

## Moving beyond linearity:

- polynomial regression
- step functions
- regression splines
- smoothing splines
- local regression
- generalized additive model

### Polynomial regression

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_d x_{id} + \epsilon_i$$

usually up until degree 3-4, after too flexible

can calculate  $\text{Var } \hat{f}(x_0)$ , as least squares gives a variance for each  $\hat{\beta}_i$   
as well as the covariances  $\rightarrow$  if  $\hat{C}$  is the covariance matrix,

$$\text{and } \hat{L}_0 = [1, x_0, x_0^2, \dots, x_0^d], \text{ then } \text{Var } \hat{f}(x_0) = \hat{L}_0^\top \hat{C} \hat{L}_0$$

to plot the 95% confidence interval, we calculate

the pairwise standard error and double it ( $\sim$  Normal assumption) on both sides of  
the curve

- if we have a binary  $y$ , e.g.  $y > 250$  or not are content if as logistic regression

$$\text{e.g. } p(y > 250 | x_i) = \frac{\exp(\beta_0 + \beta_1 x_{i1} + \dots + \beta_d x_{id})}{1 + \exp(\beta_0 + \dots + \beta_d x_{id})}$$

## 2) Step functions

break  $x$  into bins, fit a different constraint on all

~ this converts  $x$  into an ordered categorical variable

creates cut points  $c_1, \dots, c_k$  and construct  $k+1$  variables  $C_0, \dots, C_k$

$$\begin{aligned} C_0(x) &= I(x \leq c_0) \\ C_1(x) &= I(c_0 \leq x < c_1) \\ &\vdots \\ C_k(x) &= I(c_k \leq x) \end{aligned}$$

and then fit

$$y_i = \beta_0 + \beta_1 C_1(x_i) + \beta_2 C_2(x_i) + \dots + \beta_k C_k(x_i) + \varepsilon_i$$

not using  $C_0$  as redundant w/ the intercept

$\beta_0$ : mean value of  $y$  for  $x < c_0$ ,

$\beta_j$ : mean increase in  $y$  for  $c_j \leq x < c_{j+1}$  relative to  $x < c_j$

Basis function:  $y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_k b_k(x_i) + \epsilon_i$

→ poly regression:  $b_j(x) = x^j$

→ step function:  $b_j(x) = I(c_j \leq x_i < c_{j+1})$

### 3, regression splines

a) piece-wise low-degree polynomials

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \epsilon_i \quad \text{where}$$

$\beta_0, \beta_1, \beta_2$  and  $\beta_3$  differ in different ranges

parts where coeffs change are called knots

e.g. we fit a new poly b/w every knot

b) we can impose continuity at the knots

- continuity of fn

- continuity of derivatives → smoothes out the fn

- usually up until  $(d-1)$ th derivative

degrees of freedom:  $(\underbrace{\text{degree} + 1}_{\text{intercept}}) * (\underbrace{\text{knots} + 1}_{(\text{degree} - 1)}) - (\underbrace{\text{constraint on derivatives}}_{''}) * \text{knots}$

e.g. cubic spline w/  $k$  knots has  $4k+4 - 3k = k+4$  degrees of freedom.

basis representation for splines:

Cubic spline has  $k+4$  degrees of freedom, w/  $b_1, \dots, b_{k+3}$  ( $b_0$  replaces  $b_0$ )

$$h(x, \varepsilon) = (x - \varepsilon)_+^3 = \begin{cases} (x - \varepsilon)^3 & \text{if } x > \varepsilon \\ 0 & \text{otherwise} \end{cases}$$

This is continuous up until the 3rd derivative, where it will be discontinuous

so we fit  $x, x^2, x^3, h(x, \varepsilon_1), h(x, \varepsilon_2), \dots, h(x, \varepsilon_k)$   
where  $\varepsilon_i$  are the knots

$\rightarrow k+4$  degrees of freedom

natural spline: we require linearity at the two outer edges

as the variance there is usually big

(this is counted as 2 knots at the boundaries w/ another constraint to enforce linearity)

choosing the location of knots

usually at uniform quantiles of the data

choose degrees of freedom and knot number w/ CV

→ Splines are usually more stable than poly regression

less error at edges if natural

don't need  $\Rightarrow$  high degree polys

can be very flexible where data changes rapidly  
and less flexible where data more stable  
by placement of knots

