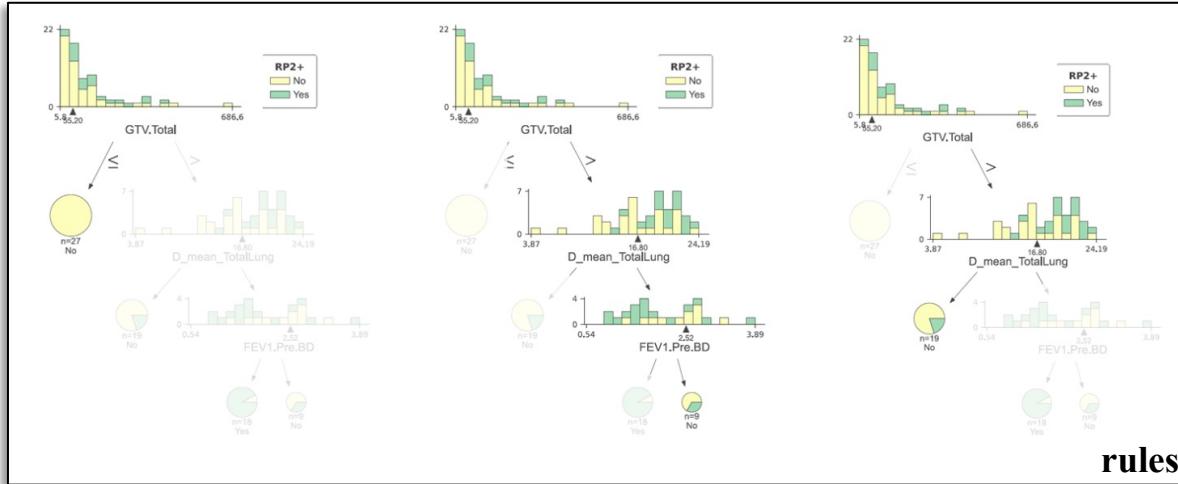
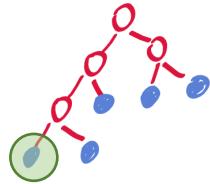
- Explainable Boosting Machines

- Decision Trees

- Rule-based Approaches

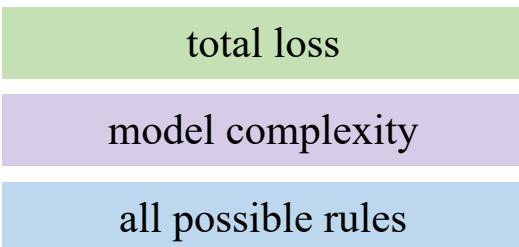


Decision Rules



- R₁:** If Marginal Adhesion ≥ 8 then Malignant
- R₂:** If Bare Nuclei ≤ 2 and Uniformity of Cell Shape ≤ 3 then Benign
- R₃:** If Marginal Adhesion ≥ 2 and Uniformity of Cell Size ≥ 5 then Malignant

Optimization



$$\min_{R \in \mathcal{R}} \ell_{\mathcal{D}}(R) + \lambda \Omega(R)$$

Boolean Decision Rules (BDR)
Rule Generation (RUG)

...

Machine Learning, 11, 63–91 (1993)
© 1993 Kluwer Academic Publishers, Boston. Manufactured in The Netherlands.

Very Simple Classification Rules Perform Well on Most Commonly Used Datasets

ROBERT C. HOLTE

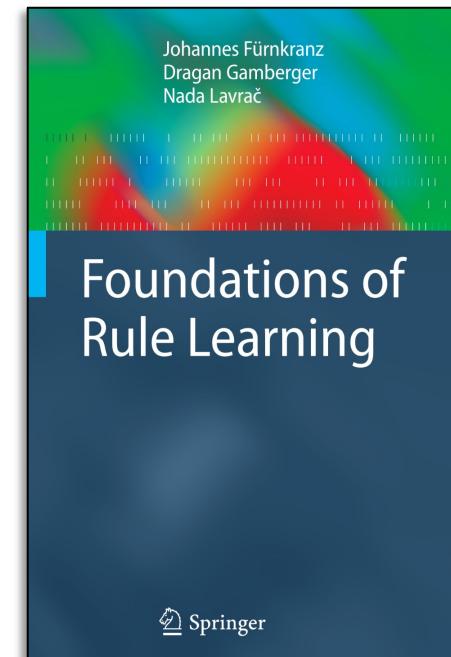
Computer Science Department, University of Ottawa, Ottawa, Canada K1N 6N5

HOLTE@CSI.UOTTAWA.CA

Editor: Bruce Porter

Abstract. This article reports an empirical investigation of the accuracy of rules that classify examples on the basis of a single attribute. On most datasets studied, the best of these very simple rules is as accurate as the rules induced by the majority of machine learning systems. The article explores the implications of this finding for machine learning research and applications.

