

Lecture 22: trees

CARTs - classification and regression trees

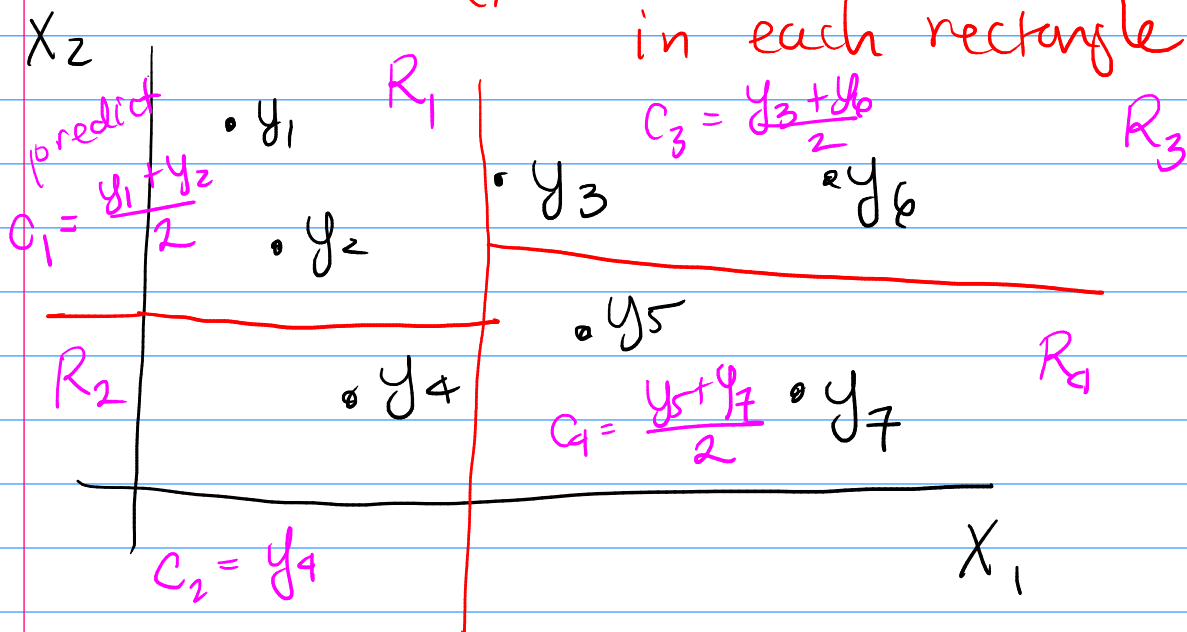
Regression Trees

Basic idea:

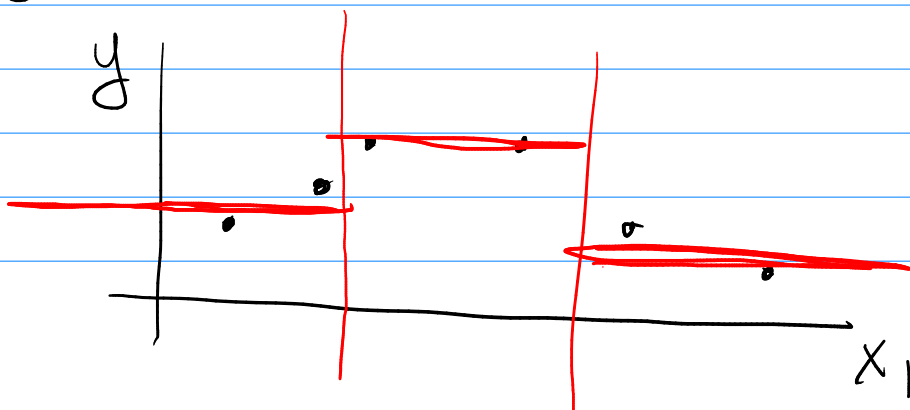
(1) break up X -space into rectangles

(2) fit a simple model on each rectangle

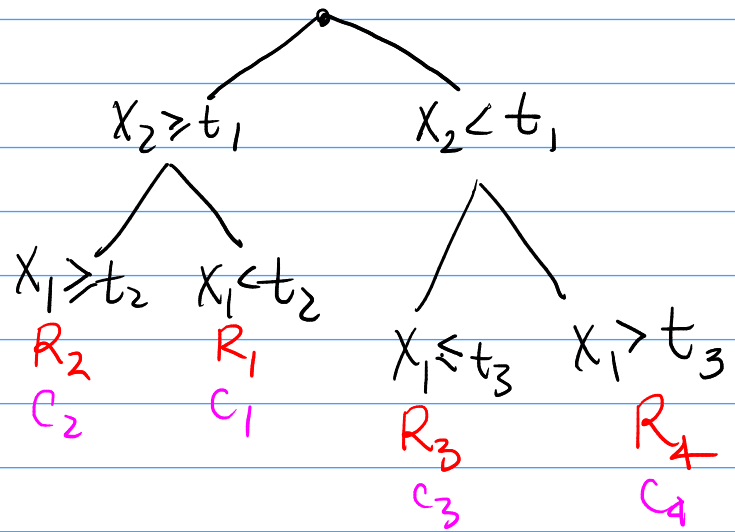
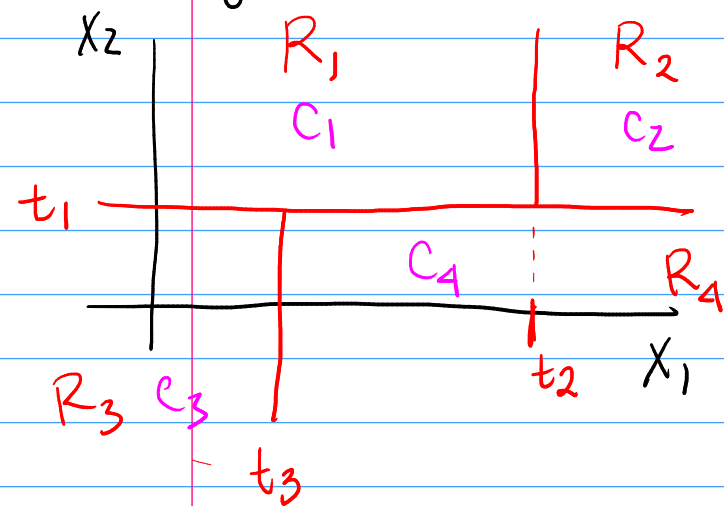
(predict the mean of training y s in each rectangle)



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



Why called a tree? Can represent as a decision tree.



Goal: build a good tree

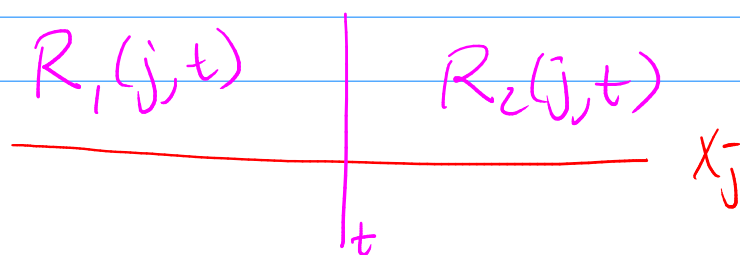
Need to decide

- ① which variable I split on
- ② where I split this variable
- ③ when do I stop splitting.

Optimal approach: try all trees - not computationally tractable

So use some greedy approach.

Define $R_1(j, t)$ and $R_2(j, t)$ to be the half-spaces if we split variable j at point t



We can define the RSS if we make this split as

$$\begin{aligned} \text{RSS}(j, t) &= \text{RSS}(R_1(j, t)) + \text{RSS}(R_2(j, t)) \\ &= \sum_{i \text{ in } R_1(j, t)} (y_i - c_1)^2 + \sum_{i \text{ in } R_2(j, t)} (y_i - c_2)^2 \end{aligned}$$

\uparrow \uparrow
 $i \text{ in } R_1(j, t)$ mean of y s in $R_1(j, t)$

\uparrow similar for $R_2(j, t)$

Algorithm to choose j and t

- ① For each var j search over possible t and calc. $\text{RSS}(j, t)$
- ② Choose j and t that minimize $\text{RSS}(j, t)$.
- ③ recursively do this for each half space.

When do we stop?

