

MY474: Applied Machine Learning for Social Science

Lecture 5: Regularization, Decision Trees

Blake Miller

29 January 2021

Agenda

1. Regularization
 - ▶ Ridge regression
 - ▶ Lasso regression
2. Tree Regression and Classification

Regularization

Regularization

- ▶ **Ridge regression** and **LASSO** are regularized linear models
- ▶ These models **constrain** or **regularize** coefficient estimates, or equivalently, **shrink** coefficient estimates towards zero.
- ▶ Why would we want to shrink our coefficients?
 1. Bias-variance tradeoff: The least squares estimates often have lower bias and higher variance, reducing variance could improve prediction.
 2. Interpretation: With too many independent variables models become unparsimonious (kitchen sink regressions) and difficult to interpret. Sometimes, regularization can help us automatically perform **feature selection**.
 3. Big p data: OLS models can not handle $p > n$ data.

Ridge Regression

- ▶ For OLS regression, we estimate β by minimizing RSS:

$$\text{RSS} = \varepsilon' \varepsilon = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$$

- ▶ The objective for ridge regression coefficient estimates $\hat{\beta}^R$ is

$$\min_{\beta} [(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta'\beta]$$

where $\lambda \geq 0$ is a **hyperparameter** controlling how much to shrink coefficients.

Deriving Ridge Regression Coefficients

- ▶ The objective is to minimize **penalized RSS (PRSS)**:

$$\min_{\beta} [\text{PRSS}(\lambda)] = \min_{\beta} [(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta'\beta]$$

$$\begin{aligned}\frac{\partial}{\partial\beta} [\text{PRSS}(\lambda)] &= \frac{\partial}{\partial\beta} \{(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) + \lambda\beta'\beta\} \\ &= \frac{\partial}{\partial\beta} \{(\mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\beta + \beta'\mathbf{X}'\mathbf{X}\beta + \beta'(\lambda\mathbf{I})\beta\} \\ &= -2\mathbf{X}'\mathbf{y} + 2(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})\beta\end{aligned}$$

- ▶ Solving for zero yields

$$\hat{\beta}^R = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$$

Shrinkage Penalty

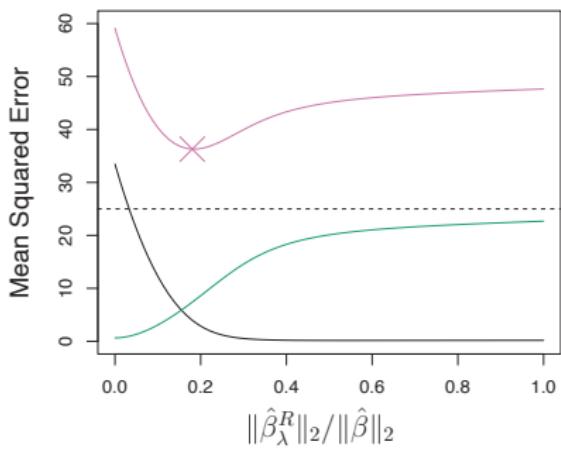
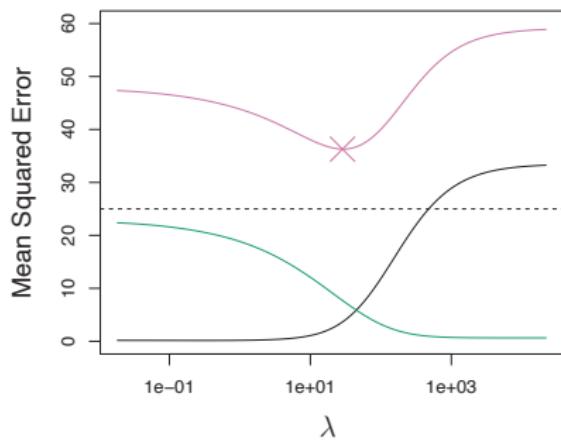
- ▶ The ℓ_2 norm is ridge regression's **shrinkage penalty** which controls the complexity of the model:

$$\|\beta\|_2 = \lambda \beta' \beta = \lambda \sum_{j=1}^p \beta_j^2$$

- ▶ When coefficients are small, $\|\beta\|_2$ is small, and vice versa.
- ▶ PRSS(\cdot) controls the tradeoff between fit to data and model complexity using λ :
 - ▶ When λ is large, we “shrink” coefficients towards zero.
 - ▶ When λ is small, it moves closer to **ordinary least squares** estimates.
- ▶ Selecting a good value for λ is critical; cross-validation is used for this.

Why Does Ridge Regression Improve Over Least Squares?

- ▶ Regularization improves predictions by **trading variance for bias**; best tradeoff determined with **cross-validation**.
- ▶ An example using simulated data with $n = 50$ observations, $p = 45$ predictors, all nonzero coefficients.



Squared bias (black), variance (green), and test mean squared error (purple) for the ridge regression predictions. Horizontal dashed lines indicate the minimum possible MSE.

LASSO (Least Absolute Shrinkage and Selection Operator)

- ▶ Ridge regression does have one obvious disadvantage: it does not select models that involve just a subset of the variables
- ▶ Ridge regression will include all p predictors in the final model
- ▶ The Lasso is a relatively recent alternative to ridge regression that overcomes this disadvantage. The LASSO coefficients, $\hat{\beta}_\lambda^L$, minimize the quantity

$$\text{PRSS}(\lambda) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^p |\beta_j|$$

- ▶ The LASSO uses an ℓ_1 rather than an ℓ_2 shrinkage penalty. The ℓ_1 norm of $\boldsymbol{\beta}$ is given by $\sum_{j=1}^p |\beta_j|$.

LASSO (Least Absolute Shrinkage and Selection Operator)

- ▶ As with ridge regression, the LASSO shrinks the coefficient estimates towards zero.
- ▶ However, in the case of the LASSO, the ℓ_1 penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter λ is sufficiently large.
- ▶ Hence, the LASSO performs **variable selection**.
- ▶ We say that the LASSO yields **sparse** models — that is, models that involve only a subset of the variables.
- ▶ As in ridge regression, selecting a good value of λ for the LASSO is critical; cross-validation is again the method of choice.

The Variable Selection Property of the Lasso

- ▶ Why does the LASSO result in coefficient estimates that are exactly equal to zero?
- ▶ Lasso and ridge regression coefficient estimation can be reformulated as constrained optimization problems:

$$\underset{\beta}{\text{minimize}} \left\{ (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \right\} \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| < s$$

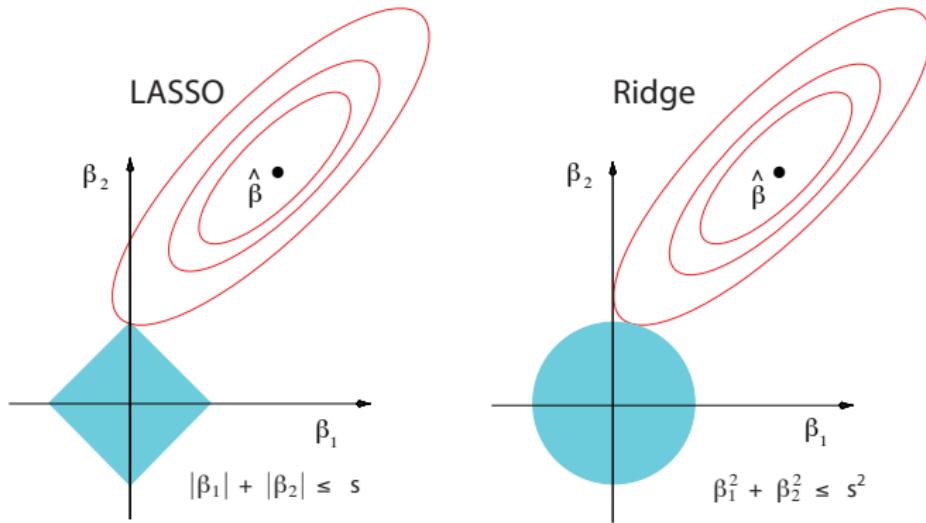
and

$$\underset{\beta}{\text{minimize}} \left\{ (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \right\} \quad \text{subject to} \quad \beta' \beta < s$$

respectively.

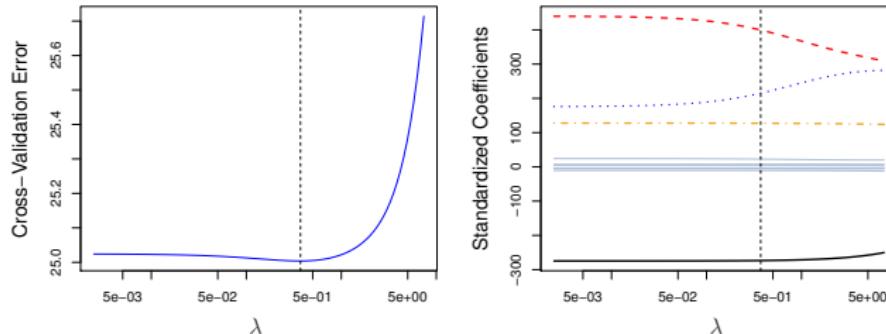
- ▶ λ is the lagrange multiplier for this problem.