#### 4) Smoothing splines

→ we want to minimize RSS but minimize overfitting by having a smooth function.

$$\text{Cost fn} = \underbrace{\sum (y_i - g(x_i))^2}_{\text{RSS loss}} + \lambda \underbrace{\int g''(t)^2 dt}_{\text{penalty}} \xrightarrow{\text{so they pass}}$$

penalizing the change in the slope

if  $\lambda=0$ , exactly fitting to the points,

if  $\lambda \rightarrow \infty$ , straight line

controls effective degrees of freedom

the function that minimises this cost fn is:

- piece-wise cubic polys w/ lots at  $x_1, \dots, x_n$  w/ continuous derivatives
- linear at the outer regions

$\Rightarrow$  THAT'S A NATURAL CUBIC SPLINE

but not the same one we got w/ the method before  
as we have  $\lambda$  controlling the shrinkage

- LOOCV is very easily computed for smoothing splines  
(or any basis for regression)

### S<sub>3</sub>: Local regression

- only use nearby datapoints,

Algo:  $\left\{ \begin{array}{l} 1, \text{ gather the fraction } s = \frac{k}{n} \text{ points who are closest to } x_0 \\ 2, \text{ assign } K_{i0} = k(x_i, x_0) \text{ to each point of the nbhd} \\ \quad \text{so closest has largest weight, furthest has smallest weight} \\ 3, \text{ fit minimizing } \sum K_{i0} (y_i - \beta_0 - \beta_1 x_i)^2 \end{array} \right.$   
for  $x=x_0$ :

4, Fitted value is  $\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0$

→ we need to do this for all datapoints

Choices :  $\beta_i k(x_i, x_0)$   
 $\downarrow$

large  $\rightarrow$  less wiggly       $0 \leq \beta \leq 1$   
 small  $\rightarrow$  more wiggly

- for multiple predictors, for excep<sup>t</sup>s can fit local in one variable, global in others (e.g. local in time, global in space)  
 $\Rightarrow$  varying coefficient model
- can use 2d / 3d neighbourhoods for spatial data
- will perform locally above-threshold as points are far away

## 5, Generalized additive models (GAM)

$$y_i = \beta_0 + \sum_{j=1}^p f_j(x_{ij}) + \varepsilon_i$$

where  $f_i(x)$  is a smooth nonlinear fn

Each predictor has a separate  $f_i$

→ Building blocks

e.g. can fit one predictor w/ splines, the other w/  
poly regression, or local regression

one big regression fit w/ least squares

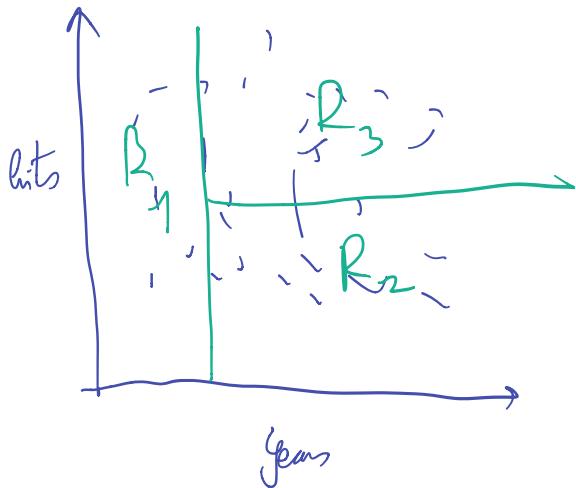
{ under smoothing splines, then w/ backfitting methods  
but difference b/w smoothing and natural splines is small}

Pros: • can use non-linear  $f_i$  for each  $x_j$   
so model non-linear relationships w/o  
fixing any variable interactions

- make accurate predictions possible b/c non-linear
- can see response to one  $x_j$  by holding all other variables fix  
→ useful for inference