- too many splits danger is over fitting
- too few danger is under fitting

Bad strategy very greedy approach and split until RSS falls below some threshold

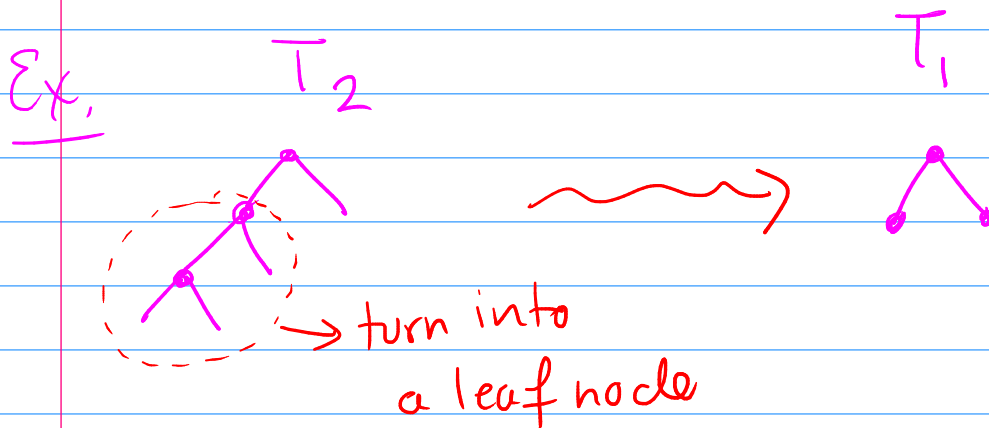
problem: a bad split might lead to

an even better split later

better approach:

- ① grow a really large tree (overfit)
 - ② reduce its size by 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



In particular use "cost-complexity pruning"

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

$\alpha > 0$

↑ size of tree
i.e. # of
leaf nodes

Given some $\alpha \geq 0$

① grow a large tree

② search over sub-trees

and choose sub-tree that minimizes
 $C_\alpha(T)$

If $\alpha = 0 \Rightarrow$ choose largest subtree i.e.
remove nothing

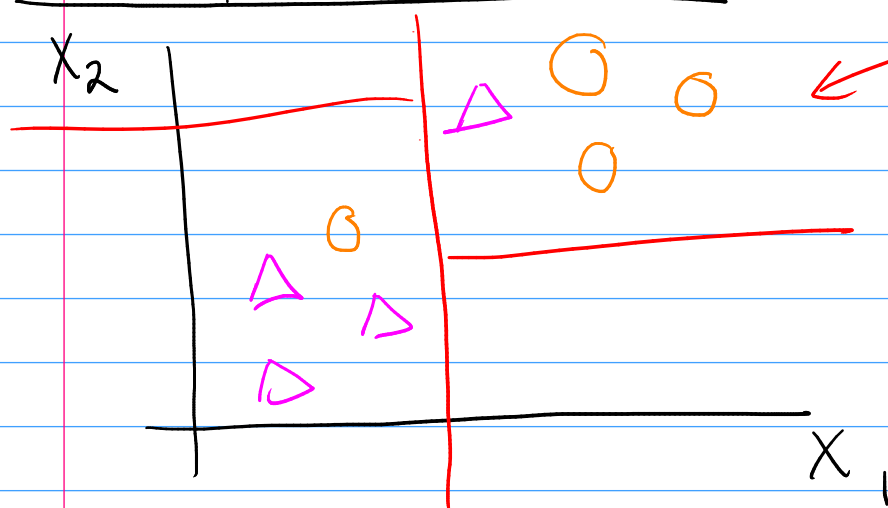
$\alpha = \infty \Rightarrow$ choose smallest subtree i.e.
remove everything and
don't split

Turns out that if $\alpha_1 \leq \alpha_2$ then the
optimal tree for α_1 (call it T_1)
and optimal tree for α_2 (call it T_2)
are related so that

$T_2 \subset T_1$
↑ sub-tree.

So as α increase we get a nested
series of trees.

Classification Trees



predict the majority (plurality) class in each rectangle

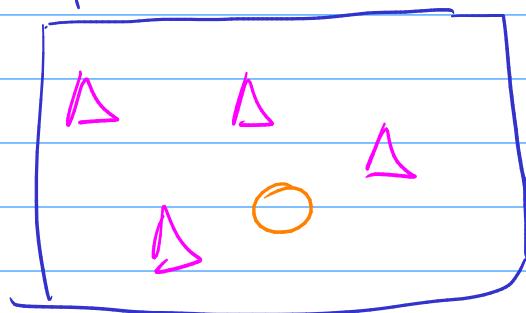
What makes a good split?

