

Discussion 11

Today:

1. TF-IDF comments from discussion 10
2. CART
3. Random Forest

1.) TF-IDF

• Bag-of-words: x_i is a vector whose entry j is the number of times word j appears in doc. i .

• TF-IDF

$$tf(x_{ij}) = \log(1 + x_{ij})$$

$$idf(j) = \log \frac{N}{1 + \sum_{i=1}^N \mathbb{I}(x_{ij} > 0)}$$

$$\text{tf-idf}(i) = \left[\text{tf}(x_{ij}) \cdot \text{idf}(j) \right]_{j=1}^M, \quad M \text{ total \# of words}$$

Kernel: cosine similarity

$$K(x_1, x_2) = \frac{x_1^T x_2}{\|x_1\|_2 \|x_2\|_2} \leadsto \text{bag-of-words}$$

$$K(x_1, x_2) = \frac{\phi(x_1) \cdot \phi(x_2)}{\|\phi(x_1)\|_2 \|\phi(x_2)\|_2}$$

$$\phi(x_i) = \text{tf-idf}(i)$$

$\leadsto \parallel \leadsto$

CART

• Growing the tree

$$(j^*, t^*) = \arg \min_{j \in \{1, \dots, D\}} \min_{t \in \mathbb{T}_j} \left[\text{cost}(\text{left node}) + \text{cost}(\text{right node}) \right]$$

Best split on best feature.

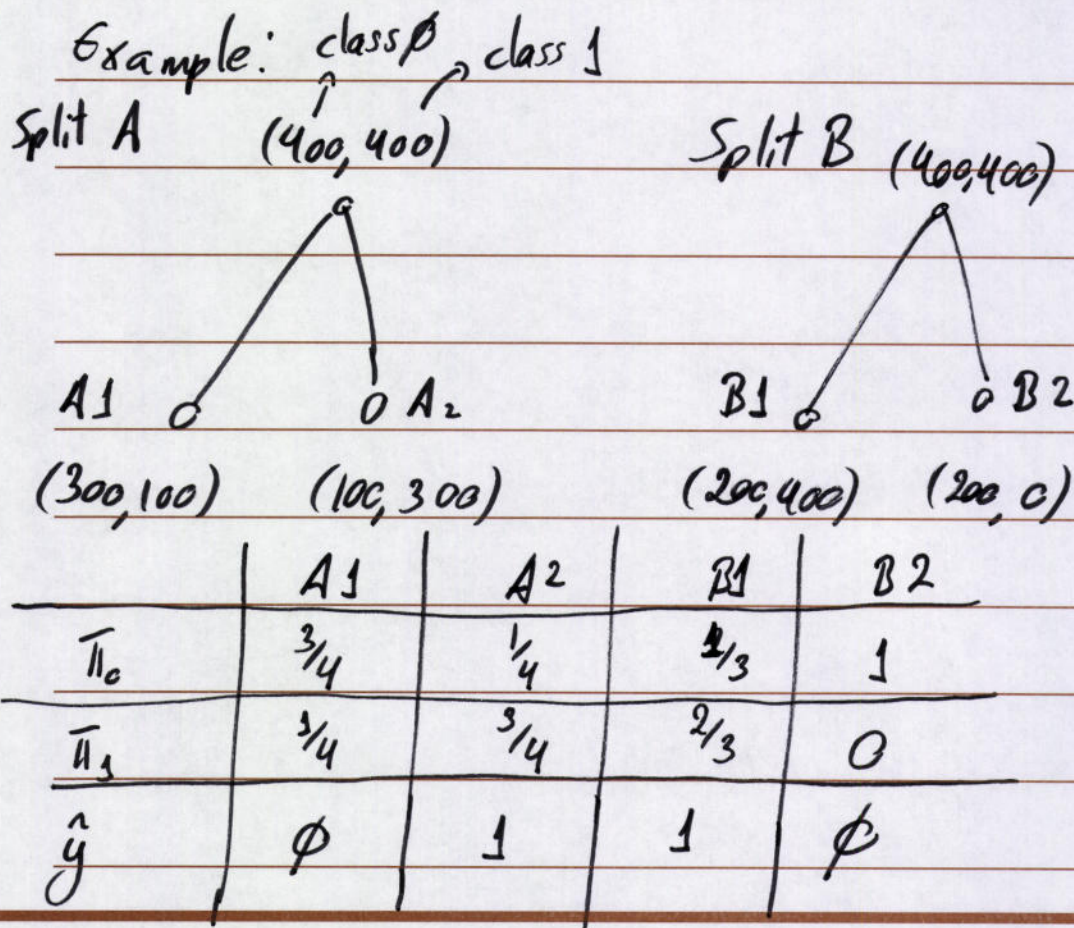
Since after the split regions are indep. of each other (greedy algo.), the order in which regions are analyzed is (usually) not important

3

. Some stopping criteria from sklearn (attention to default values)

- max_depth, default = none
- min_sample_split, default = 2
- min_samples_leaf, default = 1
- max_leaf_nodes, default = none
- min_impurity_decrease, default = 0