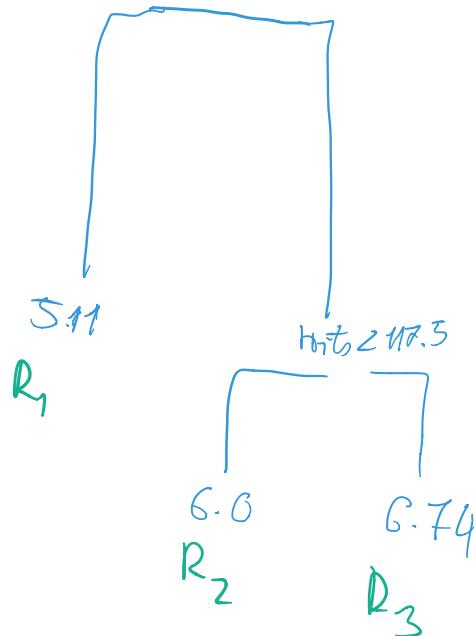
- smoother of  $f_{ij}$  can be summarised by degree of freedom
- Can add interaction terms manually eg.  $f_{ij}(x_j, x_k)$   
these can be fit w/ 2d smoothes

## Chapter 8: Tree-Based Methods



want to predict salary  
given years and bits  
(log-transformed already)

Years < 4.5



$$R_1 = \{x | \text{Years} \leq 4.5\}$$

$$R_2 = \{x | \text{Years} \leq 4.5, \text{bits} \leq 117.5\}$$

$$R_3 = \{x | \text{Years} > 4.5, \text{bits} > 117.5\}$$

$R_x$  : terminal nodes or leaves

internal nodes

branches

- divide space into regions
- predict the mean of observations in that region

How to find the regions?

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

use greedy approach "binary tree splitting" top-down

- Choose  $j$  and  $s$  such that regions  $\{\underbrace{x|x_j < s}\}_{R_1(j,s)} \text{ and } \{\underbrace{x|x_j \geq s}\}_{R_2(j,s)}$  minimize RSS

$$\text{e.g. } \sum_{i \in x_i \in R_1(j,s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i \in x_i \in R_2(j,s)} (y_i - \hat{y}_{R_2})^2$$

- Repeat the process by splitting one of the previously splitted regions

Stop at a stopping criterion (regions, # of observation in a region etc.)

- Predict my mean in that region

## Tree-pruning

Above is likely to overfit

→ want smaller trees, we will cut some branches after creating the tree

Algo:

- 1 Recurive binary splitting until each terminal node has fewer than some min. number of observations
- 2 Apply cost complexity pruning w/ tuning parameter  $\alpha$
- 3 Choose  $\alpha$  w/  $C_V$

## Complexity pruning

For each  $\alpha \in \mathbb{R}$  subtree  $T \subset T_0$  if

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

↑ # of terminal nodes

$\alpha = 0 \Rightarrow T = T_0$  original tree

$\alpha \rightarrow \infty \Rightarrow$  fewer terminal nodes (from large to less)

Nice thing: as we increase  $X$ , branches get cut in a predictable and needed way

Can make a list of branches cut as  $X \rightarrow \infty$   
subtrees left

can use CW to find right  $X$

### Classification trees

same as regression tree but for qualitative variables

- instead of mean, we use most common class for prediction
- instead of RSS, misclassification rate  
could be used but is NOT sensitive enough for pruning

We use the Gini index instead:

$\hat{P}_{mis}$ : portion of training samples in the with class  $k$  but are from the  $j$ th class

$$G = \sum_{k=1}^K \hat{P}_{mis} (1 - \hat{P}_{mis}) , \text{ a measure of total variance across classes}$$

small if all  $\hat{p}_{ms}$  are close to 0 or 1

measure of purity

- node contains observations mostly  
from a single class or small

Entropy:

$$D = - \sum_{s=1}^k \hat{p}_{ms} \log \hat{p}_{ms}$$

$D \approx 0$  if  $\hat{p}_{ms} \rightarrow 0$  or 1

$\Rightarrow$  same as Gini index

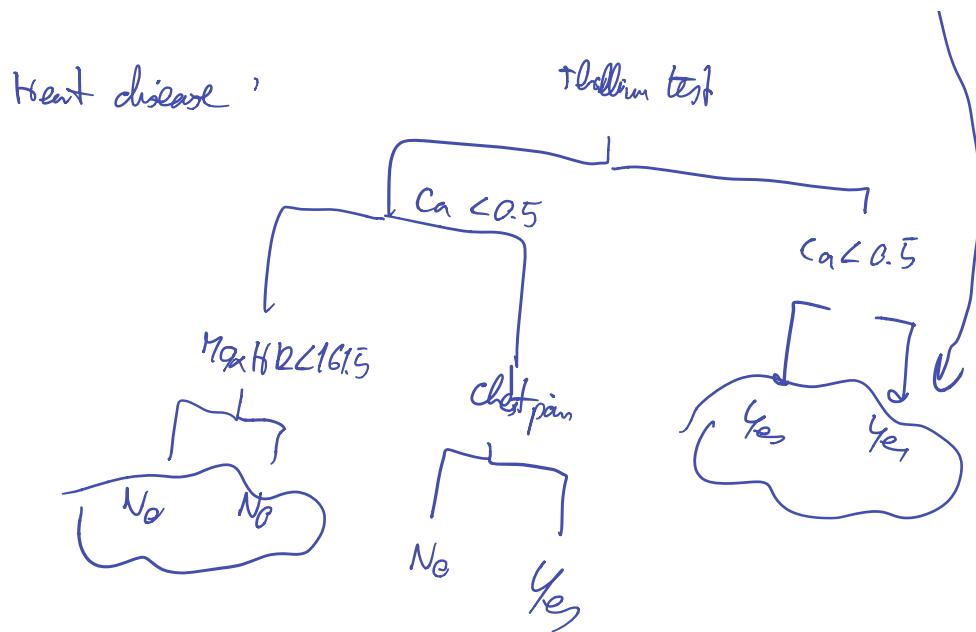
Gini index and entropy are more sensitive to node purity  
than misclassification error rate

Can we collect for pruning

Weird thing: After pruning, we can end up w/ nodes that  
have the same answer at the end nodes.

This is OK if increases node purity (doesn't change misclassification)  
ever

e.g. one node is 100% Yes, the other only 7/11  
yes  $\Rightarrow$  increases overall Gini score  $\uparrow$

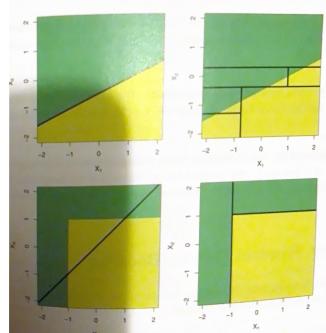


## Tree vs. Linear Models

- depends if data has linear structure or complex relationships
- trees are easier to interpret

### Advantages:

- easy to explain
- mirror human decision-making
- can be displayed graphically
- can easily handle qualitative predictors



tree boundary is linear on top  
tree boundary is regional  
on bottom row

### Disadvantages

- less accurate

- very non-robust, small changes result in a different tree

## Bagging

Trees have high variance and are not stable

the train-test split heavily influences the outcome

bagging = bootstrap aggregation

→ used to reduce variance by bootstrapping training sets and taking the mean of the predictions (or majority for classification)  
 we don't prune these trees → high variance  
 low bias

## Ensemble

On average, each bagged tree uses  $2/3$ rd of the total observations

↳ "out of Bag"

for each observation we take those trees for which it was out of bag ( $\approx 1/3$  of trees) and average their predictions to get a final

prediction. We can see the test error from the errors of these final predictions.

(Cheaper than CV)

Bagging makes us lose interpretability at the expense of accuracy.

How to find important variables:

For regression tree, RSS: total amount the RSS has decreased due to splits over a give predictor

for classification tree, Gini: if large, important predictor  
low the Gini decreased due to splits over a give predictor

## Random Forest

- de-correlates trees

- when splitting your tree, you can only choose one of  $m$  randomly selected variables at each split
  - usually  $m \approx \sqrt{p}$
- you won't always pick the same strongest variable,  
so your trees will look more different from each other  
and be less correlated
- average of trees less variable because of the less correlation
- small  $m$  is good when large number of predictors

## Boosting

growing trees sequentially, using information from previous trees

- Algo:
- 1) set  $\hat{f}(x) = 0$  and  $r_i = y_i - \hat{f}(x)$  in training set
  - for regression trees
    - 3) for  $b=1, 2, \dots, B$ :
      - a) fit tree  $\hat{f}^b(x)$  w/  $c$  splits to training data  $(x_i, r_i)$
      - b) update  $\hat{f}$  by adding a shrinked version of the new tree  

$$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x)$$
      - c) update the residuals  

$$r_i \leftarrow r_i - \lambda \hat{f}^b(x_i)$$
    - 3) output boosted model  

$$\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x)$$

→ learns slowly, we fit the trees to the residuals of the current model instead of the outcome.

we then add this decision tree to the fitted fn and update

the variables.