The Lasso Picture



- ▶ Contours represent RSS function; center represents the minimum least squares estimate of $\hat{\beta}$.
- ▶ Objective: find the minimum value of RSS that satisfies the constraint (within blue area).
- ▶ ℓ_2 norm constraint region is circular, ℓ_1 is a diamond.

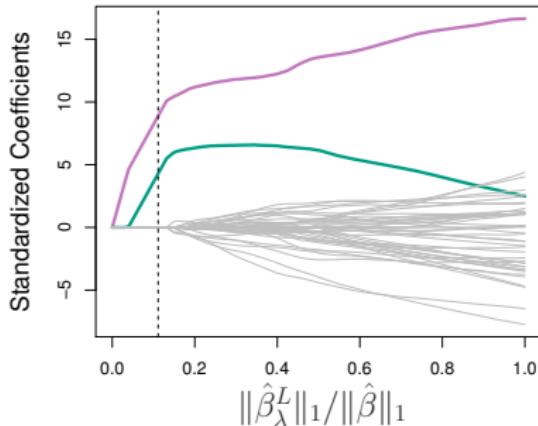
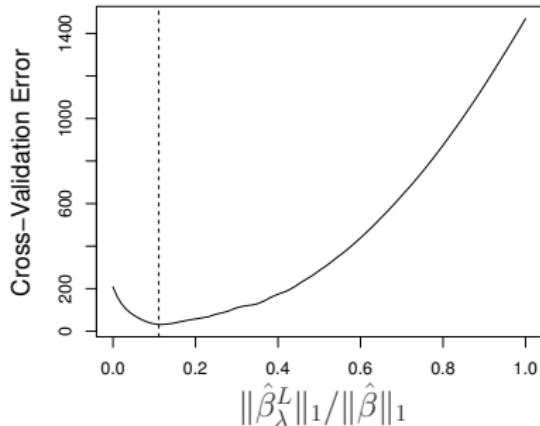
Selecting the Tuning Parameter λ



Example of ridge regression coefficients at different levels of λ . λ is chosen where CV error (left) is minimized (vertical dotted lines); corresponding coefficients on the right.

1. Choose a grid search of λ values (or equivalent method), compute the cross-validation error rate for each value of λ .
2. Select λ where cross-validation error is either 1) smallest, or 2) one standard error from the smallest cross-validation error (**one s.e. rule**).
3. Refit model to available observations with selected value of λ .

Visualizing Variable Selection in the LASSO



- ▶ Left: Ten-fold cross-validated MSE for the LASSO, applied to simulated data (only two inputs related to the output).
- ▶ Right: The corresponding LASSO coefficient estimates are displayed. The vertical dashed lines indicate the LASSO fit for which the cross-validation error is smallest.

LASSO: Advantages and Disadvantages

Advantages

- ▶ LASSO achieves high predictive accuracy when the true model is **sparse**, that is, when only a few variables matter.
- ▶ Interpretation becomes easier when models are more sparse.
- ▶ Post-LASSO gives asymptotically valid confidence intervals.
- ▶ LASSO works well for models with many features.

Disadvantages

- ▶ LASSO won't work well when true models are not sparse—many of the variables are important (ridge works better in that case).
- ▶ With high collinearity, LASSO arbitrarily selects one feature among the correlated features (fine if your goal is prediction).
- ▶ Coefficient estimates vary quite a bit with slightly different samples, however, predictions will be similar.
- ▶ LASSO can select at most n variables.

Statistical Inference with LASSO

We often care about inference rather than prediction; need confidence intervals for $\hat{\beta}$. Several options:

1. Post-Lasso (Belloni and Chernozhukov, 2013)
 - ▶ After “hard thresholding” with LASSO, take the surviving features and run OLS.
 - ▶ Will give a consistent estimate under some conditions (“approximate sparsity”)
2. Covariance test (Lockhart, Taylor, Tibshirani, 2014)
 - ▶ P-value for each variable as it is added to LASSO model
3. Bayesian LASSO (Park and Casella, 2008)
 - ▶ Bayesian model with a Laplacian prior over β_j :

$$y_i | \beta, \lambda \sim \mathcal{N}\left(\sum_{k=1}^K x_{ik} \beta_k, \sigma_\epsilon^2\right)$$
$$\beta_k | \lambda \propto \exp(-\lambda |\beta_k|)$$

- ▶ The LASSO estimator corresponds with the mode of the Bayesian posterior.

Example: Chinese Govt. Censorship of Academic Articles

Cambridge University Press accused of 'selling its soul' over Chinese censorship

Academics and activists decry publisher's decision to comply with
a Chinese request to block more than 300 articles from leading
China studies journal



▲ A list of the blocked articles, published by CUP, shows they focus overwhelmingly on topics China's one-party state regards as taboo Photograph: Nick Ansell/PA

Source: [The Guardian](#)

Example: Chinese Censorship of Academic Journal Articles

- ▶ We classify the titles of all articles published in China Quarterly
- ▶ Using the [list of censored articles](#) published by the journal, we construct a binary outcome variable $censored \in \{0, 1\}$
- ▶ Features x are counts of all words and phrases that appear in at least 15 article titles
- ▶ This classifier can help us infer the logic behind the censorship order.

Defining Features

```
library(quantada)
library(glmnet)

cq_data <- read.csv('cq.csv', stringsAsFactors = F)
cq <- corpus(cq_data, text_field = 'title')

tokens <- tokens(cq) %>% tokens_ngrams(n = 1:3)
dfm <- dfm(tokens, remove_punct=T, remove_numbers=T)
dfm <- dfm_trim(dfm, min_docfreq = 15)
dim(dfm)

## [1] 6402 891
```

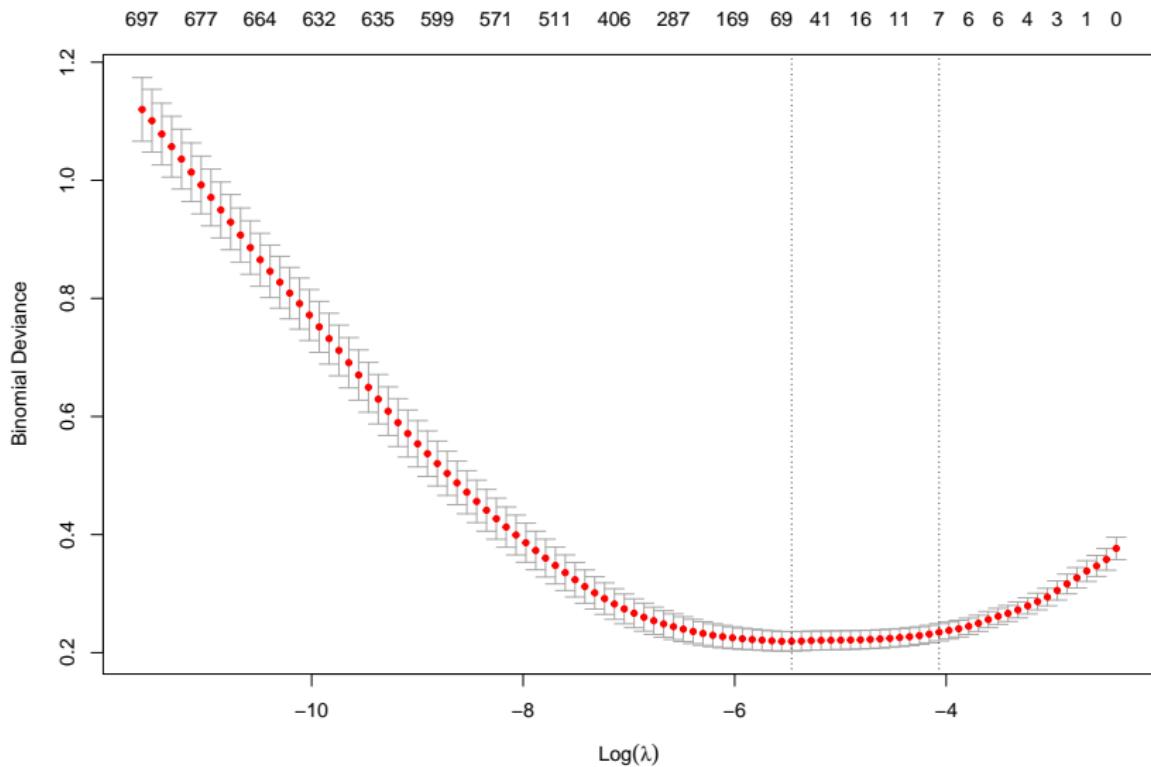
Training a Lasso Model

```
n_tr <- floor(.80 * nrow(dfm))
tr <- sample(1:nrow(dfm), n_tr)

y <- dfm$censored
mod <- cv.glmnet(dfm[tr,], y[tr], family="binomial")
```

Choosing Lambda

```
plot(mod)
```



Examining 10 Largest Coefficients

We can inspect the highest-magnitude coefficients of this regression model to infer how censorship decisions were made.

```
l <- which(mod$lambda==mod$lambda.min)
l_1se <- which(mod$lambda==mod$lambda.1se)
coefs <- mod$glmnet.fit$beta[,1]
head(coefs[order(coefs, decreasing=T)], 10)
```

##	tiananmen	tibet
##	6.357448	5.319666
##	xinjiang	cultural_revolution
##	4.911321	4.419312
##	red	the_republic_of
##	3.253431	3.145193
##	tibetan	of_"
##	2.268141	1.748898
##	taiwanese	buddhism
##	1.606889	1.596889

