

Regression trees

Model representation and loss

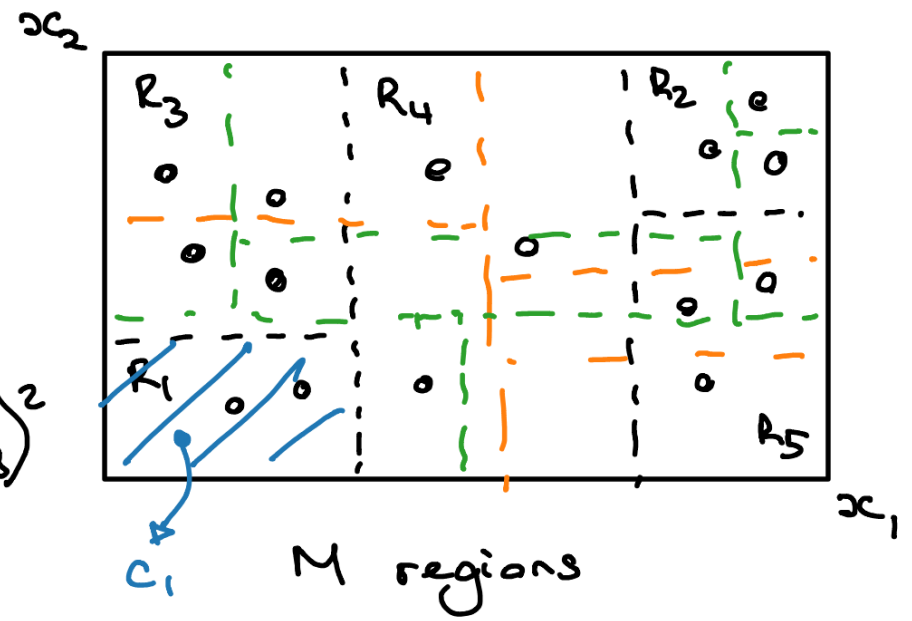
Herman Kamper

<http://www.kamperh.com/>

Regression tree model

Model: $f(\underline{x}; \underline{\theta}) = \sum_{m=1}^M c_m I\{\underline{x} \in R_m\}$

Loss: $J(\underline{\theta}) = \sum_{i=1}^n (y^{(i)} - f(\underline{x}^{(i)}))^2$
 $= \sum_{i=1}^n (y^{(i)} - \sum_{m=1}^M c_m I\{\underline{x}^{(i)} \in R_m\})^2$



Parameters: $\left. \begin{matrix} R_1, \dots, R_M \\ c_1, \dots, c_M \end{matrix} \right\} \underline{\theta}$

If we knew regions (but not the c 's):

$$\Rightarrow c_m = \frac{1}{N_m} \sum_{i: \underline{x}^{(i)} \in R_m} y^{(i)}$$

How do we learn $\underline{\theta}$?

