

Module 4: Coping with Multiple Predictors

Regression Trees

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 21st, 2013

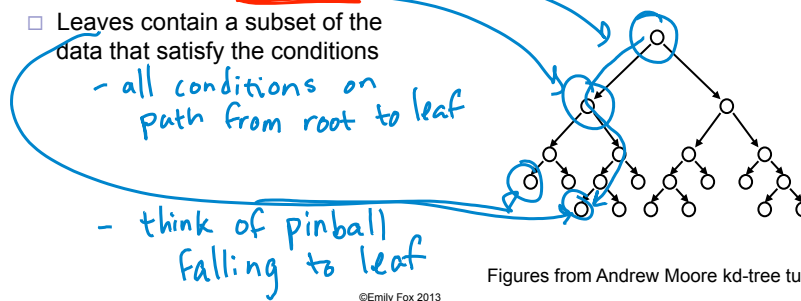
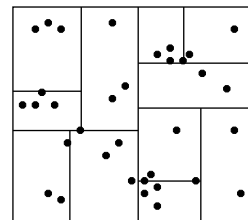
©Emily Fox 2013

1

Recursive Binary Partitions

- To simplify the process and interpretability, consider recursive binary partitions

- Described via a rooted tree
 - Every node of the tree corresponds to split decision
 - Leaves contain a subset of the data that satisfy the conditions

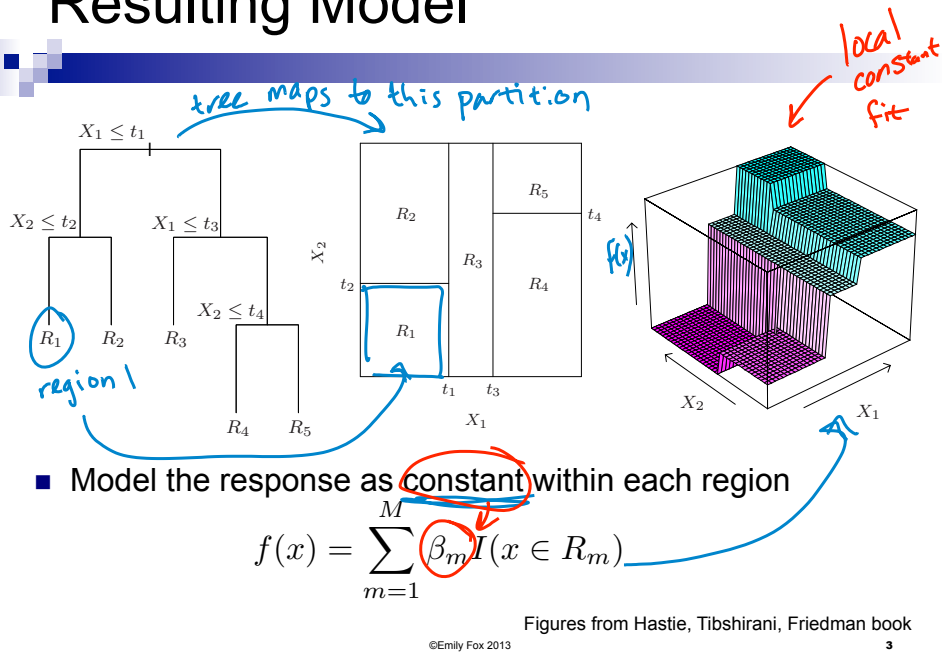


Figures from Andrew Moore kd-tree tutorial

©Emily Fox 2013

2

Resulting Model



Basis Expansion Interpretation

- Equivalent to a basis expansion

$$f(x) = \sum_{m=1}^M \beta_m h_m(x)$$

indicators on regions

- In this example:

$$h_1(x_1, x_2) = I(x_1 \leq t_1)I(x_2 \leq t_2)$$

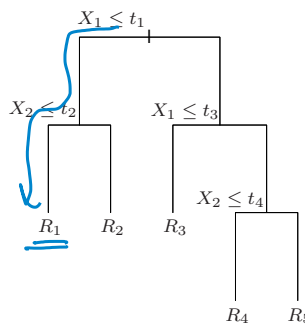
$$h_2(x_1, x_2) = I(x_1 \leq t_1)I(x_2 > t_2)$$

