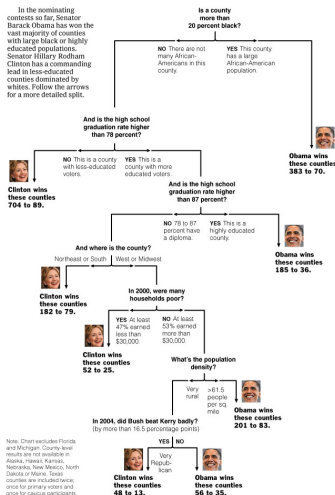


## Chapter 6

# Generalizations of regression

# REGRESSION TREES

## Decision Tree: The Obama-Clinton Divide



**Note:** a simple linear regression is too restrictive for large data sets.

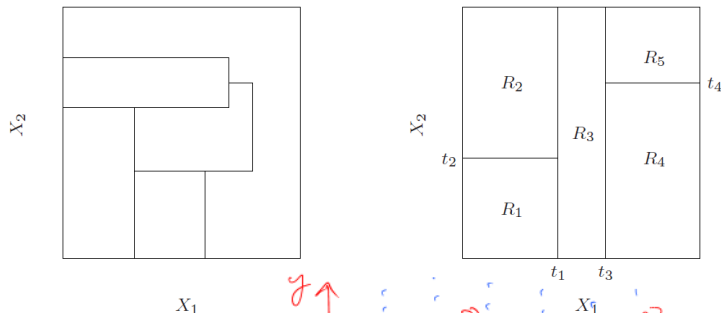
Regression trees offer a flexible technique with results, which are easy to interpret

Sources: Election results via The Associated Press; Census Bureau; Dave Leip's Atlas of U.S. Presidential Elections

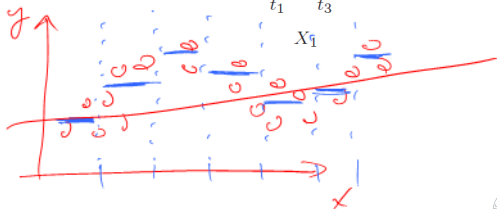
ADAM OLSZEWSKI  
THE NEW YORK TIMES

## General strategy:

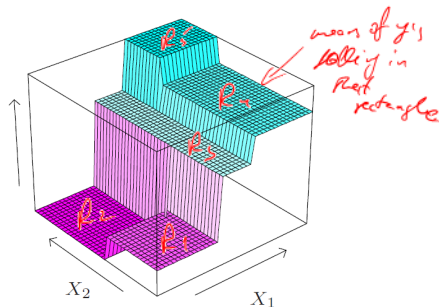
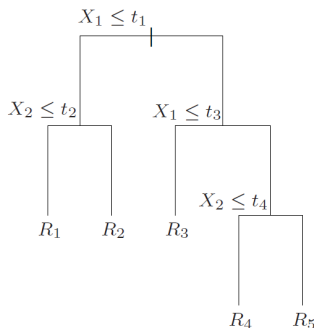
- The values of the explanatory variables are split into  $P$  disjunct regions (rectangles)  $R_1, \dots, R_P$ : **binary splitting**



Source: Hastie et al. (2001)



- In each rectangle we fit a simple model  
e.g. a constant, i.e. the forecast in rectangle  $R_p$  is the mean of all  $Y$ -values falling into this rectangle.



Source: Hastie et al. (2001)

**Question:** how to determine the regions?

OLS method:

$$\sum_{i=1}^n (y_i - \mathbf{x}_i' \mathbf{b})^2 \rightarrow \min, \quad \text{bzgl. } \mathbf{b}.$$

For regression trees:

*over rectangles*  $\sum_{p=1}^P \sum_{i \in R_p} (y_i - \hat{y}_{R_p})^2 \rightarrow \min, \quad \text{w.r.t. } R_1, \dots, R_P,$  *over observations within rectangles* *nothing new.* *we minimize w.r.t to rectangles (splitting point + order) and # of rectangles*

where  $\hat{y}_{R_p}$  is the mean of observations in the  $p$ -th rectangle.