Examining All Non-Zero Coefficients (1 S.E. rule)

```
coefs <- mod$glmnet.fit$beta[,l_1se]
coefs[coefs != 0]

##          taiwan cultural_revolution
##      0.1409574        4.0432128
##          tibet             red
##      4.0502237        2.0559735
##  the_republic_of     tiananmen
##      2.5813547        4.4856124
##          xinjiang
##      3.5183851
```

Examining Coefficients

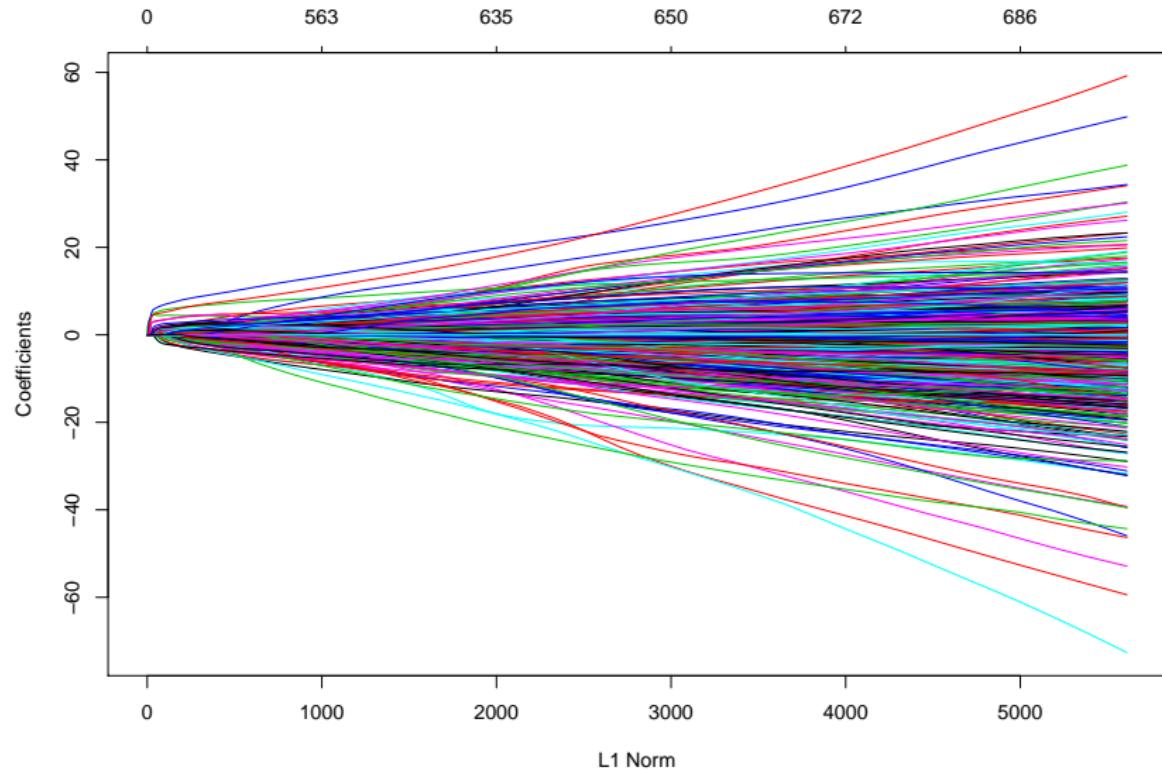


The highest-magnitude coefficients correspond to terms resembling the “Three T’s”:

1. Taiwan (Republic of China)
2. Tibet (and Xinjiang)
3. Tiananmen (student movement of 1989)

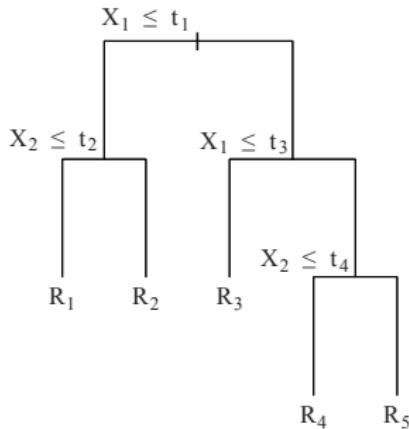
Exploring How Features were “Shrunk”

```
mod <- glmnet(dfm[tr,], y[tr], family="binomial")  
plot(mod)
```



Trees

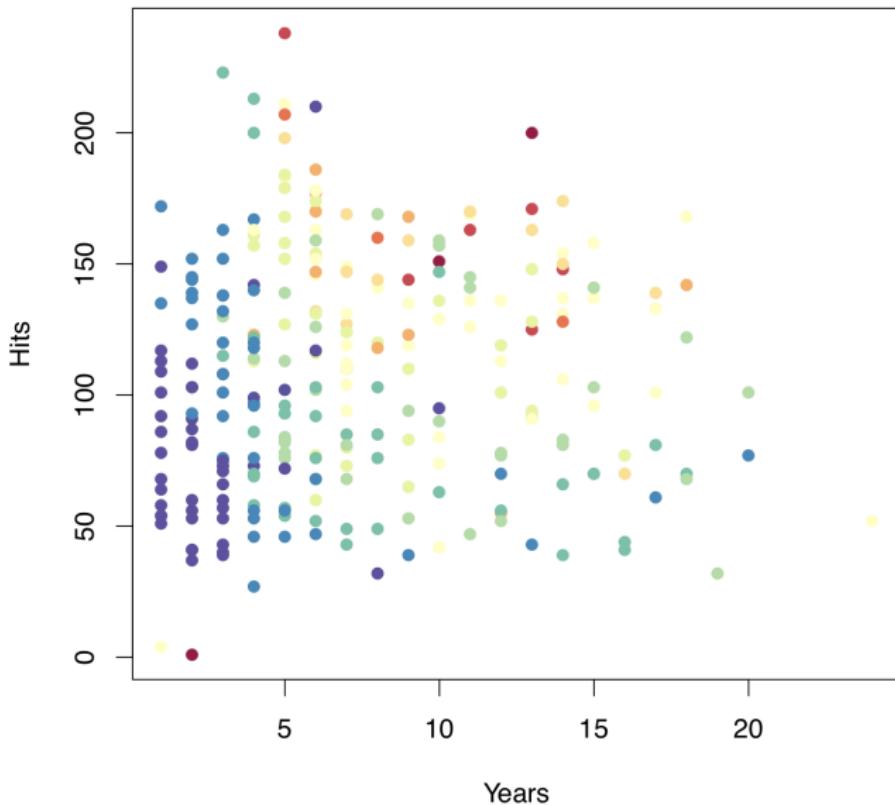
Tree-based Methods



A partition of the feature space into regions R_1, \dots, R_5 , visualized in a tree.

- ▶ **Tree-based** methods are used for regression or classification.
- ▶ Trees recursively partition the feature space \mathbf{x} into smaller regions R_1, \dots, R_J .
- ▶ Each region is created by a binary split at a value t_k of x_j .
- ▶ Each region corresponds to a single outcome of y

Baseball salary data: how would you stratify it?



Decision tree for these data

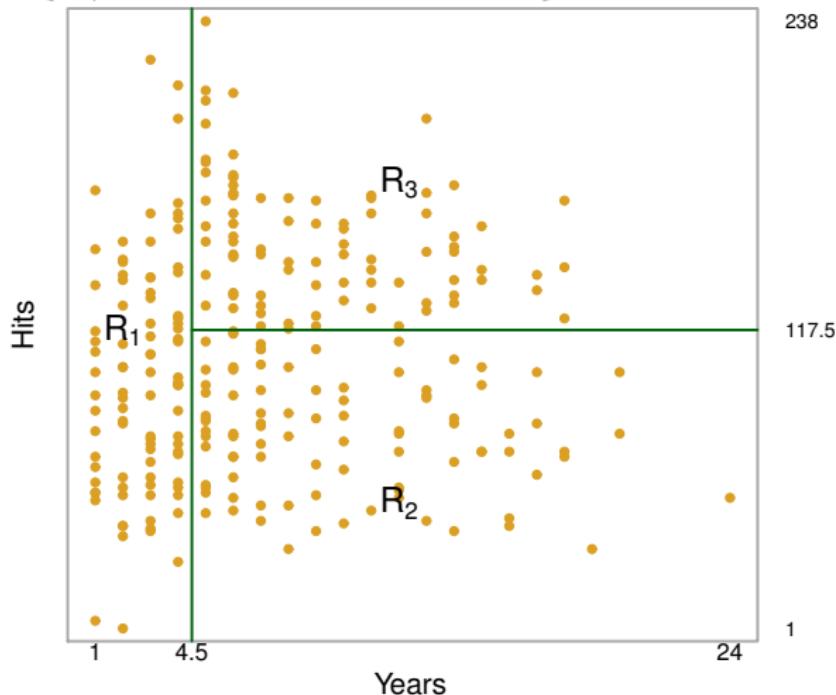


Details of previous figure

- ▶ For the Hitters data, a regression tree for predicting the log salary of a baseball player, based on the number of years that he has played in the major leagues and the number of hits that he made in the previous year.
- ▶ At a given internal node, the label (of the form $X_j < t_k$) indicates the left-hand branch emanating from that split, and the right-hand branch corresponds to $X_j \geq t_k$. For instance, the split at the top of the tree results in two large branches. The left-hand branch corresponds to Years < 4.5 , and the right-hand branch corresponds to Years ≥ 4.5 .
- ▶ The tree has two internal nodes and three terminal nodes, or leaves. The number in each leaf is the mean of the response for the observations that fall there.

