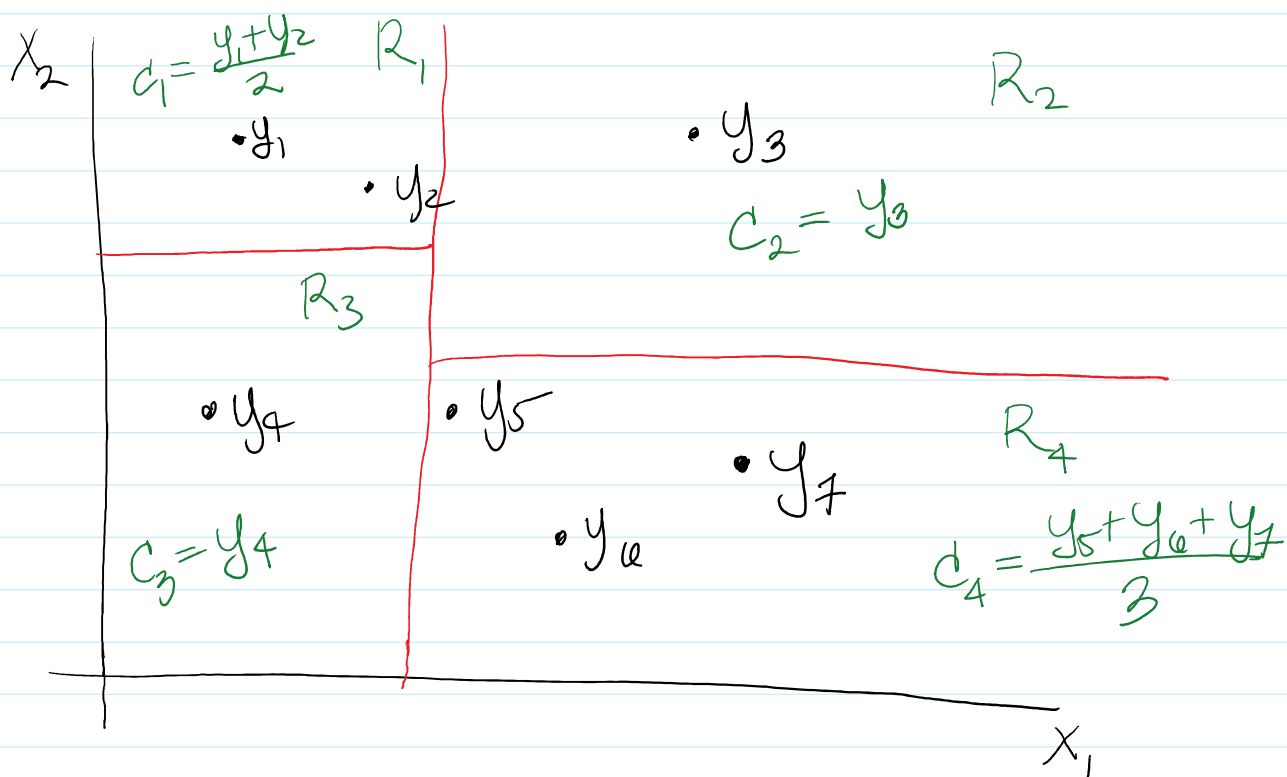


Back to Supervised Learning

Regression Trees

Basic idea

- ① break up feature space into a collection of rectangles - by successively splitting dimensions
- ② Fit a simple model to each rectangle (predict avg. value of training data in each rectangle)