Keywords: empirical learning, accuracy-complexity tradeoff, pruning, ID3



(link)



Rule-based Models

```

function LEARNRULEBASE( $\mathcal{E}$ )
  Input:
     $\mathcal{E}$  set of training examples

  Algorithm:
     $\mathcal{R} := \emptyset$ 
    for each class  $c_i$ ,  $i = 1$  to  $C$  do
       $\mathcal{P}_i := \{\text{subset of examples in } \mathcal{E} \text{ with class label } c_i\}$ 
       $\mathcal{N}_i := \{\text{subset of examples in } \mathcal{E} \text{ with other class labels}\}$ 
       $\mathcal{R}_i := \text{LEARNSETOFRULES}(c_i, \mathcal{P}_i, \mathcal{N}_i)$ 
       $\mathcal{R} := \mathcal{R} \cup \mathcal{R}_i$ 
    endfor
     $\mathcal{R} := \mathcal{R} \cup \{\text{default rule } (c_{max} \leftarrow \text{true})\}$ 
      where  $c_{max}$  is the majority class in  $\mathcal{E}$ .
  Output:
     $\mathcal{R}$  the learned rule set

```

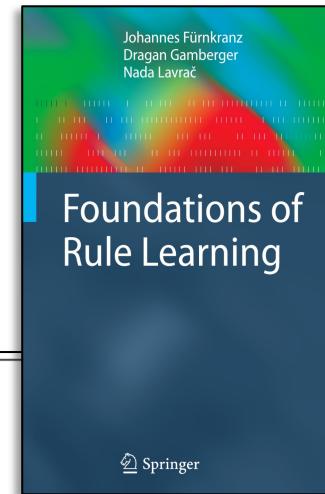
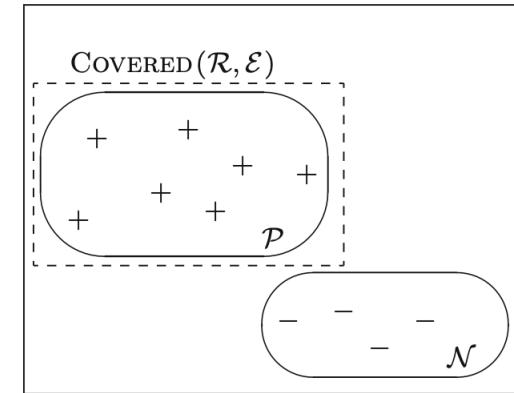
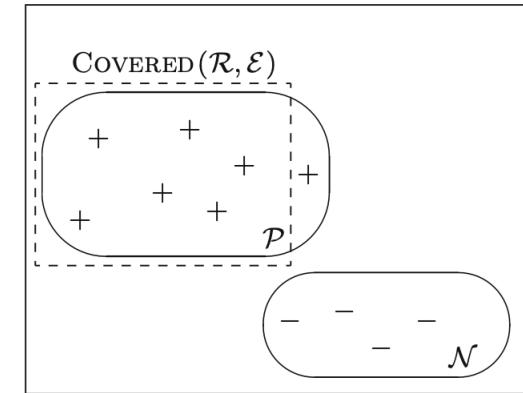


Fig. 2.10 Constructing a set of rules in a multiclass learning setting

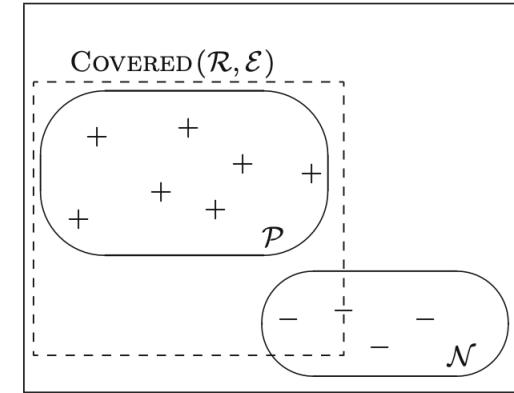
\mathcal{R} : complete, consistent



\mathcal{R} : incomplete, consistent



\mathcal{R} : complete, inconsistent



\mathcal{R} : incomplete, inconsistent

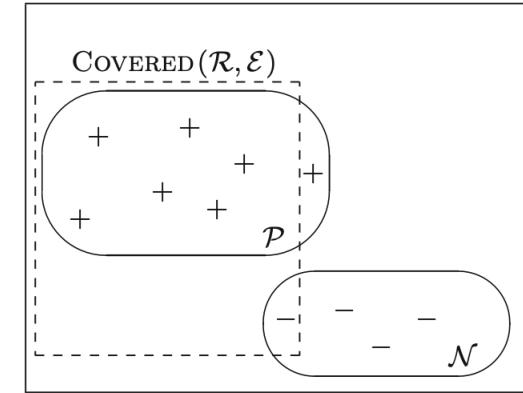


Fig. 2.3 Completeness and consistency of a hypothesis (rule set \mathcal{R})