$$h_3(x_1, x_2) = I(x_1 > t_1)I(x_1 \leq t_3)$$

$$h_4(x_1, x_2) = I(x_1 > t_1)I(x_1 > t_3)I(x_2 \leq t_4)$$

$$h_5(x_1, x_2) = I(x_1 > t_1)I(x_1 > t_3)I(x_2 > t_4)$$

reduced tensor product spline w/ step for basis



Choosing a Split Decision

- Starting with all of the data, consider splitting on variable j at point s

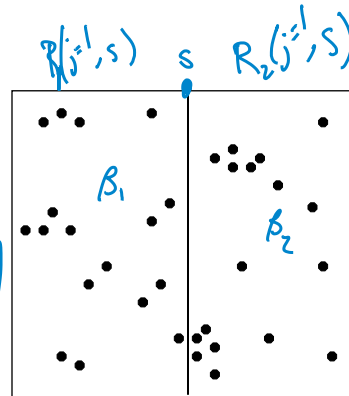
- Define

$$R_1(j, s) = \{x \mid x_j \leq s\}$$

$$R_2(j, s) = \{x \mid x_j > s\}$$

- Our objective is

$$\min_{j, s} \left[\min_{\beta_1} \sum_{x_i \in R_1(j, s)} (y_i - \beta_1)^2 + \min_{\beta_2} \sum_{x_i \in R_2(j, s)} (y_i - \beta_2)^2 \right]$$



- For any (j, s) , the inner minimization is solved by

$$\hat{\beta}_k = \text{avg}(y_i \mid x_i \in R_k(j, s)) \quad k=1, 2$$

©Emily Fox 2013

5

Cost-Complexity Pruning

- Searching over all subtrees and selecting using AIC or CV is not possible since there is an exponentially large set of subtrees

→ look at penalized RSS instead

- Define a subtree $T \subset T_0$ to be any tree obtained by pruning T_0

prune = collapse an internal node

and $|T| = \#$ of leaf nodes

$$n_m = |\{x_i \in R_m\}|$$

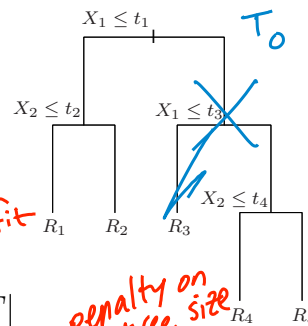
region-specific RSS

$$\hat{\beta}_m = \frac{1}{n_m} \sum_{x_i \in R_m} y_i$$

$$Q_m(T) = \frac{1}{n_m} \sum_{x_i \in R_m} (y_i - \hat{\beta}_m)^2$$

- We examine a complexity criterion

$$C_\lambda(T) = \sum_{m=1}^{|T|} n_m Q_m(T) + \lambda |T|$$



tree fit

penalty on tree size

©Emily Fox 2013

6

Cost-Complexity Pruning

compute for $\hat{\lambda}$ and all trees in sequence

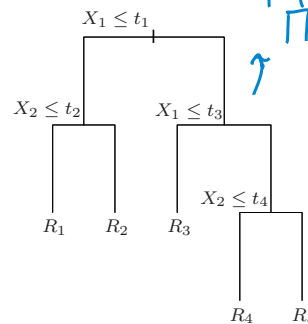
$$C_\lambda(T) = \sum_{m=1}^{|T|} n_m Q_m(T) + \lambda|T|$$

- Can find using *weakest link pruning*
 - Successively collapse the internal node that produces smallest increase in RSS

$$\sum_m n_m Q_m(t)$$

- Continue until at single-node (root) tree
- Produces a finite sequence of subtrees, which must contain T_λ
- See Breiman et al. (1984) or Ripley (1996)

- Choose λ via 5- or 10-fold CV $\rightarrow \hat{\lambda}$
- Final tree: $T_{\hat{\lambda}}$



©Emily Fox 2013

7

Issues

- Unordered categorical predictors
 - With unordered categorical predictors with q possible values, there are $2^{q-1}-1$ possible choices of partition points to consider for each variable
 - Prohibitive for large q
 - Can deal with this for binary y ...will come back to this in "classification"
- Missing predictor values...how to cope?
 - Can discard
 - Can fill in, e.g., with mean of other variables
 - With trees, there are better approaches
 - Categorical predictors: make new category "missing"
 - Split on observed data. For every split, create an ordered list of "surrogate" splits (predictor/value) that create similar divides of the data. When examining observation with a missing predictor, when splitting on that dimension, use top-most surrogate that is available instead

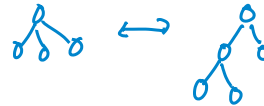
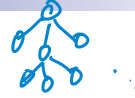
©Emily Fox 2013

8

Issues

- Binary splits

- Could split into more regions at every node
- However, this more rapidly fragments the data leaving insufficient data and subsequent levels
- Multiway splits can be achieved via a sequence of binary splits, so binary splits are generally preferred



- Instability

- Can exhibit high variance
- Small changes in the data → big changes in the tree
- Errors in the top split propagates all the way down
- **Bagging** averages many trees to reduce variance ... more later