$$\hat{y} = \hat{f}(x) = c_i \text{ where } x \in R_i$$

\xrightarrow{M}

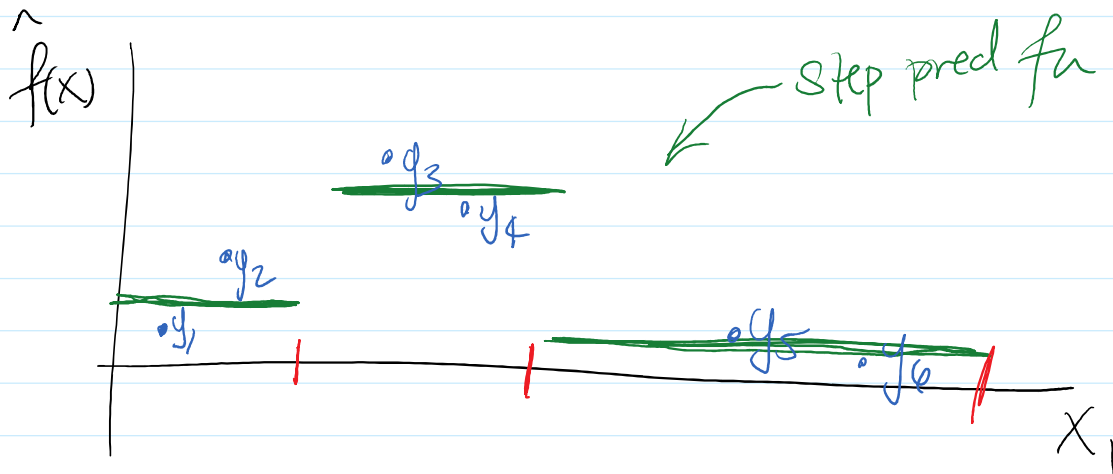
$$\hat{y} = f(x) = c_i \text{ where } x \in R_i$$

$$= \sum_{i=1}^M c_i \mathbb{I}(x \in R_i)$$

rectangles $\rightarrow M$

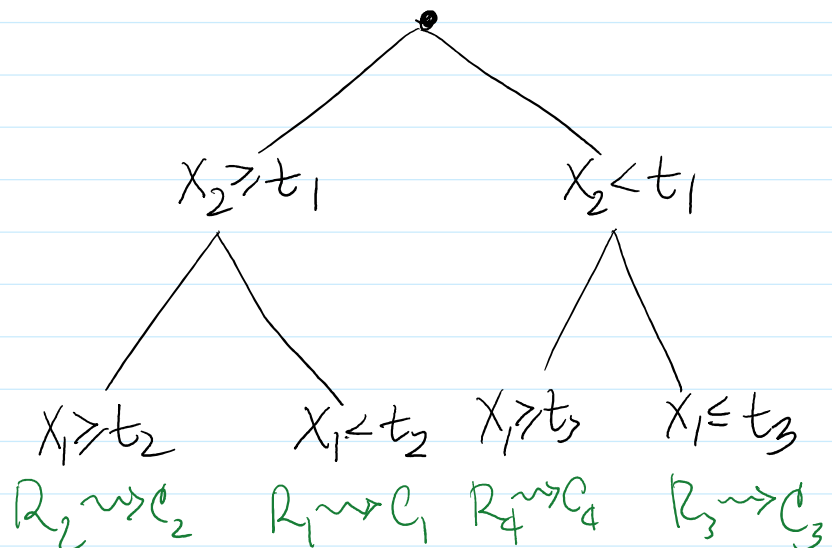
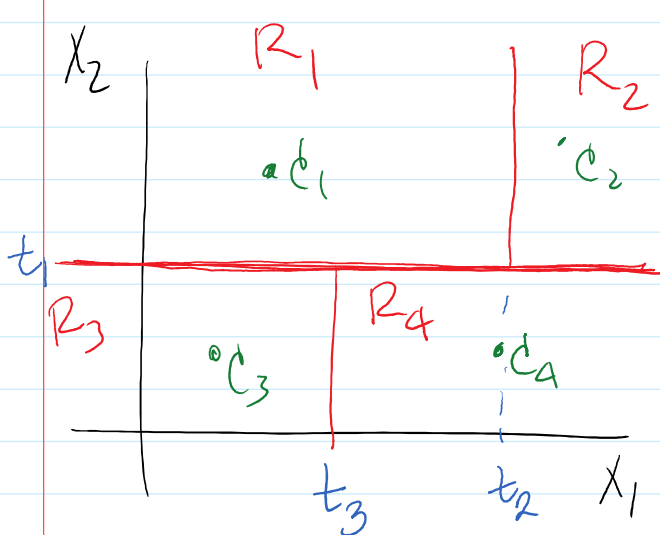
arg val. of my training y s in R_i

$\begin{cases} 1 & x \in R_i \\ 0 & \text{else} \end{cases}$



Why do we call this a tree?

Can represent as a binary decision tree.



Goal: Build a good tree.