Param:  $B$  - number of trees

if too large, can overfit (inhibit bagging and rf)

$\lambda$  - shrinkage parameter

usually 0.01 or 0.001

if very small, we need a large  $B$

$d$  - number of splits

$d=1$  can work well (stump trees)  $\Rightarrow$  additive model

controls the interaction -  $d$  splits can involve  
at most  $d$  variables

## Exercises Ch7

Slow

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 (x-\epsilon)^3$$

$\hookrightarrow$  a cubic regression spline, (ie  $x, x^2, x^3, (x-\epsilon)^3$  basis functions)

a) Find cubic poly  $f_1(x)$  if  $f(x) = f_1(x) \quad \forall x \leq \epsilon$ ,

$$(f_1(x) = a_1 + c_1 x + e_1 x^2 + d_1 x^3)$$

$$f_1(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

b)  $f_2(x) \neq f_1(x)$  if  $f(x) = f_2(x) \quad \forall x > \epsilon$ :

$$\begin{aligned} f_2(x) &= \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 (x-\epsilon)^3 \\ &= (\beta_0 - \beta_4 \epsilon^3) + (\beta_1 + \beta_4 3\epsilon^2)x + (\beta_2 - \beta_4 3\epsilon^3)x^2 + (\beta_3 + \beta_4)x^3 \end{aligned}$$

c)  $f_1(\epsilon) = f_2(\epsilon)$  :  $f(x)$  is cont.

$$\beta_4 (\epsilon - \epsilon)^3 = 0 \text{ at } x = \epsilon$$

d)  $f'_1(\epsilon) = f'_2(\epsilon)$  :  $f'(x)$  is cont.

$$3\beta_4 (\epsilon - \epsilon)^2 = 0 \text{ at } x = \epsilon$$

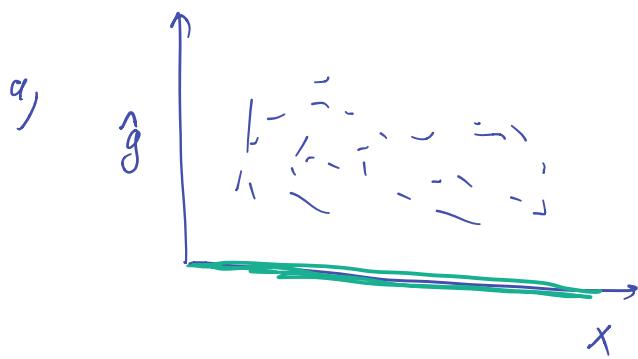
e)  $f''_1(\epsilon) = f''_2(\epsilon)$  :  $f''(x)$  is cont.

$$6\beta_4 (\epsilon - \epsilon) = 0 \text{ at } x = \epsilon$$

2)

$$\hat{g} = \arg \min_g \left( \sum (y_i - g(x_i))^2 + \lambda \int [g''(x)]^2 dx \right)$$

$\Rightarrow$  with derivative instead of road in smooth regression cost



$$\lambda = \infty, m=0$$

$$\Rightarrow \hat{g} = 0$$

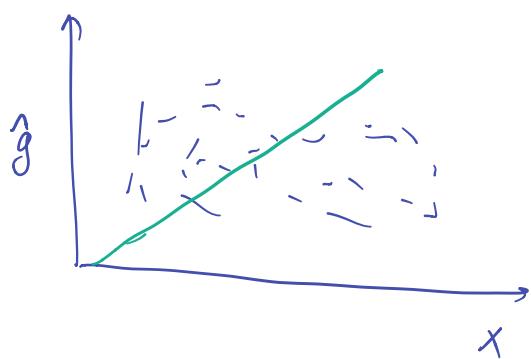
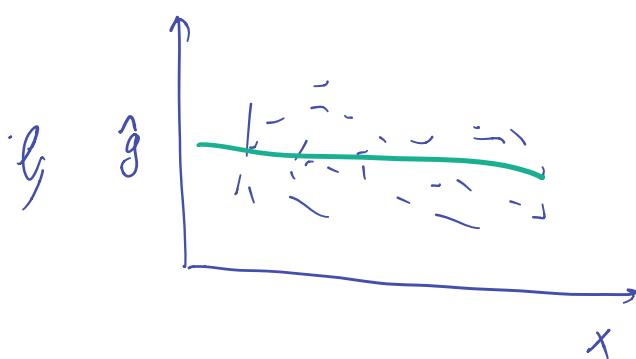
$$\text{minimizes } \int g^2 dx$$