Results

- Overall, the tree stratifies or segments the players into three regions of predictor space: $R_1 = \{X | \text{Years} < 4.5\}$, $R_2 = \{X | \text{Years} \geq 4.5, \text{Hits} < 117.5\}$, and $R_3 = \{X | \text{Years} \geq 4.5, \text{Hits} \geq 117.5\}$.



Terminology for Trees

- ▶ In keeping with the **tree** analogy, the regions R_1 , R_2 , and R_3 are known as **terminal nodes**
- ▶ Decision trees are typically drawn upside down, in the sense that the leaves are at the bottom of the tree.
- ▶ The points along the tree where the predictor space is split are referred to as **internal nodes**
- ▶ In the hitters tree, the two internal nodes are indicated by the text Years < 4.5 and Hits < 117.5.

Interpretation of Results

- ▶ Years is the most important factor in determining Salary, and players with less experience earn lower salaries than more experienced players.
- ▶ Given that a player is less experienced, the number of Hits that he made in the previous year seems to play little role in his Salary.
- ▶ But among players who have been in the major leagues for five or more years, the number of Hits made in the previous year does affect Salary, and players who made more Hits last year tend to have higher salaries.
- ▶ Surely an over-simplification, but compared to a regression model, it is easy to display, interpret and explain

Details of the tree-building process

1. We divide the predictor space — that is, the set of possible values for X_1, X_2, \dots, X_p — into J distinct and non-overlapping regions, R_1, R_2, \dots, R_J .
2. For every observation that falls into the region R_j , we make the same prediction, which is simply the mean of the response values for the training observations in R_j .

More details of the tree-building process

- ▶ In theory, the regions could have any shape. However, we choose to divide the predictor space into high-dimensional rectangles, or boxes, for simplicity and for ease of interpretation of the resulting predictive model.
- ▶ The goal is to find boxes R_1, \dots, R_J that minimize the RSS, given by

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

where \hat{y}_{R_j} is the mean response for the training observations within the j th box.

- ▶ In other words, the goal is to partition the data into boxes that balance **homogeneity** within boxes and **distinctness** among boxes.

More details of the tree-building process

- ▶ Unfortunately, it is computationally infeasible to consider every possible partition of the feature space into J boxes.
- ▶ For this reason, we take a **top-down, greedy** approach that is known as recursive binary splitting.
- ▶ The approach is **top-down** because it begins at the top of the tree and then successively splits the predictor space; each split is indicated via two new branches further down on the tree.
- ▶ It is **greedy** because at each step of the tree-building process, the **best** split is made at that particular step, rather than looking ahead and picking a split that will lead to a better tree in some future step.

Why Binary Splits?

- ▶ Multiway splits into more than two groups fragment the data too quickly, leaving insufficient data at the next level down.
- ▶ Since multiway splits can be achieved by a series of binary splits, the latter are preferred.

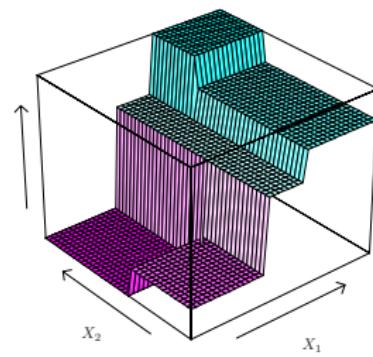
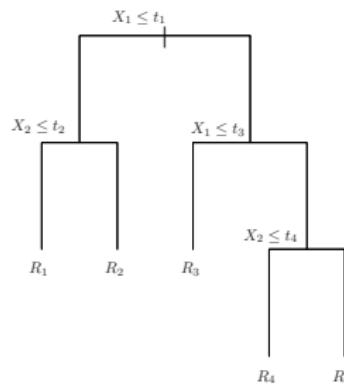
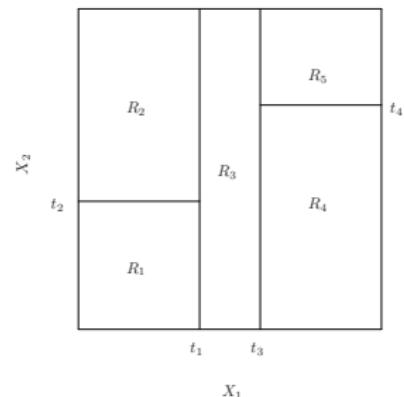
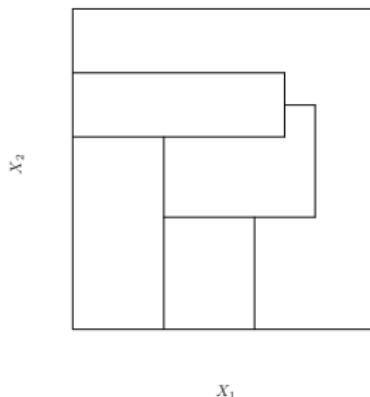
Details— Continued

- ▶ We first select the predictor X_j and the cutpoint s such that splitting the predictor space into the regions $\{X|X_j < s\}$ and $\{X|X_j \geq s\}$ leads to the greatest possible reduction in **RSS**.
- ▶ Next, we repeat the process, looking for the best predictor and best cutpoint in order to split the data further so as to minimize the RSS within each of the resulting regions.
- ▶ However, this time, instead of splitting the entire predictor space, we split one of the two previously identified regions. We now have three regions.
- ▶ Again, we look to split one of these three regions further, so as to minimize the RSS. The process continues until a stopping criterion is reached; for instance, we may continue until no region contains more than five observations.

Predictions

- ▶ We predict the response for a given test observation using the mean of the training observations in the region to which that test observation belongs.
- ▶ A five-region example of this approach is shown in the next slide.

Visualizing Tree Regression



Details of previous figure

- ▶ Top Left: A partition of two-dimensional feature space that could not result from recursive binary splitting.
- ▶ Top Right: The output of recursive binary splitting on a two-dimensional example.
- ▶ Bottom Left: A tree corresponding to the partition in the top right panel.
- ▶ Bottom Right: A perspective plot of the prediction surface corresponding to that tree.

Pruning a tree

- ▶ The process described above may produce good predictions on the training set, but is likely to **overfit** the data, leading to poor test set performance. **Why?**
- ▶ Imagine a tree with one observation in each terminal node.
What is our training error?
- ▶ A smaller tree with fewer splits (that is, fewer regions R_1, \dots, R_J) might lead to lower variance and better interpretation at the cost of a little bias.
- ▶ One possible alternative to the process described above is to grow the tree only so long as the decrease in the RSS due to each split exceeds some (high) threshold.
- ▶ This strategy will result in smaller trees, but is too **short-sighted**: a seemingly worthless split early on in the tree might be followed by a very good split — that is, a split that leads to a large reduction in RSS later on.

Pruning a tree— continued

- ▶ A better strategy is to grow a very large tree T_0 , and then **prune** it back in order to obtain a **subtree**
- ▶ **Cost complexity pruning** — also known as **weakest link pruning** — is used to do this
- ▶ we consider a sequence of trees indexed by a nonnegative tuning parameter α . For each value of α there corresponds a subtree $T \subset T_0$ such that

$$\sum_{m=1}^{|T|} \sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

is as small as possible. Here $|T|$ indicates the number of terminal nodes of the tree T , R_m is the rectangle (i.e. the subset of predictor space) corresponding to the m th terminal node, and \hat{y}_{R_m} is the mean of the training observations in R_m .

Choosing the best subtree

- ▶ The tuning parameter α controls a trade-off between the subtree's complexity and its fit to the training data.
 - ▶ $\alpha = 0$ corresponds to minimizing training error (will choose the largest tree, all pure nodes)
 - ▶ Large α results in small trees
- ▶ We select an optimal value $\hat{\alpha}$ using cross-validation.
- ▶ We then return to the full data set and obtain the subtree corresponding to $\hat{\alpha}$.