Need to decide

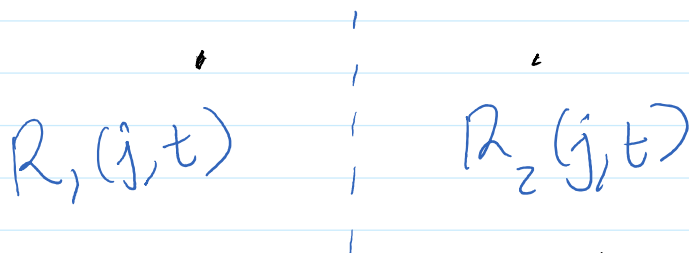
- (1) which variable to split?
- (2) where do I split this variable?
- (3) when do I stop (how many splits?)

Optimal Approach:

Consider all possible trees and choose the best.
Not computationally tractable.

Actually use is some (somewhat) greedy approach.

Define $R_1(j, t)$ and $R_2(j, t)$ to be the half-spaces formed by splitting on var X_j at point t .





furthermore let

$$\rightarrow \underline{RSS(j, t)} = RSS(R_1(j, t)) + RSS(R_2(j, t))$$

$$= \sum_{i: x_i \in R_1(j, t)} (y_i - c_1)^2 + \sum_{i: x_i \in R_2(j, t)} (y_i - c_2)^2$$

vals. predicted
in Right and Left

turns at to minimize,

$$\begin{cases} c_1 = \text{avg of } y_s \text{ in } R_1 \\ c_2 = \text{avg of } y_s \text{ in } R_2 \end{cases}$$

Algo: Need to choose \hat{j}, \hat{t} .

- (1) For each var. do a search over possible
 t and calc $RSS(j, t)$
- (2) choose \hat{j}, \hat{t} to minimize $RSS(j, t)$.

and make that split

(3) Recursively do this for split data.
(each side)

when do we stop?

→ too many splits (very flexible method)
worry: overfitting

→ too few splits (inflexible method)
worry: underfitting

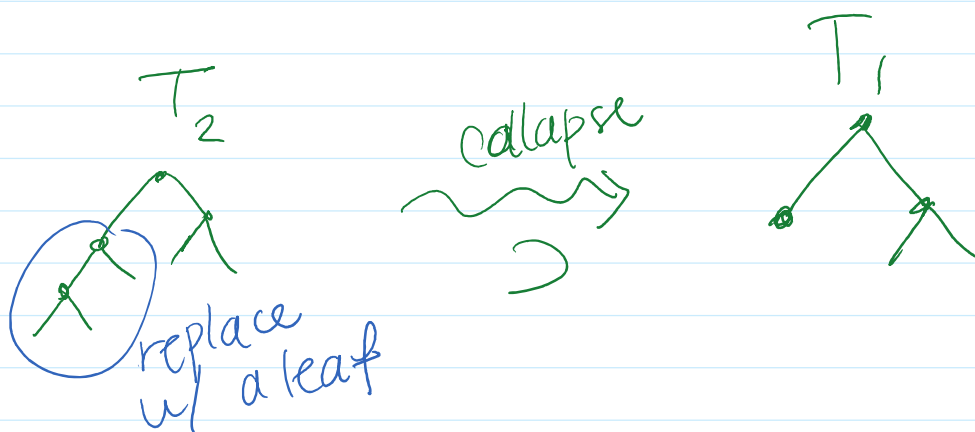
Bad strategy very greedy approach.

set some threshold and split until our
RSS decreases by less than threshold.

problem: a bad split might lead to a
better split down the road

Better: (1) grow a really large tree (overfit)
(2) reduce its size (pruning)

Aside: a tree T_1 is called a subtree of T_2 if I can get T_1 by "collapsing" part of T_2 .



$$T_1 \subset T_2.$$

Do this strategy w/ cost-complexity pruning

$$C_\alpha(T) = \text{RSS}(T) + \alpha |T|$$

Annotations:

- $\alpha \geq 0$: tuning parameter
- $|T|$: size of tree = # of leaf nodes
- T : tree T

Basically: given some $\alpha \geq 0$

- (1) grow a large tree
- (2) search over subtrees and choose (prune) subtree that minimizes $C_\alpha(T)$

note: $\alpha = 0 \Rightarrow$ choose as large a tree as poss.