training items in R_m

Regression trees

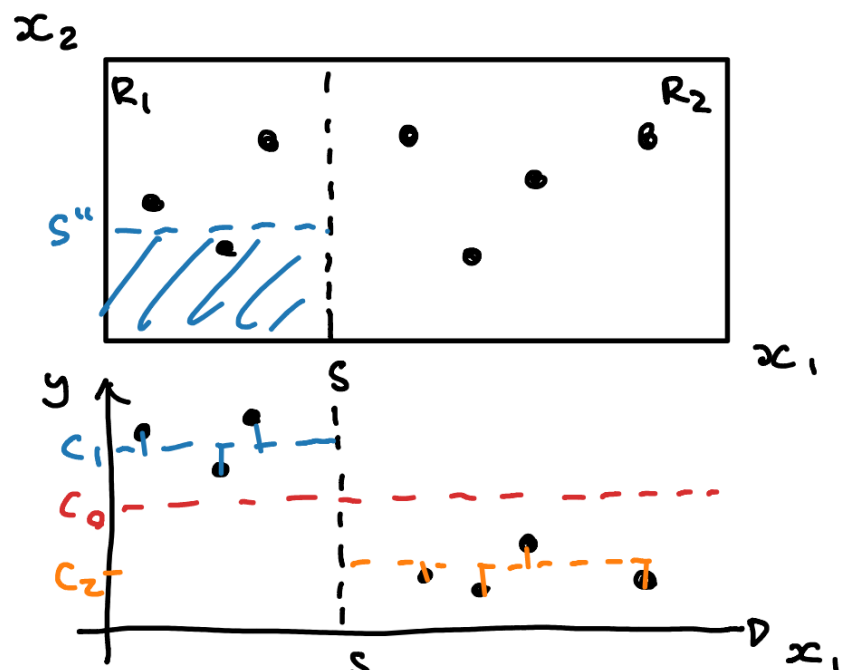
Tree building algorithm

Herman Kamper

<http://www.kamperh.com/>

Tree growing algorithm: ↙ top-down everything in

1. Start at top of tree (one region)
2. for each leaf node (region):
 for each feature x_j and split point s :
 Calculate reduction in loss if we split there
↙ greedy
3. Choose best (j, s) combination to split;
 create new child nodes (regions)
4. Repeat from (2) until stop condition is met
↖ recursive



Example:

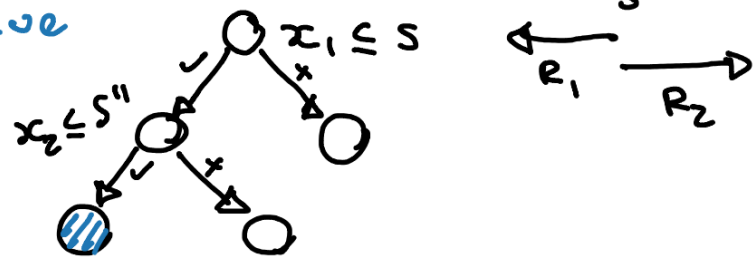
Before split: $J = \sum_{n=1}^7 (y^{(n)} - c_0)^2$

Consider splitting x_1 at s :

$$R_1 = \{ \underline{x} \mid x_1 \leq s \}$$

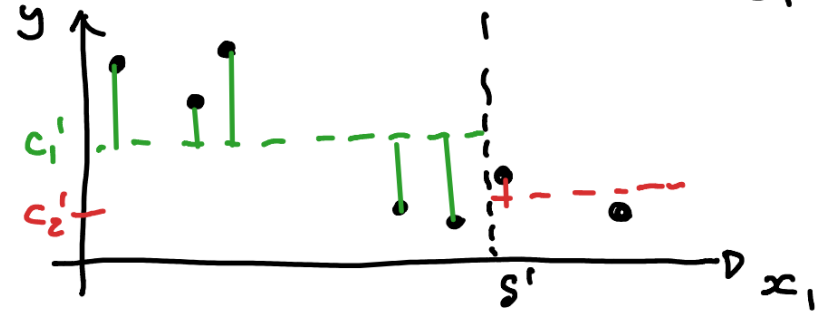
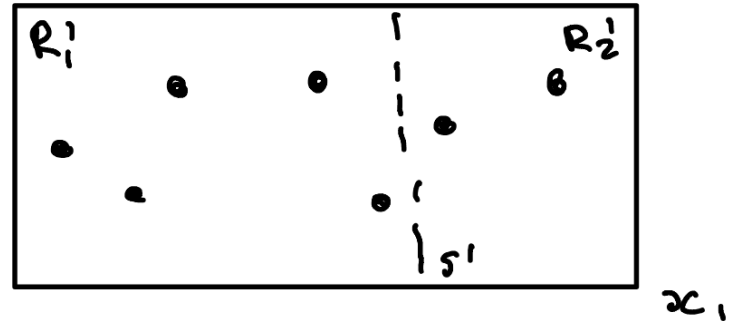
$$R_2 = \{ \underline{x} \mid x_1 > s \}$$

If we split: $J = \sum_{i: \underline{x}^{(i)} \in R_1} (y^{(i)} - c_1)^2 + \sum_{i: \underline{x}^{(i)} \in R_2} (y^{(i)} - c_2)^2$