**Note:** direct optimization is hardly possible  $\leadsto$  recursive binary splitting

$$\sum (x_i - a)^2 \rightarrow \min. \Rightarrow \hat{a} = \bar{x}.$$

## Step 1

- Find the variable  $X_j$  and the splitting point  $s$ , which separates the space into two regions:

$$R_1(j, s) = \{\mathbf{X} | X_j \leq s\} \quad \text{and} \quad R_2(j, s) = \{\mathbf{X} | X_j > s\}.$$

- $j$  and  $s$  are determined using the following objective function

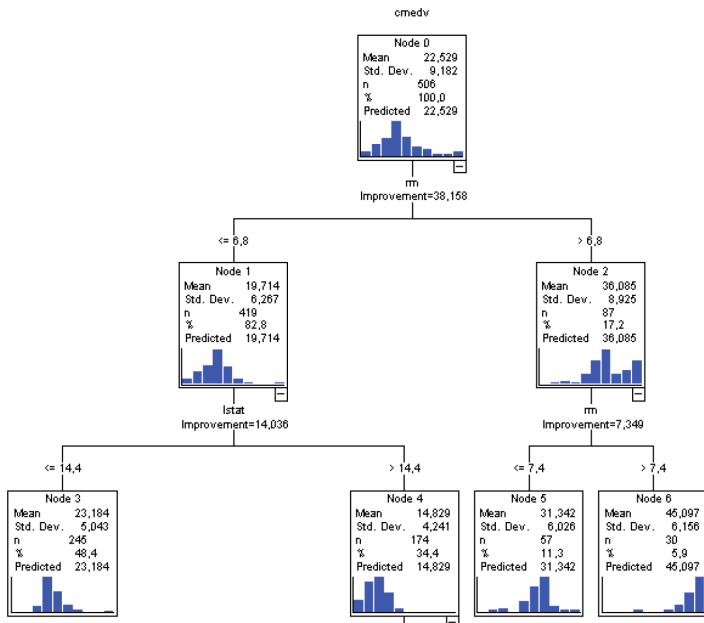
$$\sum_{i: \mathbf{x}_i \in R_1(j, s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: \mathbf{x}_i \in R_2(j, s)} (y_i - \hat{y}_{R_2})^2,$$

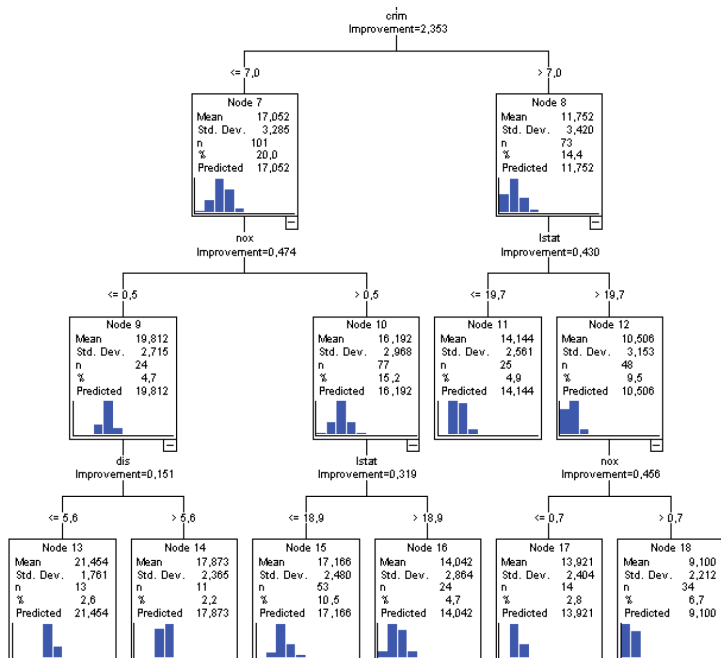
where  $\hat{y}_{R_1}$  and  $\hat{y}_{R_2}$  are averages in  $R_1$  and  $R_2$ .

line variances  
within rectangles.  
(without  $\frac{1}{n-1}$ ).

## Step 2

Repeat Step 1 to split regions  $R_1$  and  $R_2$  recursively.