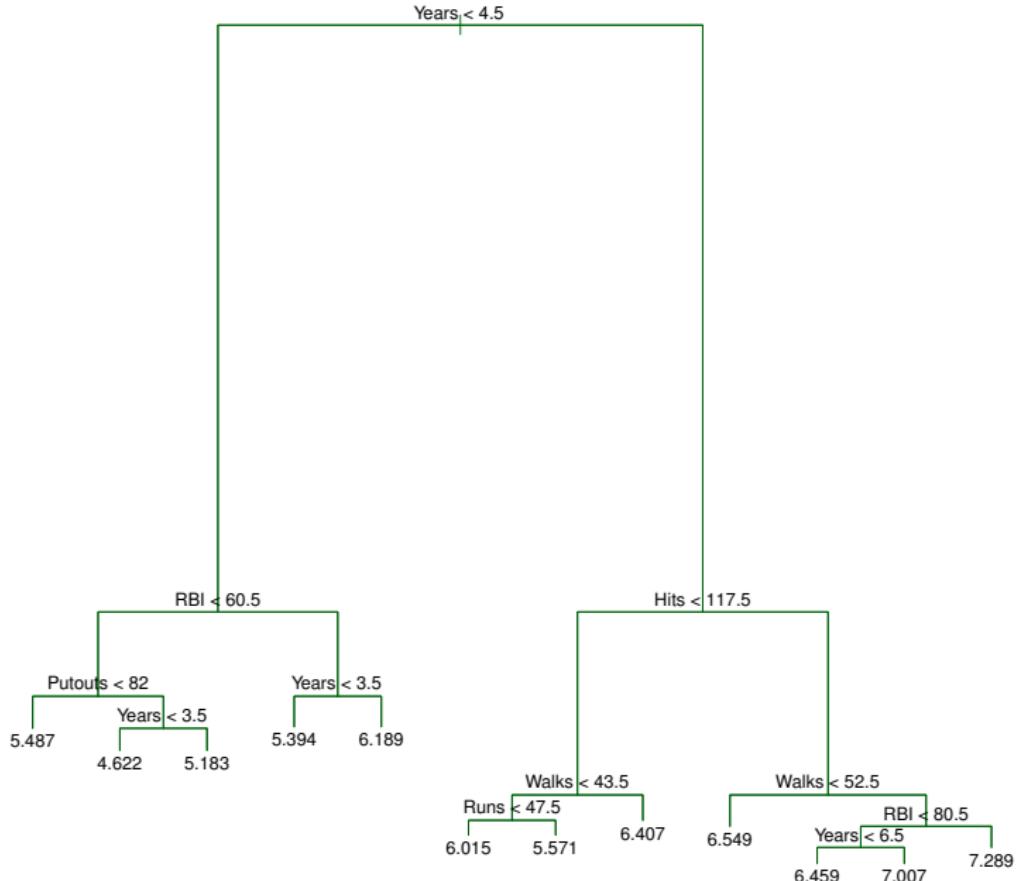
Summary: tree algorithm

1. Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
2. Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of α .
3. Use K-fold cross-validation to choose α . For each $k = 1, \dots, K$:
 - ▶ Repeat Steps 1 and 2 on the $\frac{K-1}{K}$ th fraction of the training data, excluding the k th fold.
 - ▶ Evaluate the mean squared prediction error on the data in the left-out k th fold, as a function of α . Average the results, and pick α to minimize the average error.
4. Return the subtree from Step 2 that corresponds to the chosen value of α .

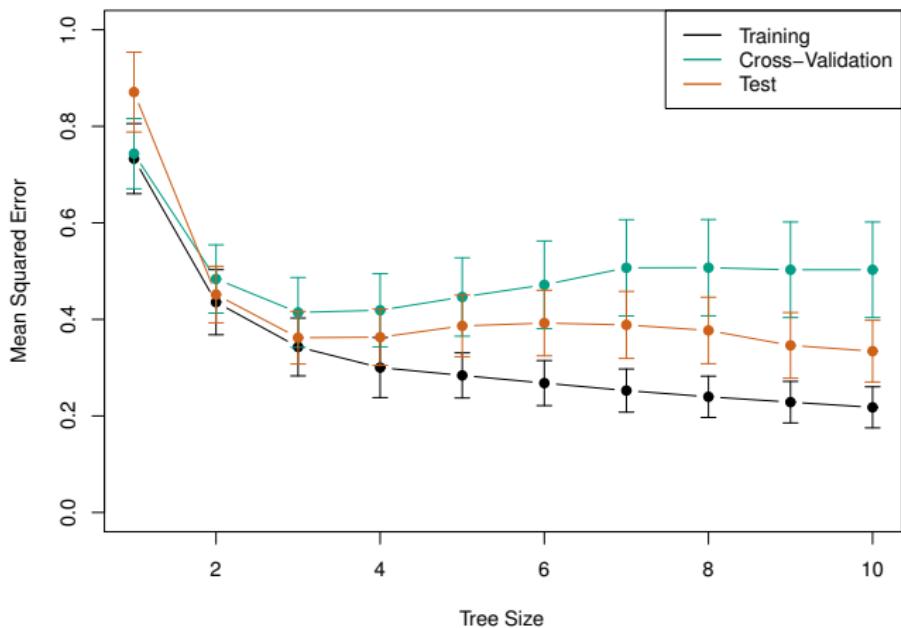
Baseball example continued

- ▶ Recall the baseball example from before. How did we decide on that tree?
- ▶ First, we randomly divided the data set in half, yielding 132 observations in the training set and 131 observations in the test set.
- ▶ We then built a large regression tree on the training data and varied α in order to create subtrees with different numbers of terminal nodes.
- ▶ Finally, we performed six-fold cross-validation in order to estimate the cross-validated MSE of the trees as a function of α .

Baseball example continued



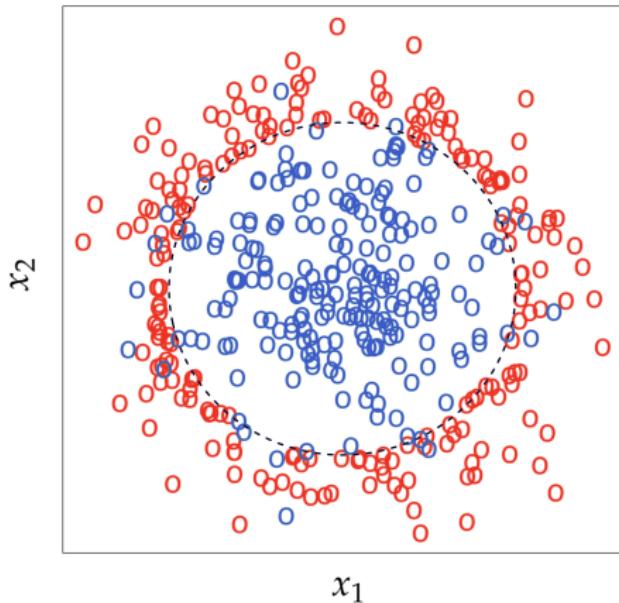
Baseball example continued



Classification Trees

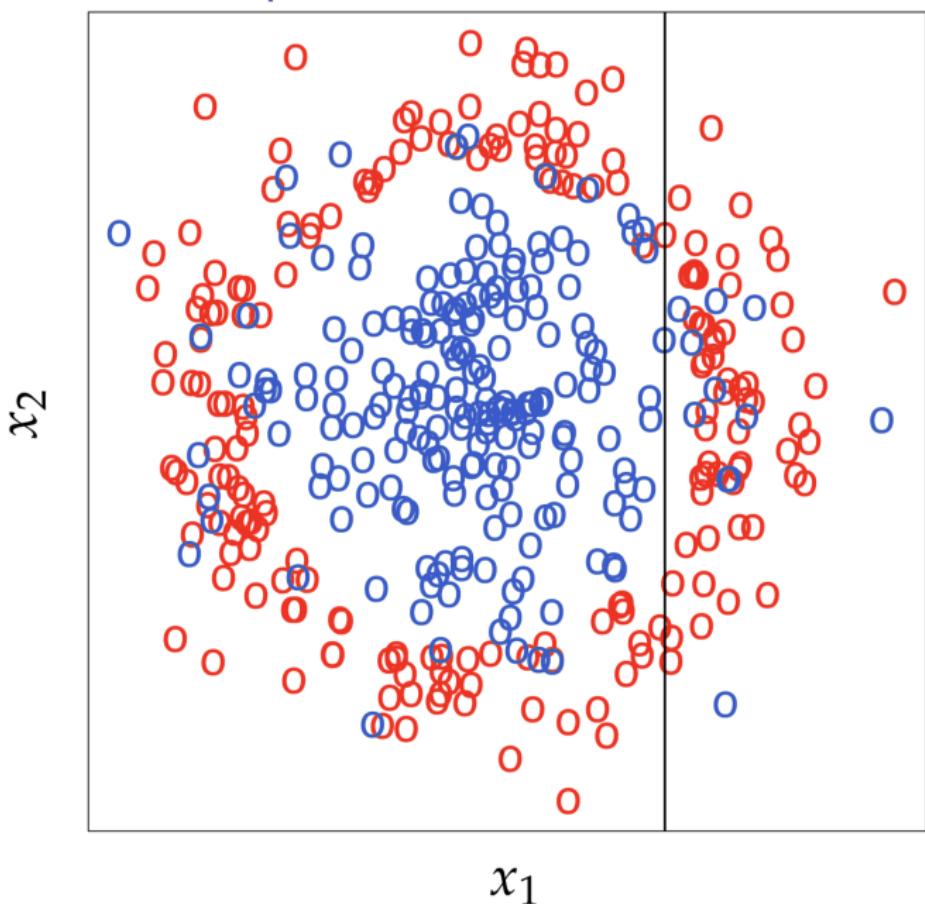
- ▶ Very similar to a regression tree, except that it is used to predict a qualitative response rather than a quantitative one.
- ▶ For a classification tree, we predict that each observation belongs to the most commonly occurring class of training observations in the region to which it belongs.

