

Module 8: Tree-based methods

05.03.2018
M8W1

Regression trees $(x_i, y_i) \quad i=1, \dots, n$
↑ p -dim predictor
↑ univariate continuous random variable

$$Y = f(x) + \epsilon$$

1) Divide predictor space into J non-overlapping regions,
 R_1, \dots, R_J

2) Prediction in R_j is $\hat{y}_j = \text{mean of training obs that fall into } R_j$

How to decide on R_1, \dots, R_J ?

Recursive binary splitting:

at the current node split into $R_1(j, s) = \{x \mid x_j < s\}$

$$R_2(j, s) = \{x \mid x_j \geq s\}$$

and choose j and s to

$$\text{minimize} \quad \sum_{i: x_i \in R_1} (y_i - \hat{y}_{R_1})^2 + \sum_{i: x_i \in R_2} (y_i - \hat{y}_{R_2})^2$$

3zone ← match R 's and tree nodes.

Classification trees

K classes

$\{1, \dots, K\}$

1) Prediction in R_j : majority rule or chosen cut-off on

$$\hat{p}_{jk} = \frac{1}{n_j} \sum_{i: x_i \in R_j} I(y_i = k)$$

region class

2) Splitting criterion: minimize misclassification rate: $1 - \max_k \hat{p}_{jk}$

When to split? \rightarrow too crude

Q: are the child nodes on average "purer" than their parents? For region j

If one $p = 1 \rightarrow$ very pure: 0 | $-\sum_k p_k \log_j p_k$ cross entropy
all $p = \frac{1}{K} \rightarrow$ not pure: large | $\sum p_{jk}(1-p_{jk})$ Gini
 \uparrow expected error rate if label is

statisticians prefer Gini \rightarrow

chosen at random from the class distribution at the node