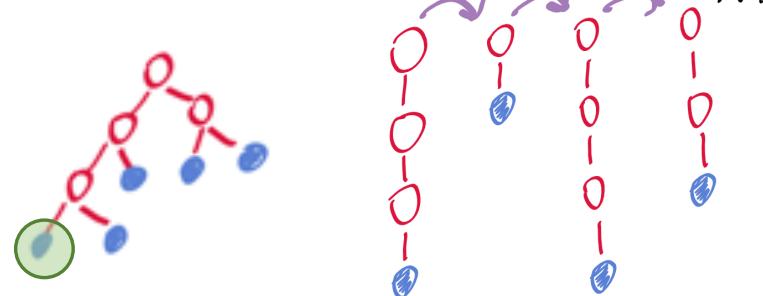
Rule Generation

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

$\mathcal{I} = \{1, \dots, n\}$, sample indices

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

multi-class



$$\min_{w_R \geq 0, R \in \mathcal{R}} \underbrace{\sum_{i \in \mathcal{I}} \ell \left(\sum_{R \in \mathcal{R}} \hat{y}_{iR} w_R, y_i \right)}_{\ell_{\mathcal{D}}(R)} + \lambda \underbrace{\sum_{R \in \mathcal{R}} c_R w_R}_{\lambda \Omega(R)}$$

(Bin)OCT, FSOT, BDR,...

Integer Programming

- do not scale well
- mostly for binary classification

$$\min_{R \in \mathcal{R}} \ell_{\mathcal{D}}(R) + \lambda \Omega(R)$$

total loss
model complexity
all possible rules

prediction of rule for a sample

weight of the rule

cost of a rule

RUG
Linear Programming



Scalable Rule Generation for Learning



Interpretation

~450K samples

XAI models for Loan Prediction - Model Overview

6 rules using 11 of 42 features

confusion matrix

Accuracy	0.9759
F1 score	0.8804
Precision score	0.8811
Recall score	

	Predicted 0	Predicted 1
Actual - 0	70148	942
Actual - 1	955	6985

global interpretation

XAI models for Loan Prediction

Customer Id : 182977 ▾

Model View

Rule 0

Months since issue date
≥ 92.5

→ Months since last payment date
≤ 67.5

→ W = 1.0

P

Koç Digital - XAI models for Loan Prediction

3 rules using 5 of 42 features

Total payment

8436.7

P

Total recurring int.

1936.7

P

Months since issue date

62

P

Last payment amount

1936.7

P

Months since last payment date

62

P

Outstanding principal

0.0

Annual income

46.000

Total current balance

113.000

Employment length

11

Term

36

Home ownership

yes

Purpose - small business

yes

Purpose - moving

no

Purpose - dept consolidation

yes

Purpose - home improvement

no

* Please hover over class symbols to see their explanations

Rule Breakdown

Rule 0

Total payment
≥ 7134.1

↓

P

W = 1.0

P

Rule 1

Months since issue date
≥ 92.5

↓

P

W = 1.0

P

Rule 4

Last payment amount
≤ 1257.8

↓

P

W = 0.5

D

Final Prediction



Performing is the final prediction of this specific observation (Customer id: 182977). According to the majority voting among the results of the 3 applied rules, it received **80%** of the total vote (**2 out of 2.5**) and is chosen as the final prediction. This explanation is generated on 2022-06-07 18:03:05