Tree growing algorithm:

1. Start at top of tree (one region)
2. for each leaf node (region):
for each feature x_j and split point s :
Calculate reduction in loss if we split there
3. Choose best (j, s) combination to split;
create new child nodes (regions)
4. Repeat from (2) until stop condition is met

 x_2 

$$\min_{j, s} \left\{ \sum_{i: \underline{x}^{(i)} \in R_1} (y^{(i)} - c_1)^2 + \sum_{i: \underline{x}^{(i)} \in R_2} (y^{(i)} - c_2)^2 \right\}$$

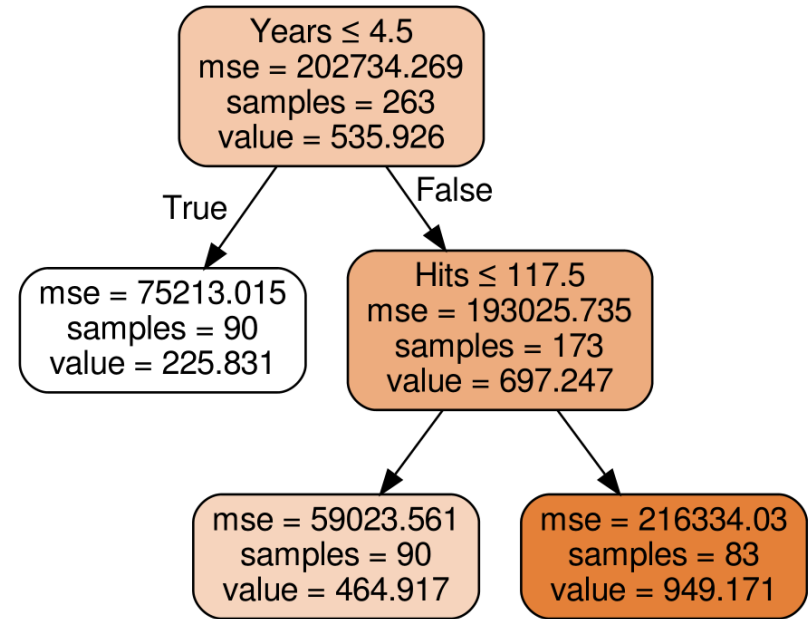
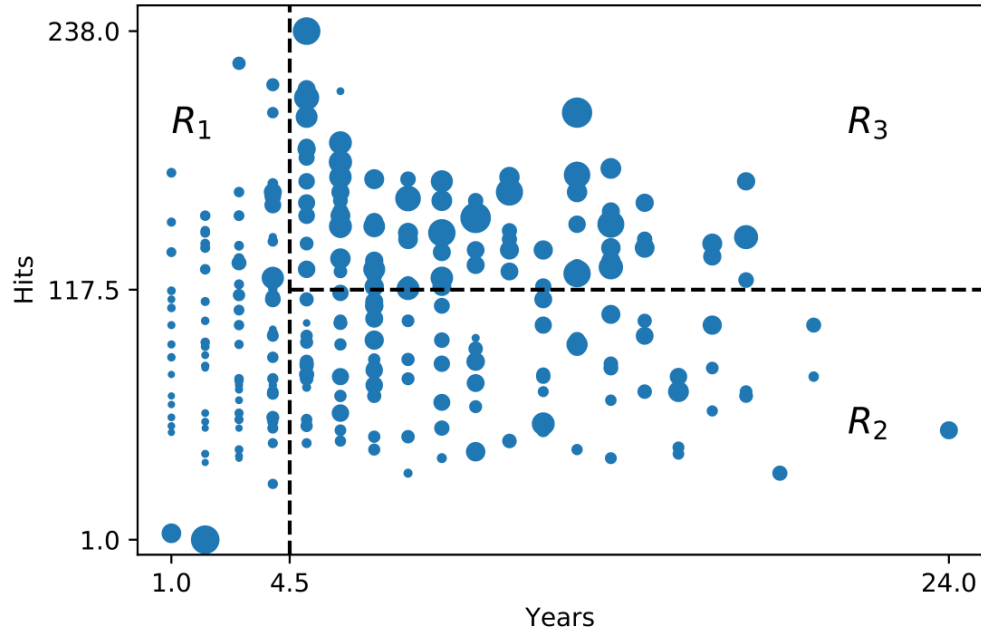
Regression trees