- Inference

- Hard...need to account for stepwise search algorithm

©Emily Fox 2013

9

Issues

- Lack of smoothness

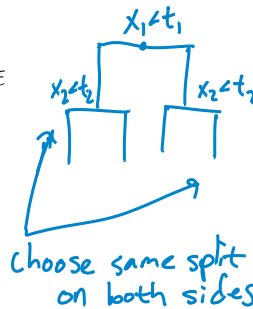
- Fits piecewise constant models...unlikely to believe this structure
- **MARS** address this issue (can view as modification to CART)

↳ later this lecture

- Difficulty in capturing additive structure

- Imagine true structure is
- $$y = \beta_1 I(x_1 < t_1) + \beta_2 I(x_2 < t_2) + \epsilon$$
- No encouragement to find this structure

- hard w/o sufficient data
 - this is just w/ 2 additive effects. Harder to happen or notice w/ more.



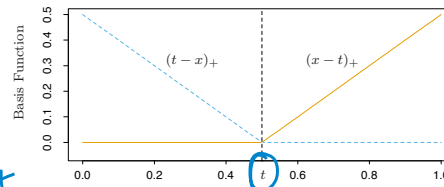
©Emily Fox 2013

10

Multiple Adaptive Regression Splines

- MARS is an adaptive procedure for regression
 - Well-suited to high-dimensional covariate spaces
- Can be viewed as:
 - Generalization of step-wise linear regression
 - Modification of CART
- Consider a basis expansion in terms of piecewise linear basis functions (linear splines)

$(x-t)_+$
 $(t-x)_+$ } "reflected pair"
 ↑ piecewise linear w/ knot @ t



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

11

Multiple Adaptive Regression Splines

- Take knots at all observed x_j

$$C = \{(x_j - t)_+, (t - x_j)_+\}_{j=1, \dots, d}$$

$t \in \{x_{1j}, \dots, x_{dj}\}$

 - If all locations are unique, then $2 \cdot n \cdot d$ basis functions
 - Treat each basis function as a function on x , just varying with x_j

$$h_m(x) = (x_j - t)_+ \text{ for example}$$

- The resulting model has the form

$$f(x) = \beta_0 + \sum_{m=1}^M \beta_m h_m(x) \quad \text{LBE}$$

↑ $h_m \in C$ or their products

- Built in a forward stepwise manner in terms of this basis

©Emily Fox 2013

12

MARS Forward Stepwise

- Given a set of h_m , estimation of β_m proceeds as with any linear basis expansion (i.e., minimizing the RSS)

- How do we choose the set of h_m ?

- Start with $h_0(x) = 1$ and $M=0$
- Consider product of all h_m in current model with reflected pairs in C
 - Add terms of the form

$$\hat{\beta}_{M+1} h_\ell(x) (x_j - t)_+ + \hat{\beta}_{M+2} h_\ell(x) (t - x_j)_+$$

Handwritten notes: $h_\ell \in \mathcal{H}$, $\hat{\beta}_{M+1}, \hat{\beta}_{M+2}$ are est. using LS + all other terms in model

- Increment M and repeat
- Stop when preset M is hit
- Typically end with a large (overfit) model, so backward delete
 - Remove term with smallest increase in RSS
 - Choose model based on generalized CV

©Emily Fox 2013

13

MARS Forward Stepwise Example

general terms: $\hat{\beta}_{M+1} h_\ell(x) (x_j - t)_+ + \hat{\beta}_{M+2} h_\ell(x) (t - x_j)_+$

- At the first stage, add term of form

$$\beta_1 (x_j - t)_+ + \beta_2 (t - x_j)_+$$

with the optimal pair being

$$\hat{\beta}_1 (x_2 - x_{72})_+ + \hat{\beta}_2 (x_{72} - x_2)_+$$

- Add pair to the model and then consider including a pair like

$$\beta_3 h_m(x) (x_j - t)_+ + \beta_4 h_m(x) (t - x_j)_+$$

with choices for h_m being:

$$\begin{aligned} h_0(x) &= 1 \\ h_1(x) &= (x_2 - x_{72})_+ \\ h_2(x) &= (x_{72} - x_2)_+ \end{aligned}$$