P - Performing

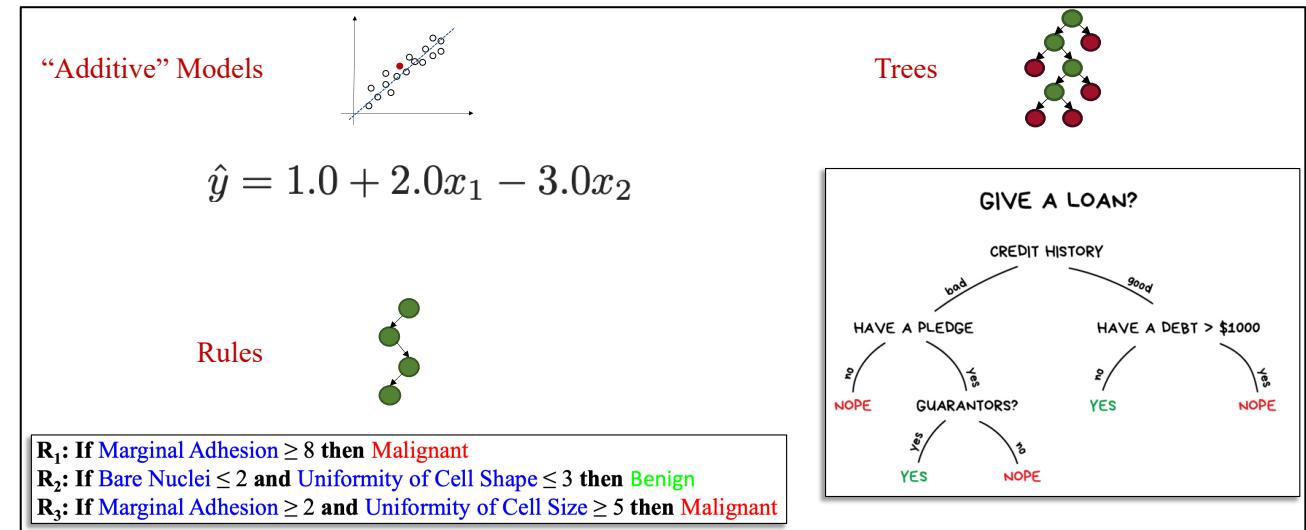
D - Default

local interpretation

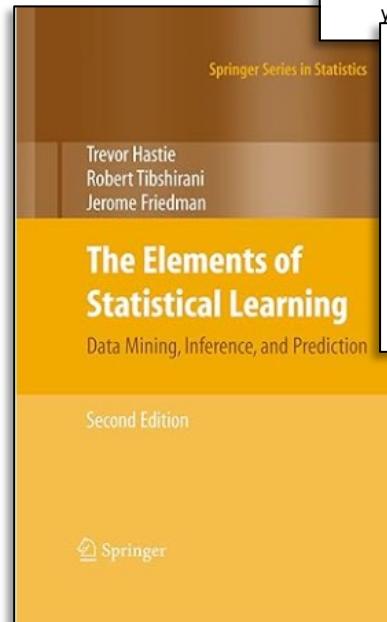


Big Picture

- Linear Models
- Generalized Additive Models
- Explainable Boosting Machines
- Decision Trees
- Rule-based Approaches



Reading Material



Intelligible Models for Classification and Regression

Yin Lou
Dept. of Computer Science
Cornell University
yinlou@cs.cornell.edu

Rich Caruana
Microsoft Research
Microsoft Corporation
rcaruana@microsoft.com

Johannes Gehrke
Dept. of Computer Science
Cornell University
johannes@cs.cornell.edu

Accurate Intelligible Models with Pairwise Interactions

Yin Lou
Dept. of Computer Science
Cornell University
yinlou@cs.cornell.edu

Rich Caruana
Microsoft Research
Microsoft Corporation
rcaruana@microsoft.com
Giles Hooker
Dept. of Statistical Science
Cornell University
giles.hooker@cornell.edu

Johannes Gehrke
Dept. of Computer Science
Cornell University
johannes@cs.cornell.edu

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

Optimal classification trees

Dimitris Bertsimas¹ · Jack Dunn²

Machine Learning, 11, 63–91 (1993)
© 1993 Kluwer Academic Publishers, Boston. Manufactured in The Netherlands.

Very Simple Classification Rules Perform Well on Most Commonly Used Datasets

ROBERT C. HOLTE
Computer Science Department, University of Ottawa, Ottawa, Canada K1N 6N5
[HOLTc@cs.uottawa.ca](mailto:holtc@cs.uottawa.ca)

Editor: Bruce Porter

Abstract. This article reports an empirical investigation of the accuracy of rules that classify examples on the basis of a single attribute. On most datasets studied, the best of these very simple rules is as accurate as the rules induced by the majority of machine learning systems. The article explores the implications of this finding for machine learning research and applications.

Keywords: empirical learning, accuracy-complexity tradeoff, pruning, ID3

Optimal Sparse Decision Trees

Xiyang Hu¹, Cynthia Rudin², Margo Seltzer^{3*}

¹Carnegie Mellon University, xiyanghu@cmu.edu

²Duke University, cynthia@cs.duke.edu

³The University of British Columbia, mseltzer@cs.ubc.ca

Rule Generation for Classification: Scalability, Interpretability, and Fairness

Adia Lumadjieng*, Tabea Röber**, M. Hakan Akyüz*, §. İlker Birbil**

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

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