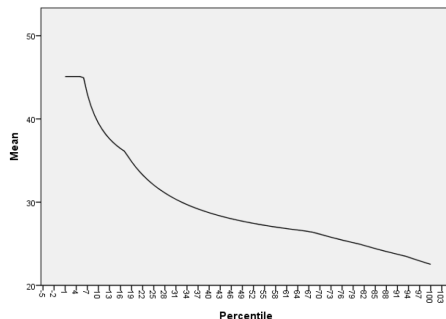






Gain Summary for Nodes

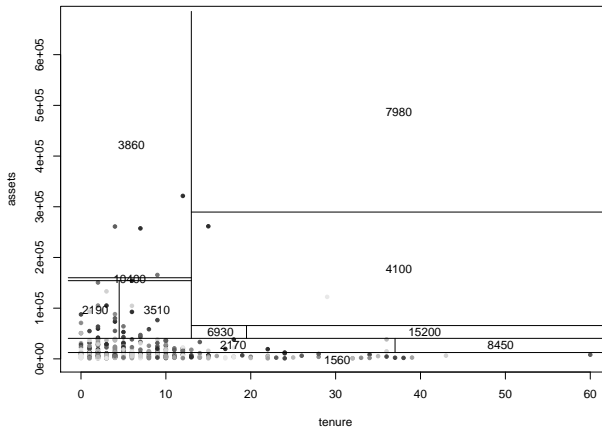
Node	Node-by-Node			Cumulative		
	N	Percent	Mean	N	Percent	Mean
6	30	5,9%	45,10	30	5,9%	45,10
5	57	11,3%	31,34	87	17,2%	36,09
3	245	48,4%	23,18	332	65,6%	26,56
13	13	2,6%	21,45	345	68,2%	26,37
14	11	2,2%	17,87	356	70,4%	26,11
15	53	10,5%	17,17	409	80,8%	24,95
11	25	4,9%	14,14	434	85,8%	24,33
16	24	4,7%	14,04	458	90,5%	23,79
17	14	2,8%	13,92	472	93,3%	23,50
18	34	6,7%	9,10	506	100,0%	22,53



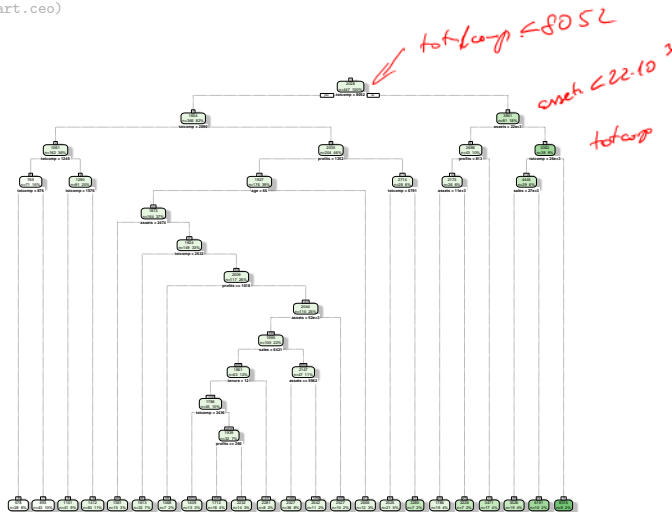
```

> library("tree")
> tree.ceo = tree(salary ~ tenure + assets, data=ceo)
> plot(ceo$tenure,ceo$assets, type="p", pch=20, xlab="tenure", ylab="assets")
> partition.tree(tree.ceo, ordvars=c("tenure","assets"), add=TRUE)

```



```
> library("rpart")
> rpart.ceo = rpart(salary ~ ., data=ceo, control=rpart.control(cp = 0.001))
> fancyRpartPlot(rpart.ceo)
```



Rattle 2017–Nov–16 09:30:27 okhrinya

**Note:** Using CART we can grow the tree to saturation.

- Fix the maximal number of splittings and a lower bound for the number of observations per region.
- Fix the minimal change in the objective function.
- **tree pruning**: after the optimal tree is found, it is shortened

$$R_{\alpha}(T) = \frac{1}{\sum_i (y_i - \bar{y})^2} \sum_{m=1}^{|T|} \sum_{i: \mathbf{x}_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

where  $|T|$  is the number of terminal nodes in a tree and  $\alpha$  is the **complexity parameter**.

## Key properties of CARTs

- For given  $\alpha$  it is possible to determine the tree  $T(\alpha)$  with the smallest  $R_{\alpha}(T)$  uniquely
- If  $\alpha > \beta$  then  $T(\alpha) = T(\beta)$  or  $T(\alpha)$  is a strict subtree of  $T(\beta)$ .

```
> printcp(rpart.ceo)
> printcp(rpart.ceo)
```

Regression tree:

```
rpart(formula = salary ~ ., data = ceo, control = rpart.control(cp = 0.001, xval = 10))
```