Visualizing Classification Trees: Donut Data

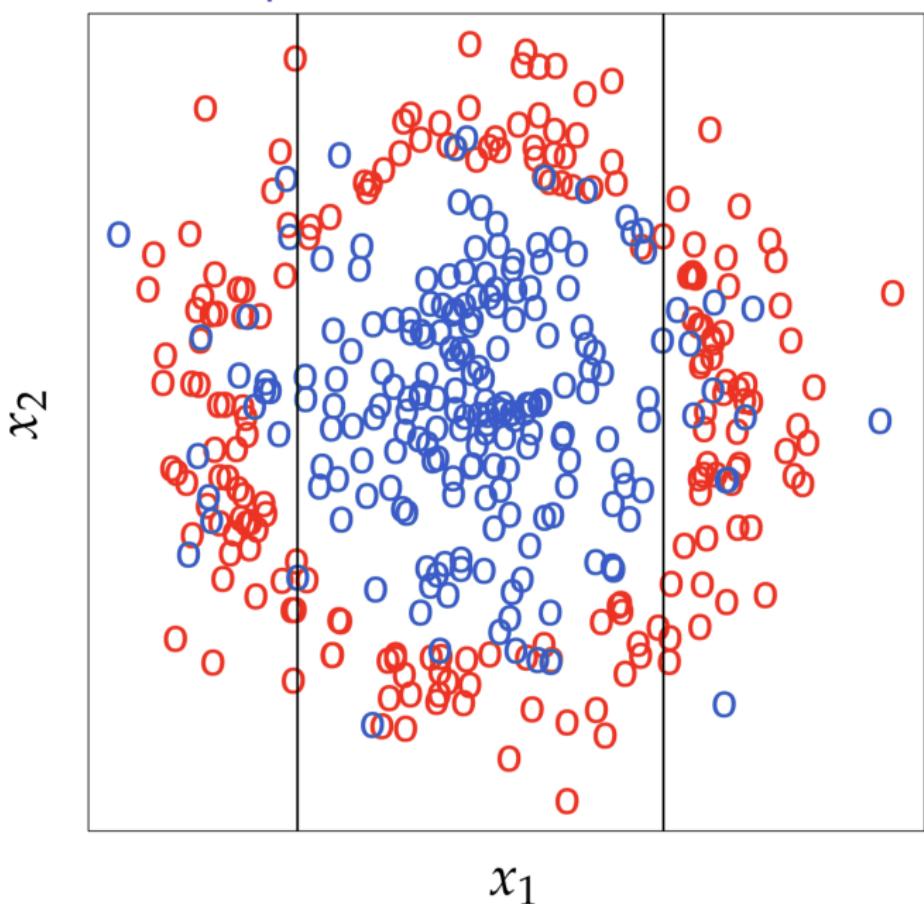


- ▶ Green class: two independent standard normal inputs X_1, X_2
- ▶ Red class: conditional on $X_1^2 + X_2^2 \geq 4.6$