$$\lambda = \infty, m=1$$

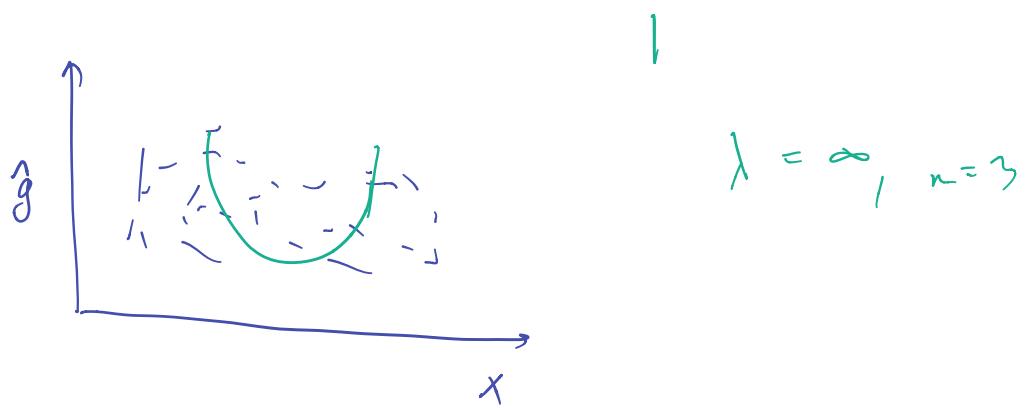
$$\text{minimizes } \int (g')^2 dx$$

$$g' = 0 \rightarrow$$

$$g(x) = C$$



$$\lambda = \infty, m=2$$



$$3) \quad l_1(x) = x$$

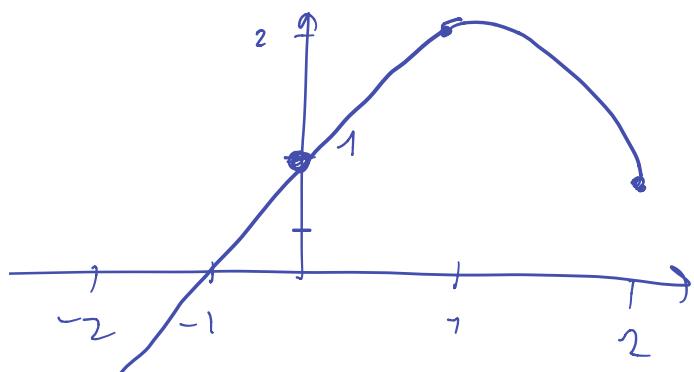
$$l_2(x) = (x-1)^2 \mathbb{I} (x \geq 1)$$

$$y = \beta_0 + \beta_1 l_1(x) + \beta_2 l_2(x) + \varepsilon$$

$$\hat{\beta}_0 = 1$$

$$\hat{\beta}_1 = 1$$

$$\hat{\beta}_2 = -2$$



$$y = \begin{cases} 1 + x - 2(x-1)^2 & \text{if } x \geq 1 \\ 1 + x & \text{if } x < 1 \end{cases}$$