$\alpha = \infty \Rightarrow$ choose as small as possible

Fact: $\alpha_1 \leq \alpha_2$ then the tree that minimizes $C_{\alpha_1}(T)$ — called T_1 — and corresp T_2

minimizes $C_{\alpha_2}(T)$ then

$$T_2 \subset T_1.$$

So optimal trees over values of α are nested.

If I have a seq

$$\alpha_1 \leq \alpha_2 \leq \alpha_3 \leq \dots \leq \alpha_N$$

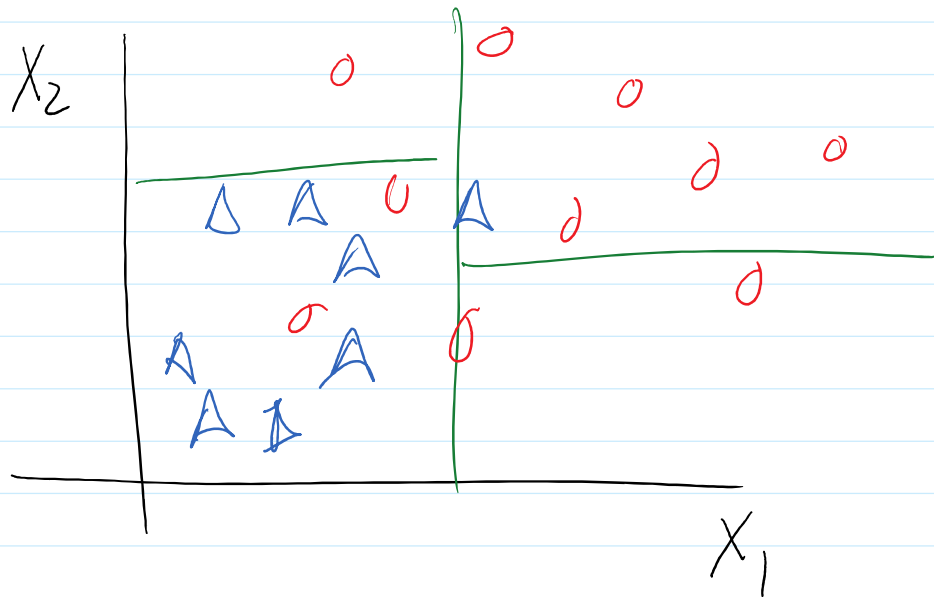
$$T_1 \supset T_2 \supset T_3 \supset \dots \supset T_N$$

↑ nested seq of optimal trees over values of α

Can easily compute optimal CC trees over a seq of α s.

Classification Trees

Classification trees



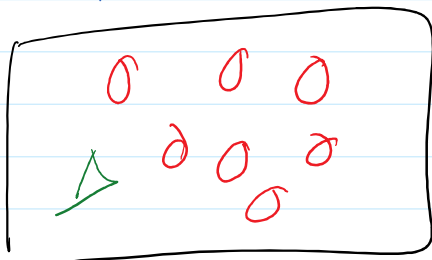
predict the majority class in each rectangle

what makes a good split?

Regression tree: RSS (reduce RSS)

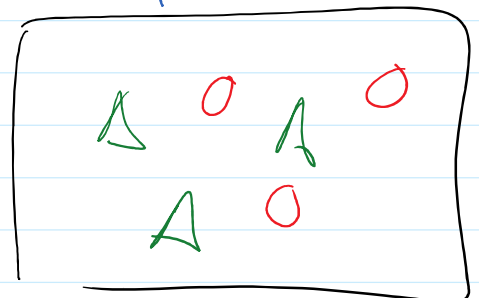
Classification Tree: node purity/impurity
(reduce node impurity)

Ex. pure node



almost all red 0

impure node



50/50 split

almost all red 0

50/50 split

Node Impurity Measures

① mis. class. rate : $1 - \hat{p}_k^1$ where \hat{p}_k^1 is pct. in node of class k
is class k

very pure ≈ 0
impure $\approx 1/K$

$K = \# \text{ classes}$

② Gini - Index

$$\sum_{k=1}^K \hat{p}_k (1 - \hat{p}_k)$$

pure ≈ 0
impure $\approx 1/K$

③ Entropy

$$-\sum_{k=1}^K \hat{p}_k \log \hat{p}_k$$

pure ≈ 0
impure $\approx 1/K$

two class

$K=2$

$1/K$