Regression trees in practice and tree pruning

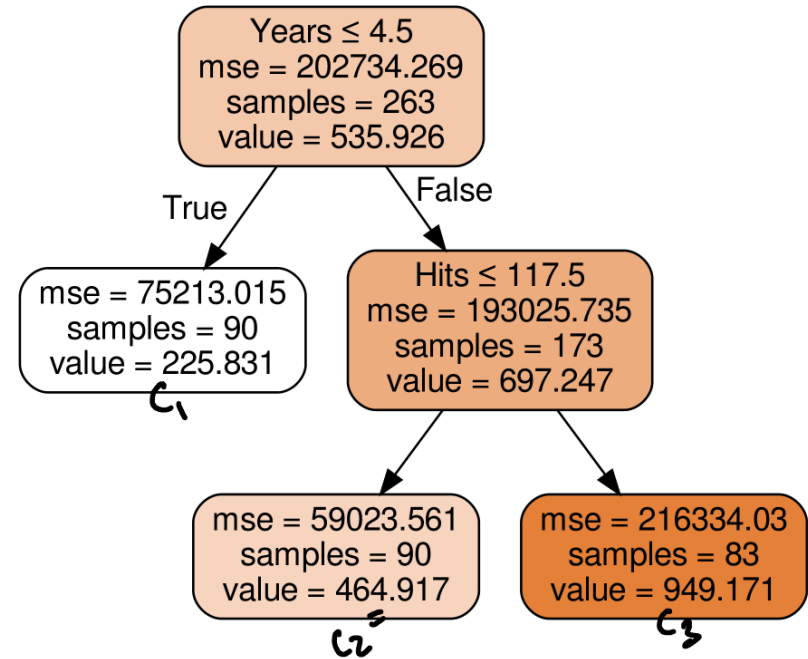
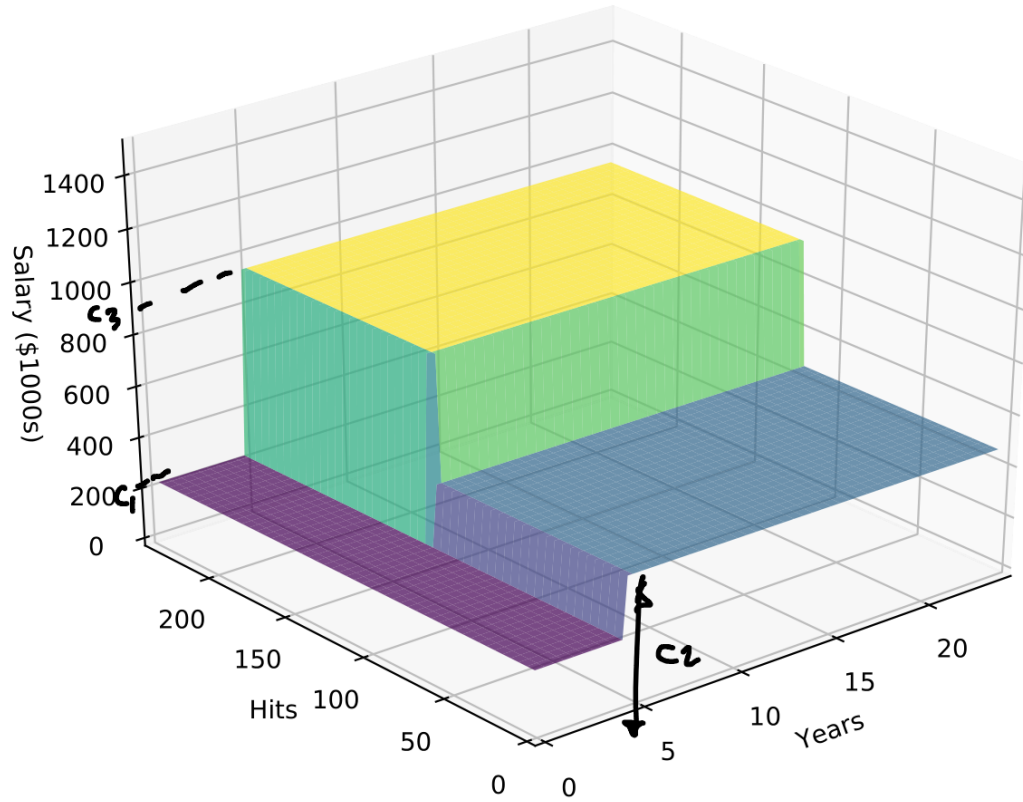
Herman Kamper