The term $(x_1 - x_{51})_+ + (x_{72} - x_2)_+$ considered... looks like

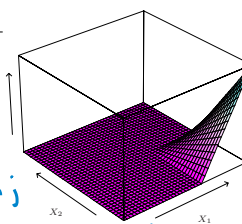


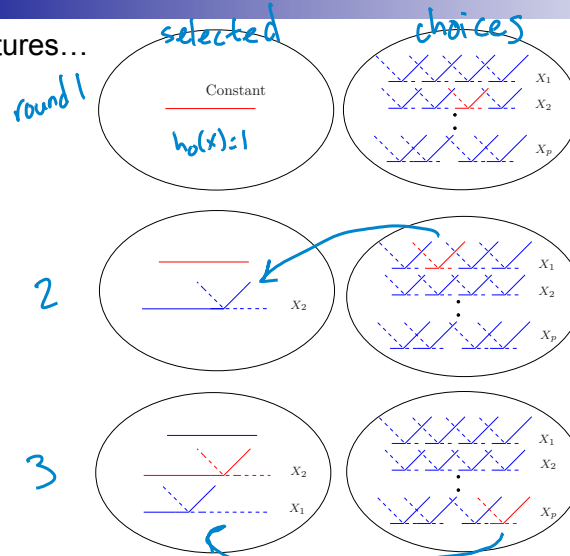
Figure from Hastie, Tibshirani, Friedman book

©Emily Fox 2013

14

MARS Forward Stepwise

- In pictures...



From
Hastie,
Tibshirani,
Friedman
book

©Emily Fox 2013

15

Why MARS?

- Why these piecewise linear basis functions?
 - Ability to operate locally
 - When multiplied, non-zero only over small part of the input space
 - Resulting regression surface has local components and only where needed (spend parameters carefully in high dims)
 - Computations with linear basis are very efficient
 - Naively, we consider fitting n reflected pairs for each input x_j
 - $\rightarrow O(n^2)$ operations
 - Can exploit simple form of piecewise linear function
 - Fit function with rightmost knot. As knot moves, basis functions differ by 0 over the left and by a constant over the right
 - \rightarrow Can try every knot in $O(n)$

©Emily Fox 2013

16

Why MARS?

- Why forward stagewise?
 - Hierarchical in that multiway products are built from terms already in model (e.g., 4-way product exists only if 3-way already existed)
 - Higher order interactions tend to only exist if some of the lower order interactions exist as well
 - Avoids search over exponentially large space

- Notes:
 - Each input can appear at most once in a product...Prevents formation of higher-order powers of an input
 - Can place limit on order of interaction. That is, one can allow pairwise products, but not 3-way or higher.
 - Limit of 1 → additive model

©Emily Fox 2013

17

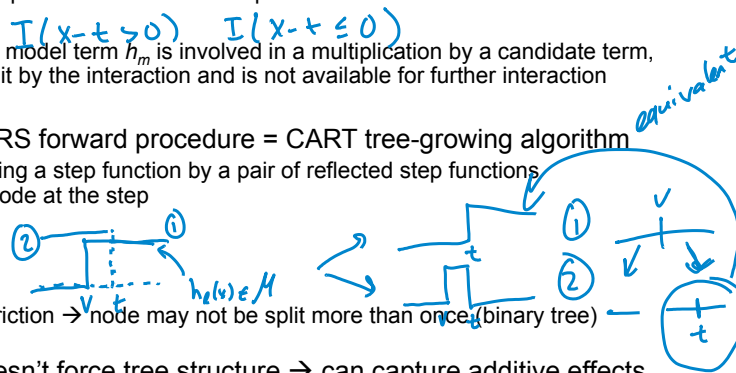
Connecting MARS and CART

- MARS and CART have lots of similarities

- Take MARS procedure and make following modifications:
 - Replace piecewise linear with step functions
 - When a model term h_m is involved in a multiplication by a candidate term, replace it by the interaction and is not available for further interaction

- Then, MARS forward procedure = CART tree-growing algorithm
 - Multiplying a step function by a pair of reflected step functions = split node at the step

- 2nd restriction → node may not be split more than once (binary tree)
- * MARS doesn't force tree structure → can capture additive effects



©Emily Fox 2013

18

What you need to know

- Regression trees provide an adaptive regression method
- Fit constants (or simple models) to each region of a partition
- Relies on estimating a binary tree partition
 - Sequence of decisions of variables to split on and where
 - Grown in a greedy, forward-wise manner
 - Pruned subsequently
- Implicitly performs variable selection
- MARS is a modification to CART allowing linear fits

©Emily Fox 2013

19

Readings

- Wakefield – 12.7
- Hastie, Tibshirani, Friedman – 9.2.1-9.2.2, 9.2.4, 9.4
- Wasserman – 5.12

©Emily Fox 2013

20

Module 4: Coping with Multiple Predictors

A Short Case Study

STAT/BIOSTAT 527, University of Washington

Emily Fox

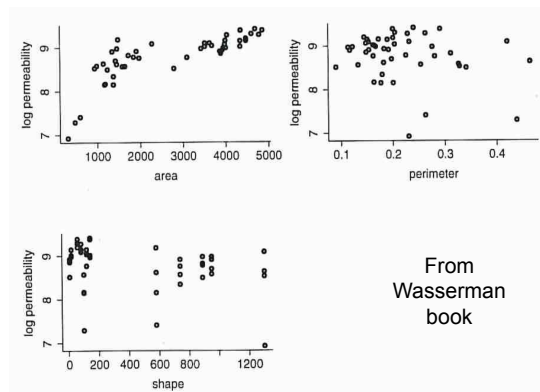
May 21st, 2013

©Emily Fox 2013

21

Rock Data

- 48 rock samples from a petroleum reservoir
- Response = permeability
- Covariates = area of pores, perimeter, and shape