. Cost choices: accuracy, entropy, gini index



Misclassification rate:

$$e_m = 1 - \pi \hat{y}$$

Gini index:

$$e_g = 2\pi \hat{y}(1 - \pi \hat{y})$$

	A1	A2	B1	B2
e_m	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{3}$	0
e_g	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{4}{9}$	0

$$\text{cost}(\text{split}) = \frac{e_1 \cdot N_1 + e_2 \cdot N_2}{N}$$

~~cost(split A)~~

Misclassification:

$$\text{cost}(\text{split A}) = \frac{\frac{1}{4} \cdot 400 + \frac{1}{4} \cdot 400}{800} = \frac{1}{4}$$

$$\text{cost}(\text{split B}) = \frac{\frac{1}{3} \cdot 600 + \phi}{800} = \frac{1}{4}$$

Gini index:

$$\text{cost}(\text{split A}) = \frac{\frac{3}{8} \cdot 400 + \frac{3}{8} \cdot 400}{800} = \frac{3}{8}$$

$$\text{cost}(\text{split B}) = \frac{\frac{4}{9} \cdot 600 + \phi}{800} = \frac{1}{3}$$

cost(B) < cost(A)

Random Forest

1. Bagging: grow different trees based on subsamples of your dataset with replacement

2. Random Forest: not only random subsets ^{of data} but also random subset of features at each node.

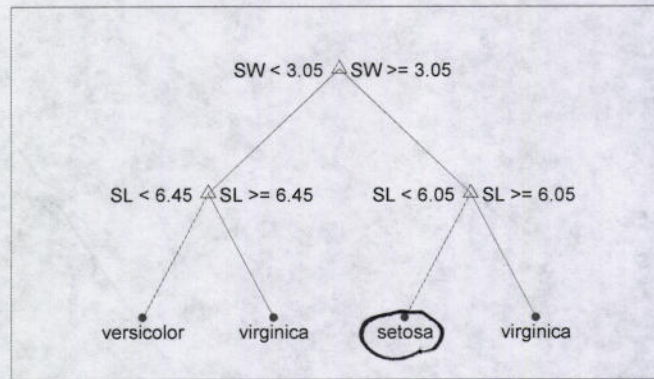
Some parameters:

- all from CART

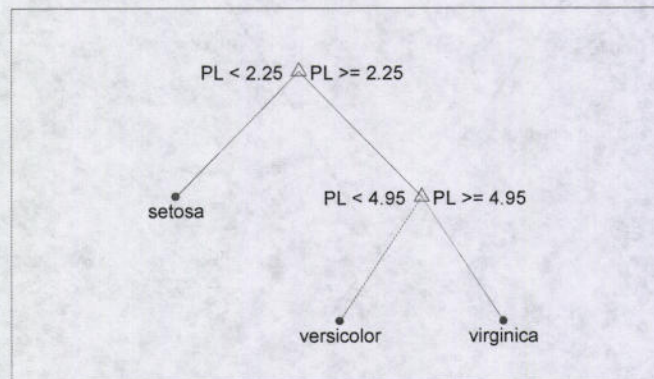
- n_estimators, default=10: # of trees

- max_features, default=auto: # of considered features at each split

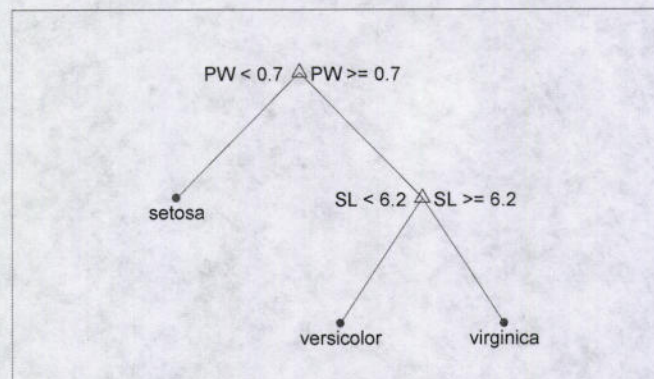
$$X = \begin{bmatrix} 5 \\ 3 \\ 2 \\ 1 \end{bmatrix} \begin{matrix} SW \\ SL \\ PW \\ PL \end{matrix}$$



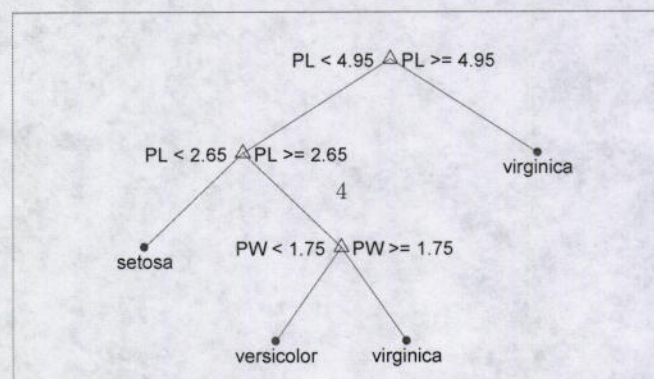
(a) Tree 1



(b) Tree 2



(c) Tree 3



(d) Tree 4

Figure 1: Random Forest simple example