<http://www.kamperh.com/>

Regression tree on hitters data

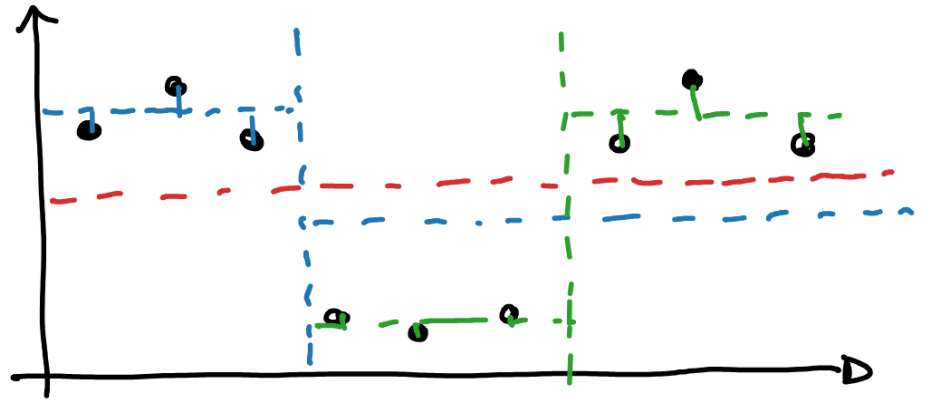
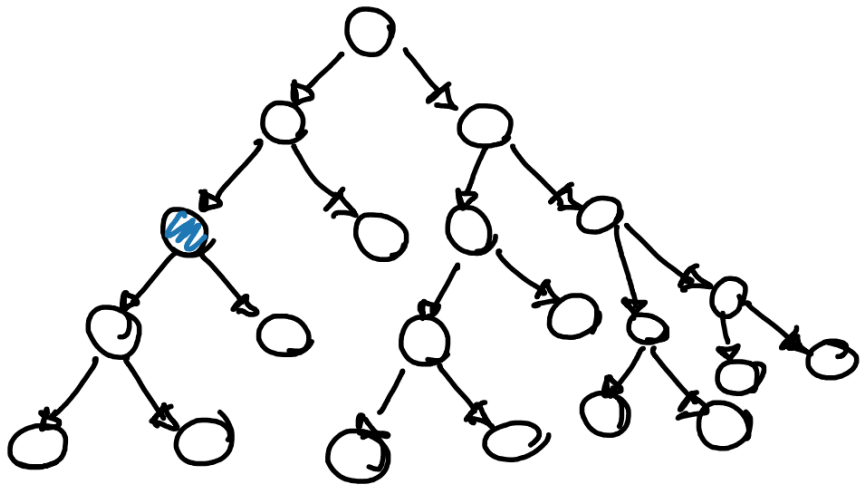


Regression tree on hitters data



Regression trees in practice:

- Can easily overfit
- How do we regularise?
 - Can stop if gain is small
 - But this can be short-sighted
- Can use tree pruning



$$J = \sum_{n=1}^N \sum_{i: x^{(i)} \in R_m} (y^{(i)} - c_m)^2$$