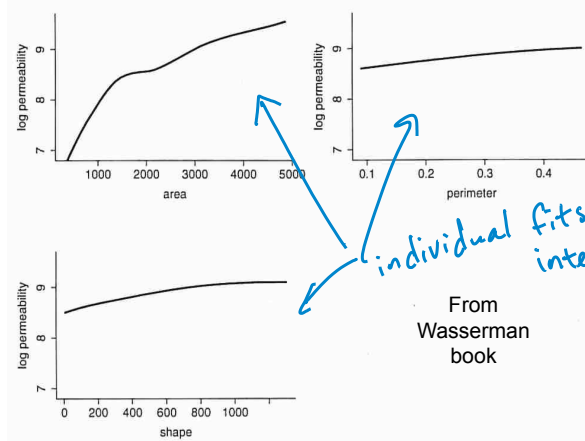
©Emily Fox 2013

22

Generalized Additive Model

- Fit a GAM:

$$\text{permeability} = f_1(\text{area}) + f_2(\text{perimeter}) + f_3(\text{shape}) + \epsilon$$



individual fits are very interpretable

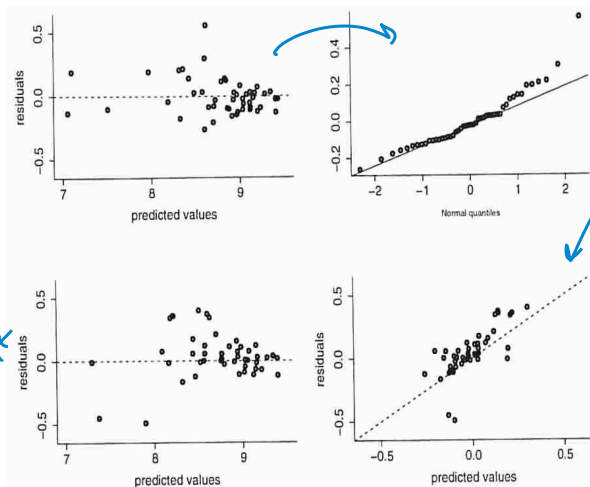
From Wasserman book

©Emily Fox 2013

23

GAM vs. Local Linear Fits

- Comparison to a 3-dimensional local linear fit



*Similar fits
=> GAM is sufficient*

Residuals of GAM vs. local linear
From Wasserman book

©Emily Fox 2013

24

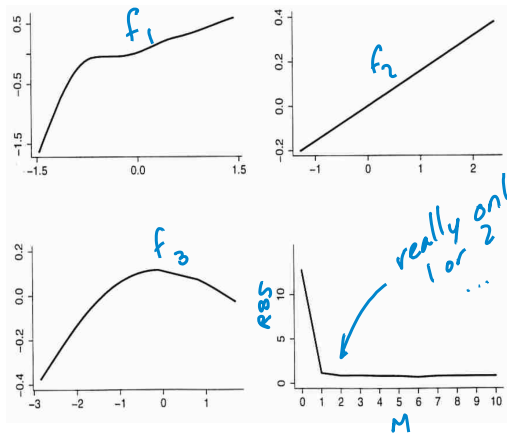
Projection Pursuit

$$f(x_1, \dots, x_d) = \alpha + \sum_{m=1}^M f_m(w_m^T x)$$

- Applying projection pursuit with $M = 3$ yields

$$w_1 = (.99, .07, .08)^T, w_2 = (.43, .35, .83)^T, w_3 = (.74, -.28, -.61)^T$$

$V_1 = 0.99$ area
 +0.07 perim.
 +0.08 shape



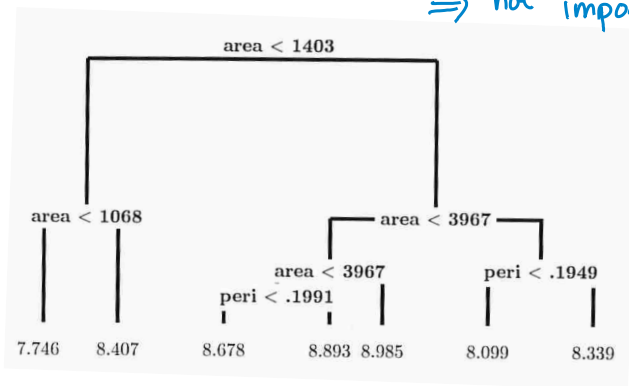
really only need
 1 or 2 terms
 ... capture interactions
 through proj.