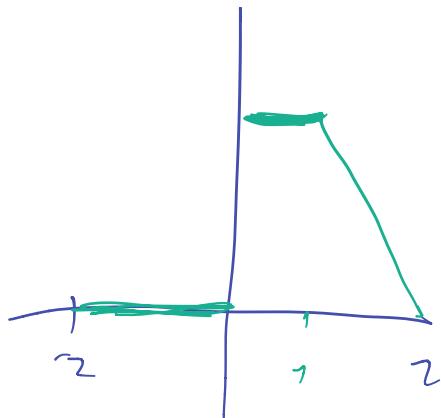
4)  $b_1(x) = \mathbb{I}(0 \leq x \leq 2) - (x-1)\mathbb{I}(1 \leq x \leq 2)$   
 $b_2(x) = (x-3)\mathbb{I}(3 \leq x \leq 4) + \mathbb{I}(4 < x \leq 5)$

$$y = \beta_0 + \beta_1 b_1(x) + \beta_2 b_2(x) + \varepsilon$$

$$\hat{\beta}_0 = 1$$

$$\hat{\beta}_1 = 1$$

$$\hat{\beta}_2 = 3$$



$$y = 1 \quad \text{for } x < 0$$

$$y = 1 + 1 = 2 \quad \text{for } 0 \leq x < 1$$

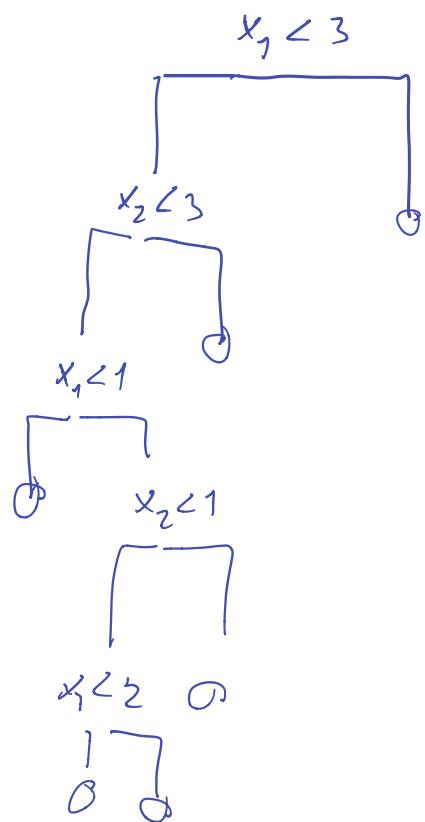
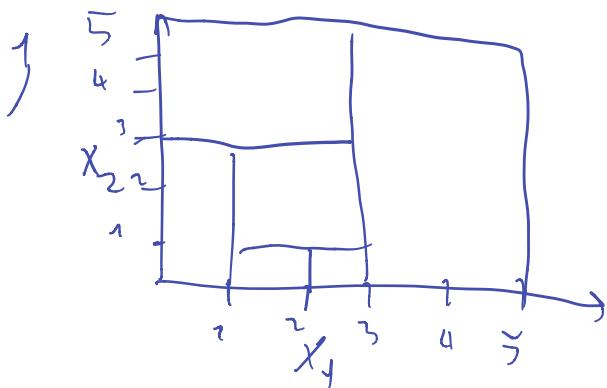
$$\begin{aligned} y &= 1 - (x-1) && \text{for } 1 \leq x < 2 \\ &= 3-x \end{aligned}$$

5) a)  $\hat{g}_2$  have null RSS on training set  
→ higher order poly

b) either, but probably  $\hat{g}_1$  as overfits less

c)  $\hat{g}_1 = \hat{g}_2$  if  $\lambda = 0$

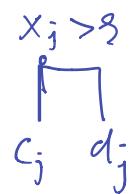
### Chapter 8 exercises



3) Steps is an additive model

$$\hat{f}(x) = \sum \lambda_j \hat{f}_j(x)$$

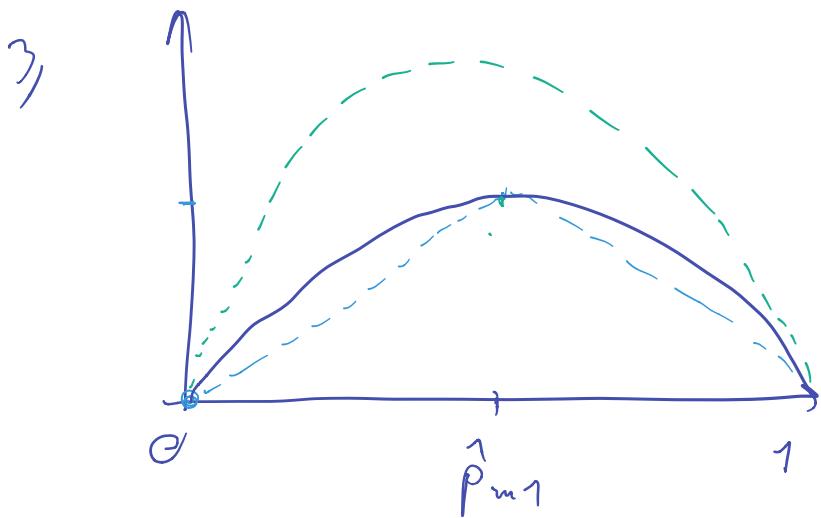
Step:



where

$$\hat{f}_j(x) = c_j \mathbb{I}(x_j < \xi_j) + d_j \mathbb{I}(x_j \geq \xi_j)$$

so  $\lambda_j \hat{f}_j(x) = f(x_j)$  for some  $j$ , if only depends on one predictor



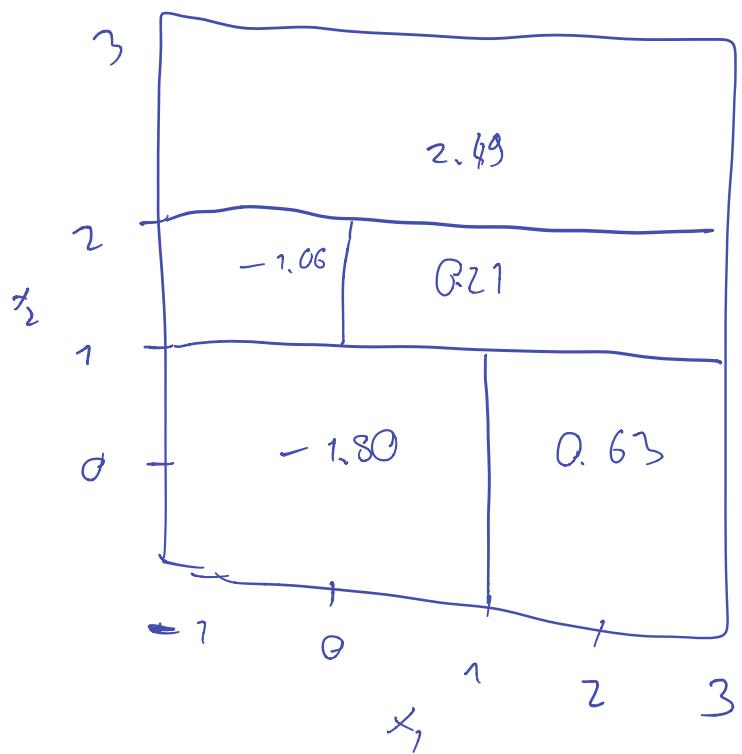
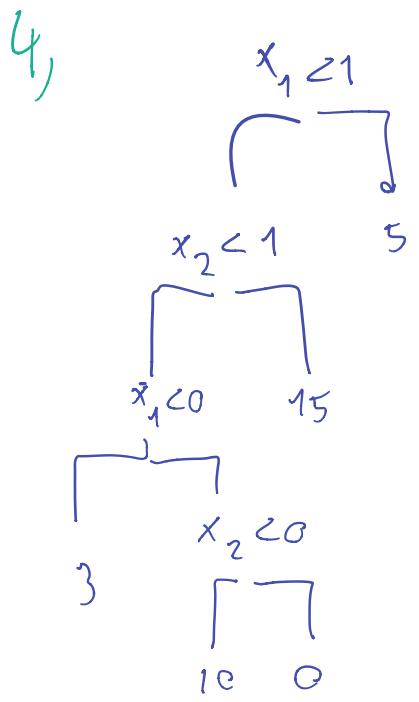
Gini  
entropy  
Classification error

two classes  
 $\hat{p}_{m2} = 1 - \hat{p}_{m1}$

$$\text{Gini} : 2\left(\hat{p}_{m1}(1 - \hat{p}_{m1})\right)^{\otimes}$$

$$\text{Entropy} : -\hat{p}_{m1} \log \hat{p}_{m1} - (1 - \hat{p}_{m1}) \log(1 - \hat{p}_{m1})$$

$$\text{Classification error} : 1 - \max(\hat{p}_{m1}, \hat{p}_{m2})$$



5) Majority vote: Ned

Arg:  $0.43 \rightarrow$  Greens

6)