Regression Trees: RSS

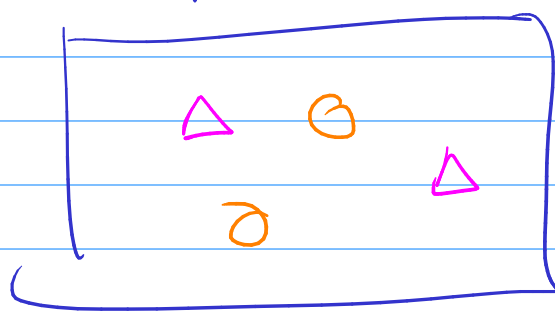
Classification Trees: node (rectangle) impurity

Ex,

pure node



impure node



Now we split to reduce node impurity

Node Impurity Measures If p_k = pct. of class k in node

① mis. class. rate : $1 - \hat{p}_{\hat{k}}$ when \hat{k} = most common class

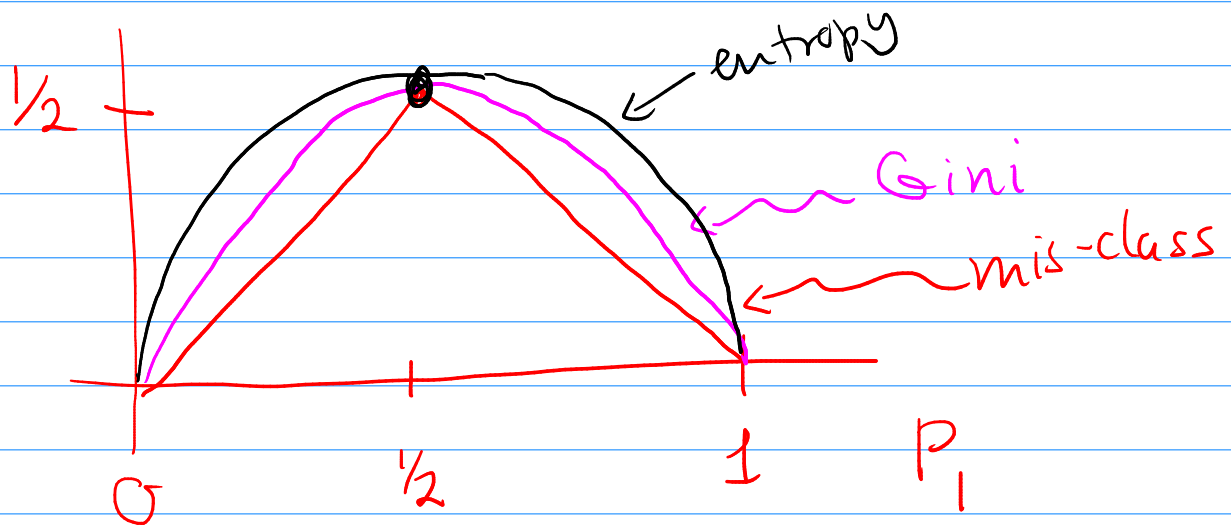
② Gini - Index

$$\sum_k p_k(1-p_k)$$

For all:
pure ≈ 0
impure $\approx 1/K$

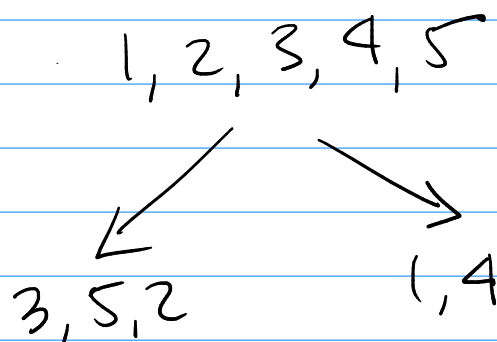
③ Entropy: $\sum_k p_k \log p_k$

$K=2$ (two-class problem)



What about Categorical vars?

Splitting a cat. var. means dividing the cats into two groups:



Problem! If I have g levels in my cat. var then I have $2^g - 1$

possible splits.

CARTs can deal w/ missing data nicely.

Cat. vars just add a "missing" category

Numeric vars. Keep track of "surrogate"
splits i.e. splits using other
vars that have similar divisions

idea: use surrogate if var is missing.