From Wasserman book

Regression Trees

- Fit a regression tree to the rock data
- Note that the variable "shape" does not appear in the tree

\Rightarrow not important



From Wasserman book

Module 5: Classification

A First Look at Classification: CART

STAT/BIOSTAT 527, University of Washington

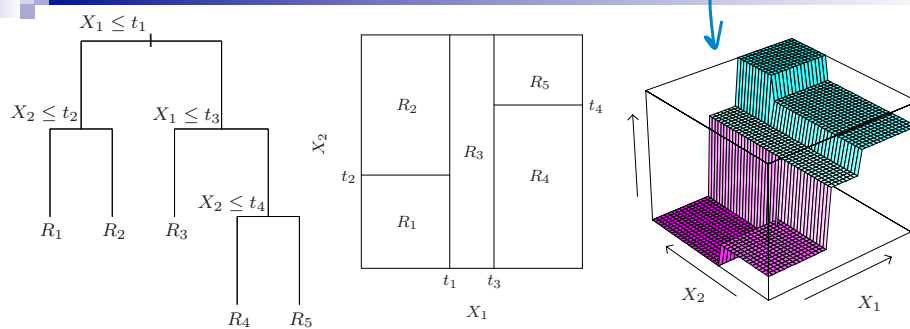
Emily Fox

May 21st, 2013

©Emily Fox 2013

27

Regression Trees



- So far, we have assumed continuous responses y and looked at regression tree models:

$$f(x) = \sum_{m=1}^M \beta_m I(x \in R_m)$$

Figures from Hastie, Tibshirani, Friedman book

©Emily Fox 2013

28

Classification Trees

- What if our response y is **categorical** and our goal is classification?

$y \in \{\text{'email'}, \text{'spam'}\} \rightarrow \{0, 1\}$

More generally, $y \in \{0, \dots, k\} \rightarrow \{1, \dots, k\}$

of classes

- Can we still use these tree structures? **Yes!**

- Recall our **node impurity** measure

$$Q_m(T) = \frac{1}{n_m} \sum_{x_i \in R_m} (y_i - \hat{\beta}_m)^2$$

local RSS to region m

- Used this for growing the tree

$$\min_{j,s} \left[\sum_{x_i \in R_1(j,s)} (y_i - \hat{\beta}_1)^2 + \sum_{x_i \in R_2(j,s)} (y_i - \hat{\beta}_2)^2 \right]$$

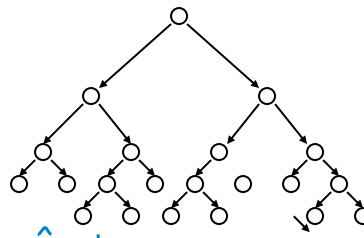
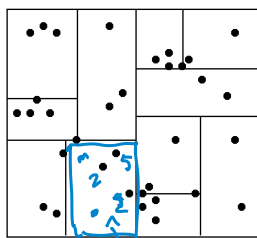
- As well as pruning $C_\lambda(T) = \sum_{m=1}^{|T|} n_m Q_m(T) + \lambda|T|$

- Clearly, squared-error is not the right metric for classification

©Emily Fox 2013

29

Classification Trees



- First, what is our decision rule at each leaf?

- Estimate probability of each class given data at leaf node:

$$\hat{p}_{mk} = \frac{1}{n_m} \sum_{x_i \in R_m} \mathbb{1}(y_i = k)$$

- Majority vote:

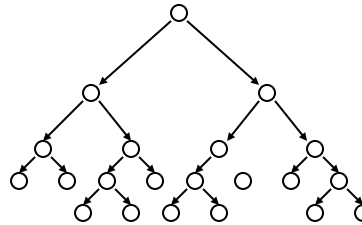
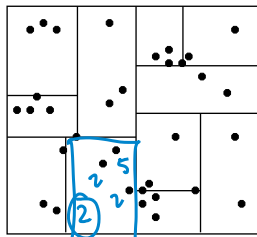
$$k(m) = \operatorname{argmax}_k \hat{p}_{mk}$$