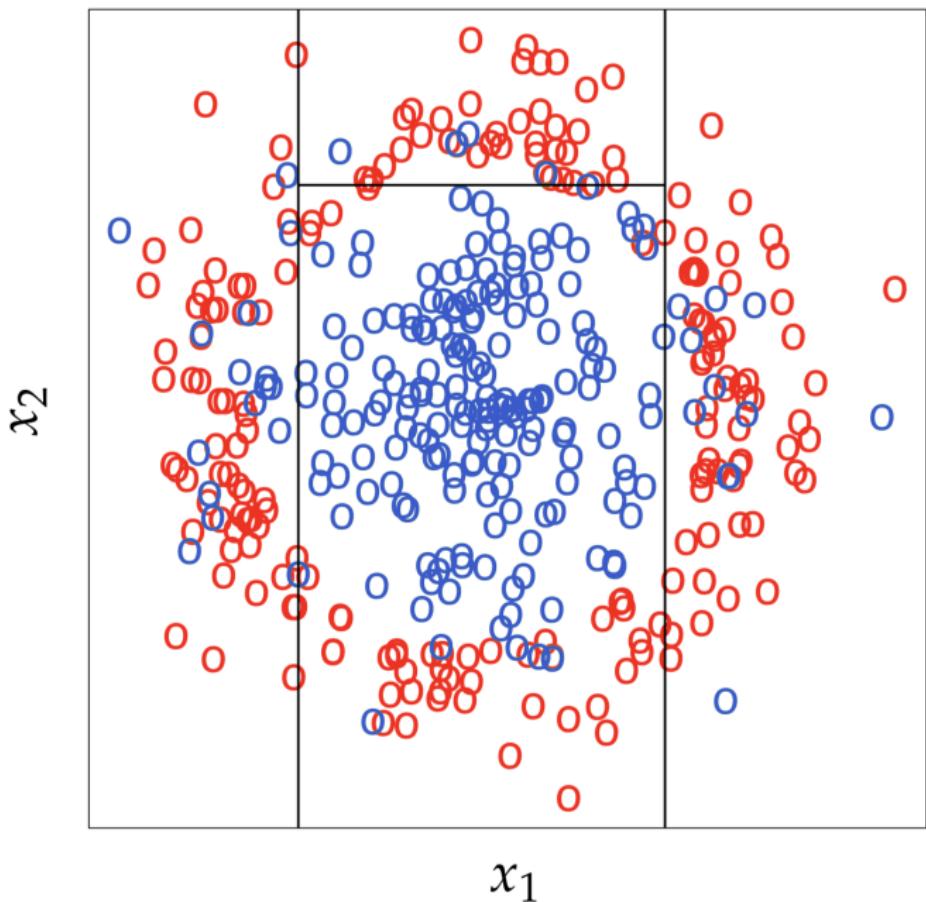
Classification tree: split 1



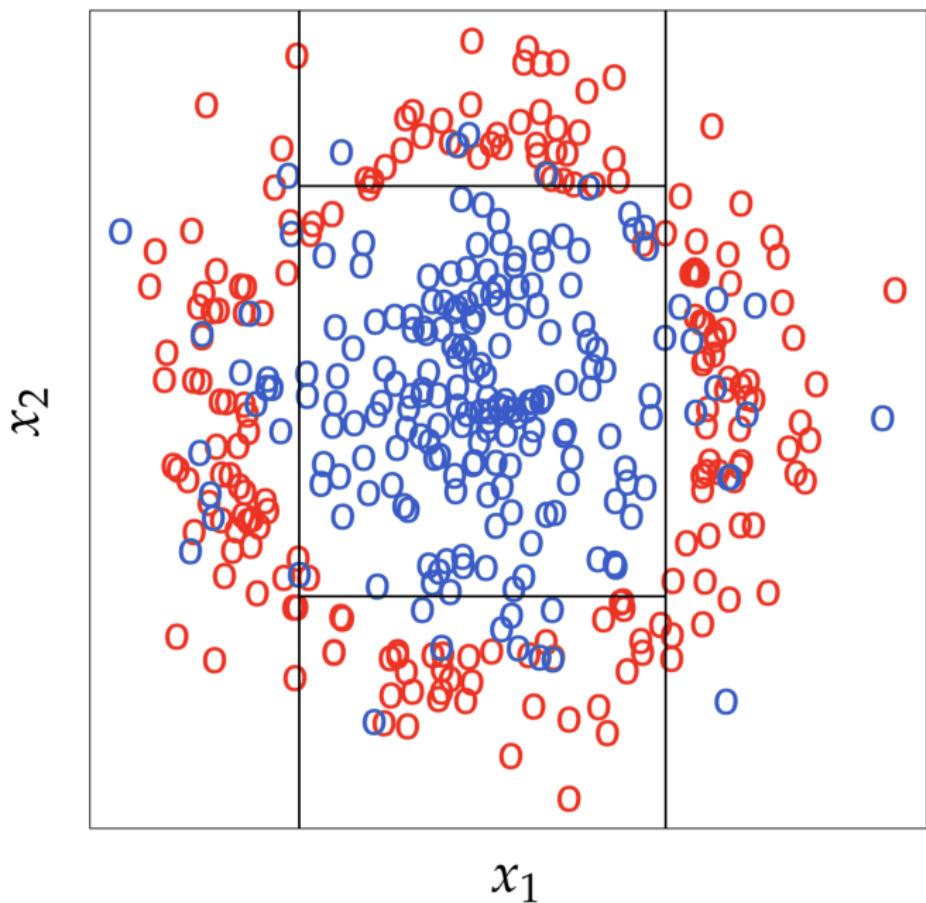
Classification tree: split 2



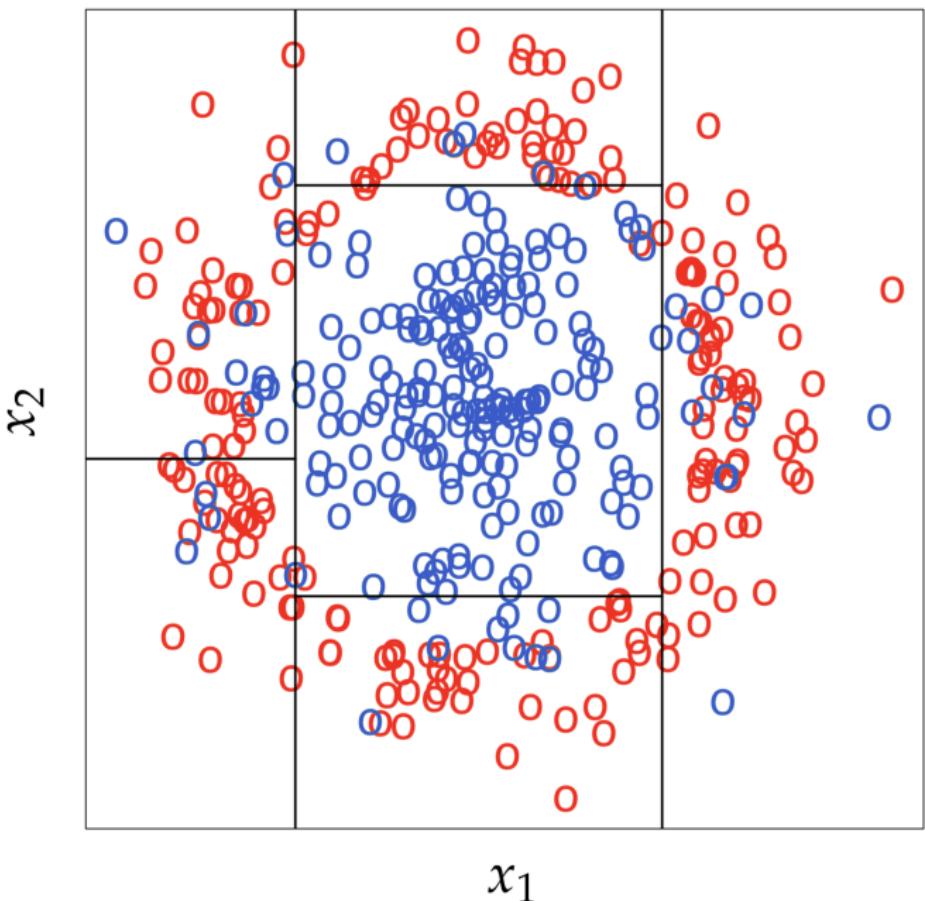
Classification tree: split 3



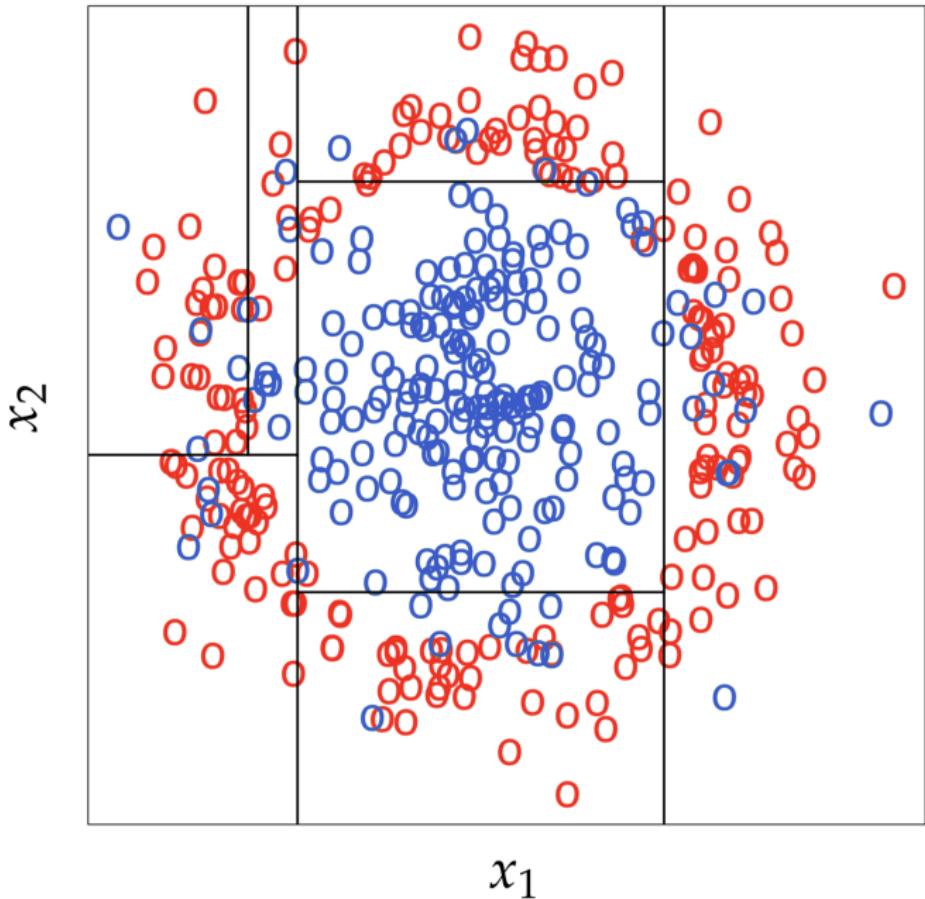
Classification tree: split 4



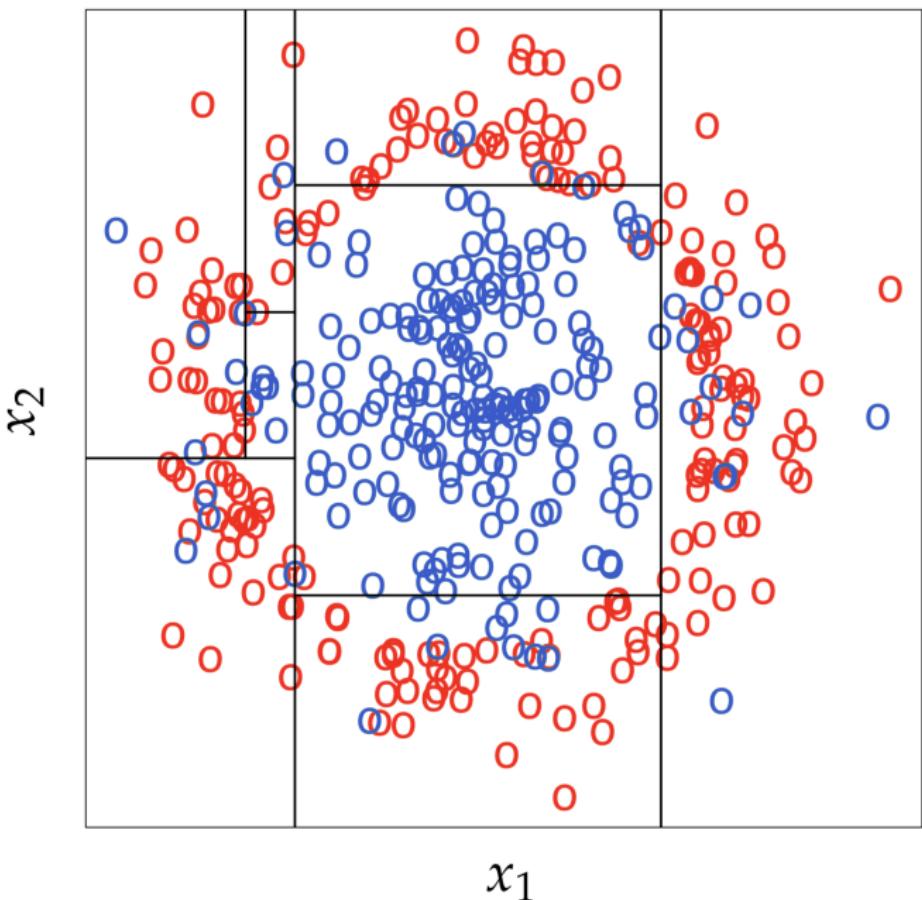
Classification tree: split 5



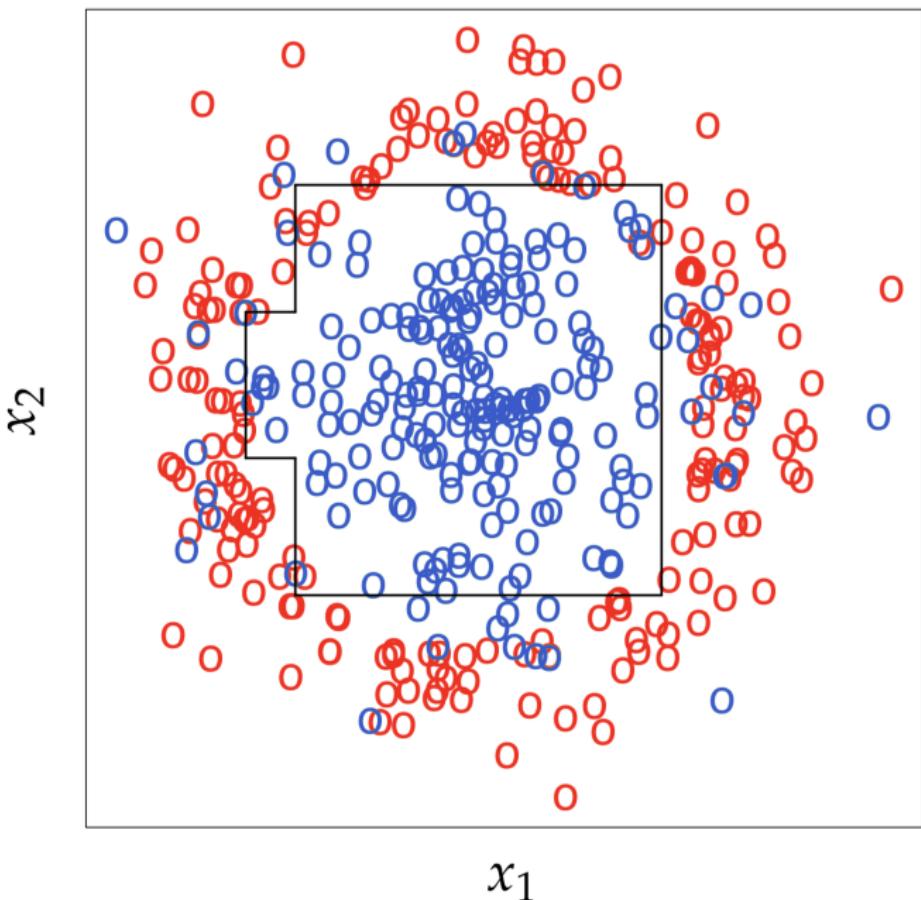
Classification tree: split 6



Classification tree: split 7



Classification tree: final decision boundary



Details of classification trees

- ▶ Just as in the regression setting, we use recursive binary splitting to grow a classification tree.
- ▶ In the classification setting, RSS cannot be used as a criterion for making the binary splits
- ▶ A natural alternative to RSS is the **classification error rate**. this is simply the fraction of the training observations in that region that do not belong to the most common class:
$$E = 1 - \max_k(\hat{p}_{mk})$$
. Here \hat{p}_{mk} represents the proportion of training observations in the mth region that are from the kth class.
- ▶ However classification error is not sufficiently sensitive for tree-growing, and in practice two other measures are preferable.

