

Classification and Regression Trees

Mihaela van der Schaar

Department of Engineering Science
University of Oxford

March 1, 2017

Algorithm for Regression Trees

- Start with $\mathcal{R}_1 = \mathbb{R}^d$
- For each feature $j = 1, \dots, d$, for each value $v \in \mathbb{R}$ that we can split on:

- Split the data set:

$$I_{<} = \{i : x_{ij} < v\} \text{ and } I_{>} = \{i : x_{ij} \geq v\}$$

- Estimate parameters:

$$\beta_{<} = \frac{\sum_{i \in I_{<}} y_i}{|I_{<}|} \text{ and } \beta_{>} = \frac{\sum_{i \in I_{>}} y_i}{|I_{>}|}$$

- Quality of split is measured by the squared loss:

$$\sum_{i \in I_{<}} (y_i - \beta_{<})^2 + \sum_{i \in I_{>}} (y_i - \beta_{>})^2$$

- Choose split with minimal loss.
- Recurse on both children, with $(x_i, y_i)_{i \in I_{<}}$ and $(x_i, y_i)_{i \in I_{>}}$.