Figures from Andrew Moore kd-tree tutorial

©Emily Fox 2013

30

Classification Trees



■ How do we measure **node impurity** for this fit/decision rule? $Q_m(\tau)$

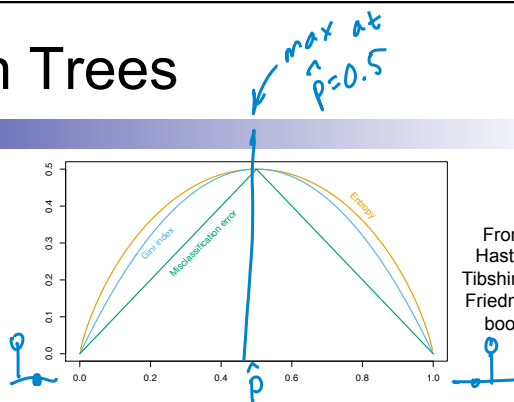
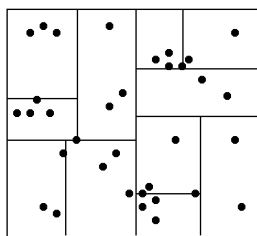
- Misclassification error: $\frac{1}{n_m} \sum_{i \in R_m} \mathbb{1}(y_i \neq k(m)) = 1 - \hat{p}_{mk(m)}$
 \uparrow in ex. 2
- Gini index: $\sum_k \sum_{k' \neq k} \hat{p}_{mk} \hat{p}_{mk'} = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk})$
- Cross-entropy or deviance: $-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$

Figures from Andrew Moore kd-tree tutorial

©Emily Fox 2013

31

Classification Trees



From Hastie, Tibshirani, Friedman book

■ How do we measure **node impurity** for this fit/decision rule?

- Misclassification error (K=2): $1 - \max(\hat{p}, 1 - \hat{p})$ $\hat{p} = \text{prop. in class 2}$
- Gini index (K=2): $2\hat{p}(1 - \hat{p})$
- Cross-entropy or deviance (K=2): $-\hat{p} \log \hat{p} - (1 - \hat{p}) \log(1 - \hat{p})$
 largest at $\hat{p} = 1/2$... lots of uncertainty

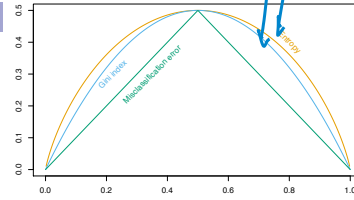
©Emily Fox 2013

32

Notes on Impurity Measures

■ Impurity measures

- Misclassification error: $1 - \hat{p}_{mk(m)}$
- Gini index: $\sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$
- Cross-entropy or deviance: $-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$



From Hastie, Tibshirani, Friedman book

■ Comments:

- Differentiability *Gini + cross-entropy*
- Sensitivity to changes in node probabilities

1 400 | 400 ... 0 truth ← misclass. = 0

300 | 100 | 100 | 300 ← misclass. = 0.25

300 | 400 | 200 ← misclass. = 0.25

pure node, so prefer this Gini + entropy are lower here

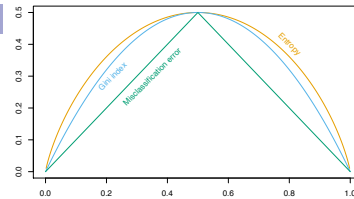
- Often use Gini or cross-entropy for growing tree, and misclass. for pruning

can use any

Notes on Impurity Measures

■ Impurity measures

- Misclassification error: $1 - \hat{p}_{mk(m)}$
- Gini index: $\sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$
- Cross-entropy or deviance: $-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$



From Hastie, Tibshirani, Friedman book

■ Other interpretations of Gini index:

- Instead of majority vote, classify observations to class k with prob. \hat{p}_{mk}

$$\text{Error} = \sum_k \sum_{k'} P(g(x) = k) P(+ (x) = k') \mathbb{1}(k \neq k')$$

↑ classifier ↑ truth

$$\text{training error} = \sum_k \sum_{k' \neq k} \hat{p}_{mk} \hat{p}_{mk'}$$

- Code each observation as 1 for class k and 0 otherwise
 - Variance: *1 against all* $\hat{p}_{mk}(1 - \hat{p}_{mk})$

- Summing over k gives the Gini index

Classification Tree Issues

- Unordered categorical predictors
 - With unordered categorical predictors with q possible values, there are $2^{q-1}-1$ possible choices of partition points to consider for each variable
 - For binary (0-1) outcomes, can order predictor classes according to proportion falling in outcome class 1 and then treat as ordered predictor
 - Gives optimal split in terms of cross-entropy or Gini index
 - Also holds for quantitative outcomes and square-error loss...order predictors by increasing mean of the outcome
 - No results for multi-category outcomes
- Loss matrix
 - In some cases, certain misclassifications are worse than others
predicting no disease when disease
 - Introduce **loss matrix** ...more on this soon
 - See Tibshirani, Hastie and Friedman for how to incorporate into CART