Gini index and Deviance

- The **Gini index** is defined by

$$G = \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$$

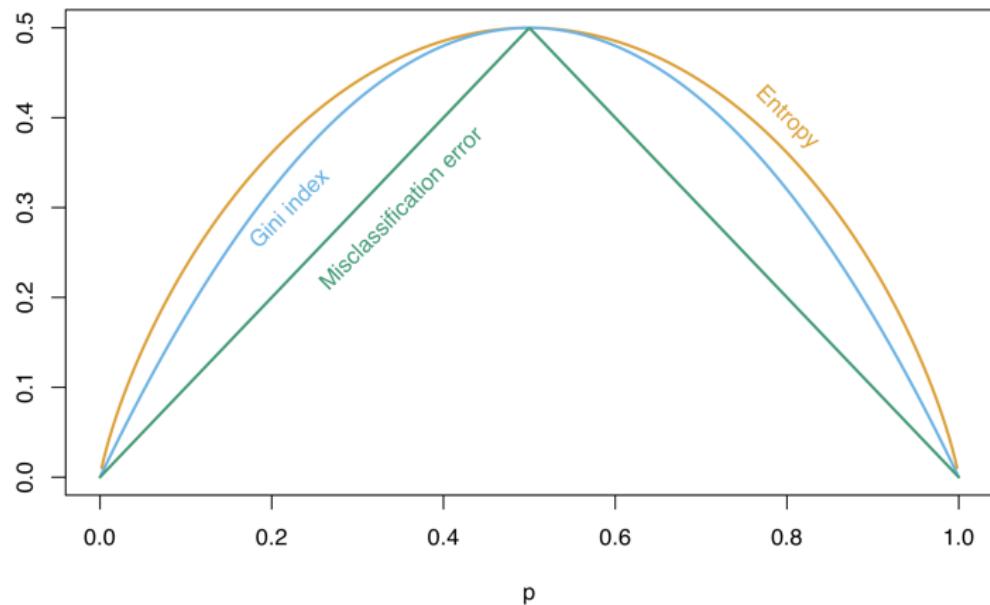
a measure of total variance across the K classes. Gini takes on a small value if all of the \hat{p}_{mk} 's are close to 0 or 1

- Gini represents the variance of a binomial distribution; the **diagonal of the multinomial variance covariance matrix** for > 2 classes
- For this reason the Gini index is referred to as a measure of **node purity** — a small value indicates that a node contains predominantly observations from a single class.
- Alternatively, **cross-entropy** or **deviance**, given by

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$$

- It turns out that the Gini index and the cross-entropy are very similar numerically.

Node impurity measures

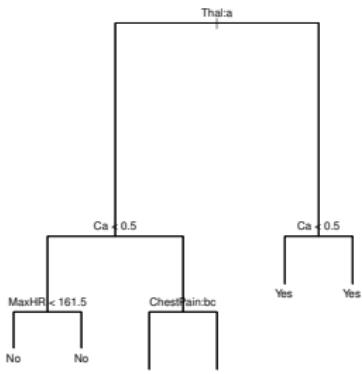
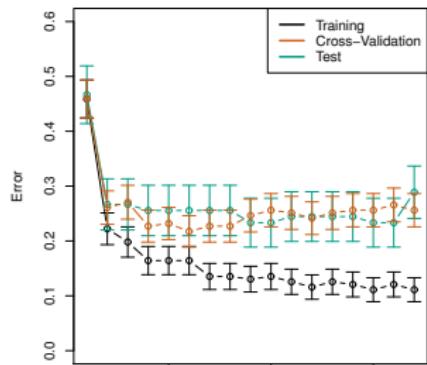
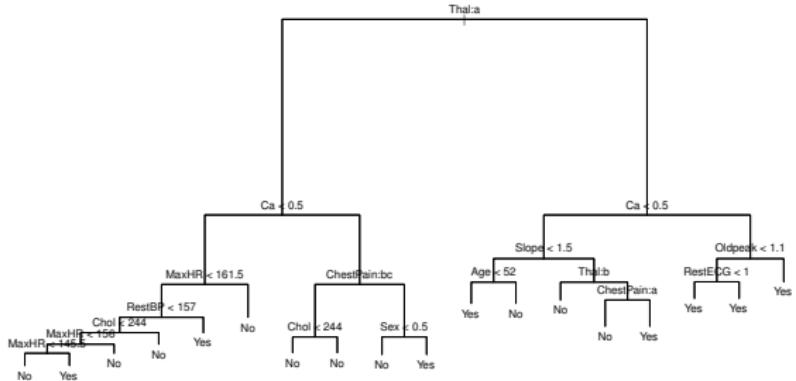


Node impurity measures for two-class classification, as a function of the proportion p in class 2. Cross-entropy has been scaled to pass through $(0.5, 0.5)$.

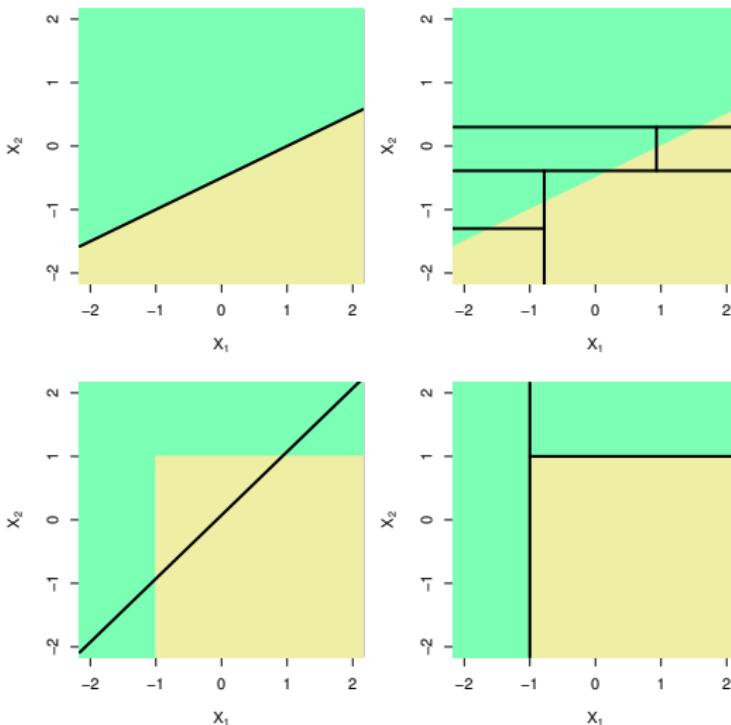
Example: heart data

- ▶ These data contain a binary outcome HD for 303 patients who presented with chest pain.
- ▶ An outcome value of Yes indicates the presence of heart disease based on an angiographic test, while No means no heart disease.
- ▶ There are 13 predictors including Age, Sex, Chol (a cholesterol measurement), and other heart and lung function measurements.
- ▶ Cross-validation yields a tree with six terminal nodes. See next figure.

Example: heart data



Trees Versus Linear Models



Top Row: True linear boundary; Bottom row: true non-linear boundary.
Left column: linear model; Right column: tree-based model

Advantages of Trees

- ▶ Trees are very easy to explain to people. In fact, they are even easier to explain than linear regression!
- ▶ Some people believe that decision trees more closely mirror human decision-making than do the regression and classification approaches seen in previous chapters.
- ▶ Trees can be displayed graphically, and are easily interpreted even by a non-expert (especially if they are small).
- ▶ Trees can easily handle qualitative predictors without the need to create dummy variables.

Disadvantages of Trees

- ▶ Unfortunately, trees generally do not have the same level of predictive accuracy as some of the other regression and classification approaches seen in this book.
- ▶ **High variance** in trees means a small change in the data can result in a very different series of splits; error in the top split is propagated down to all of the splits below it.
- ▶ **Bagging** and **Random Forests** reduce this variance by averaging many trees
- ▶ **Lack of smoothness** of the prediction surface. In classification with 0/1 loss, this doesn't hurt much, but can affect performance in the regression setting.
- ▶ By aggregating many decision trees, the predictive performance of trees can be substantially improved.