©Emily Fox 2013

35

Classification Tree Spam Example

- Example: *predicting spam*
 - Data from UCI repository *0* *1*
 - Response variable: *email* or *spam*
 - 57 predictors:
 - 48 quantitative – percentage of words in email that match a give word such as “business”, “address”, “internet”,...
 - 6 quantitative – percentage of characters in the email that match a given character (; , [! \$ #)
 - The average length of uninterrupted capital letters: CAPAVE
 - The length of the longest uninterrupted sequence of capital letters: CAPMAX
 - The sum of the length of uninterrupted sequences of capital letters: CAPTOT
- Looked at this w/ GAMs*

©Emily Fox 2013

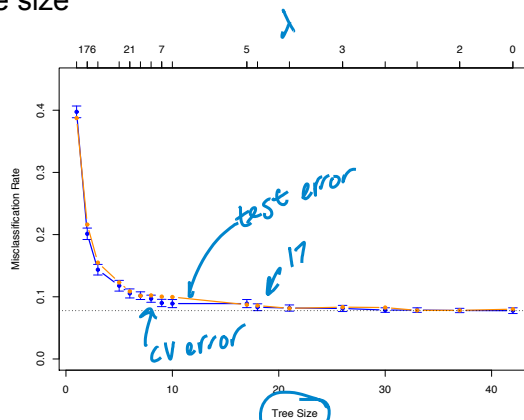
36

Classification Tree Spam Example

- Used cross-entropy to grow tree and misclassification to prune

- 10-fold CV to choose tree size

- CV indexed by λ
- Sizes refer to $|T_\lambda|$
- Error rate flattens out around a tree of size 17



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

37

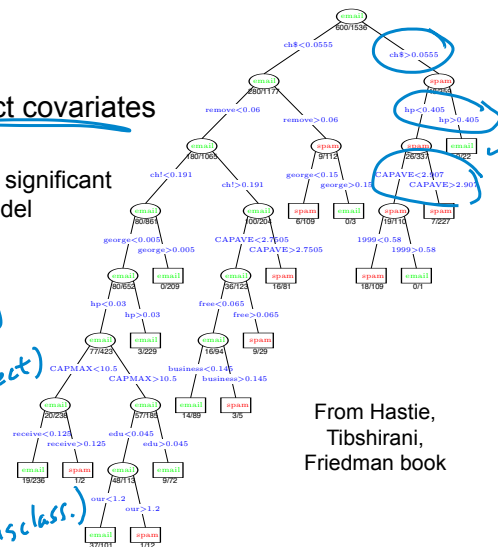
Classification Tree Spam Example

- Resulting tree of size 17

- Note that there are 13 distinct covariates split on by the tree

- 11 of these overlap with the 16 significant predictors from the additive model previously explored

$\$ > 5.59 \rightarrow$ spam warning
 if hp is freq., then 'email' (22 cases correct)
 if not, and CAPAVE > 2.9 then 'spam' (7 of 227 misclass.)



From Hastie, Tibshirani, Friedman book

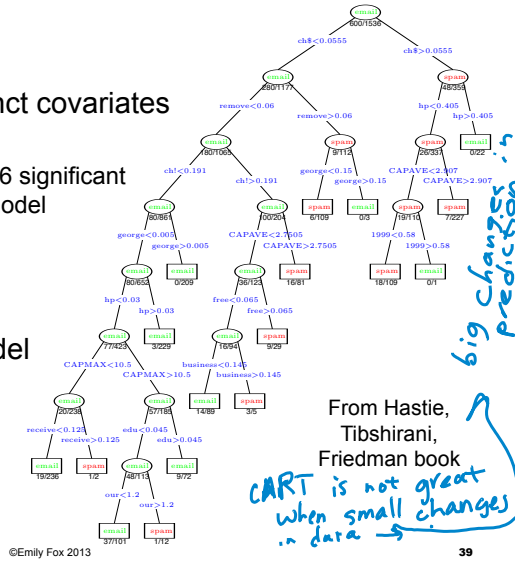
©Emily Fox 2013

38

Classification Tree Spam Example

- Resulting tree of size 17
- Note that there are 13 distinct covariates split on by the tree
 - 11 of these overlap with the 16 significant predictors from the additive model previously explored
- Overall error rate (9.3%) is higher than for additive model

True	Predicted	
	email	spam
email	57.3%	4.0%
spam	5.3%	33.4%



©Emily Fox 2013

39

What you need to know

- Classification trees are a straightforward modification to the regression tree setup
- Just need new definition of node impurity for growing and pruning tree
- Decision at the leaves is a simple majority-vote rule

©Emily Fox 2013

40

Readings

- Wakefield – 10.3.2, 10.4.2, 12.8.4
- Hastie, Tibshirani, Friedman – 9.2.3, 9.2.5, 2.4