Variables actually used in tree construction:

```
[1] age      assets  profits sales  tenure totcomp
```

```
Root node error: 1323386794/447 = 2960597
```

n= 447

	CP	nsplit	rel error	xerror	xstd
1	0.2738266	0	1.00000	1.00703	0.19979
2	0.1091070	1	0.72617	0.77938	0.13638
3	0.0777412	2	0.61707	0.83535	0.15451
4	0.0646524	3	0.53933	0.77968	0.15159
5	0.0351651	4	0.47467	0.70797	0.15182
6	0.0130789	5	0.43951	0.73676	0.17353
7	0.0113130	6	0.42643	0.78445	0.17517
8	0.0081763	7	0.41512	0.79755	0.17644
9	0.0080167	8	0.40694	0.79045	0.17654
10	0.0052976	9	0.39892	0.78938	0.17620
11	0.0032733	10	0.39363	0.78670	0.17590
12	0.0029769	11	0.39035	0.77301	0.17377
13	0.0022593	12	0.38737	0.77146	0.17367
14	0.0017931	13	0.38512	0.77539	0.17310
15	0.0016171	15	0.38153	0.77520	0.17313
16	0.0016099	17	0.37829	0.77603	0.17313
17	0.0012678	20	0.37346	0.77635	0.17313
18	0.0012449	21	0.37220	0.77348	0.17303
19	0.0010000	22	0.37095	0.77403	0.17301

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

1 - R<sup>2</sup> → Cross-validation

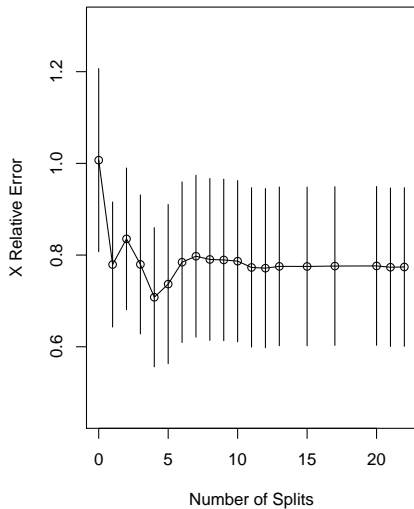
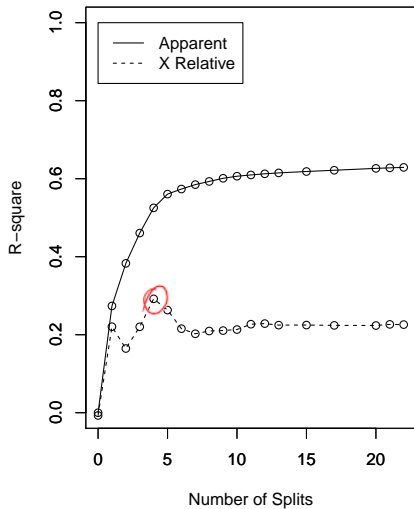
R<sup>2</sup> ≈ 0.63

Q: How to choose the overall optimal  $\alpha$  or subtree?  $\rightsquigarrow$  cross-validation

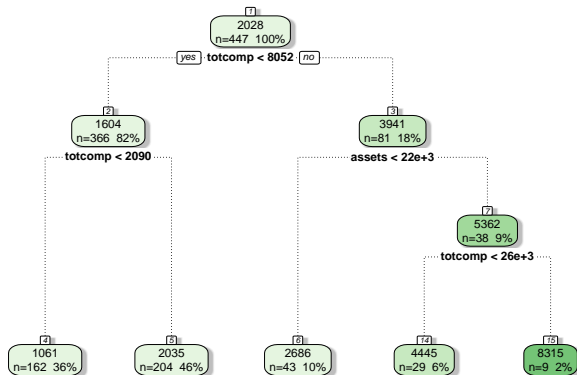
- The sequence of trees  $T_0$  (no splits) to  $T_m$  ( $m$  splits) uniquely determines the sequence of possible  $\alpha$ 's

$$\infty, \alpha_1, \dots, \alpha_{m-1}, \alpha_{min}$$

- Any  $\alpha$  between  $(\alpha_i, \alpha_{i+1}]$  leads to the same optimal subtree
- Define  $\beta_i = \sqrt{\alpha_i \alpha_{i+1}}$  as an “average” CP for every interval
- Split the data into  $B$  subsets  $G_1, \dots, G_B$  (10 by default)
  - For every subset excluding the  $G_i$ 's determine  $T_{\beta_1}, \dots, T_{\beta_m}$
  - Compute the relative MSE as the forecast loss for elements in  $G_i$
- Compute the average loss over all  $G_i$ 's and choose  $\beta$  (and thus the optimal subtree) which corresponds to the smallest one.



```
> cp.min = which.min(rpart.ceo$cptable[,4]);
> rpart.ceo.prune=prune(rpart.ceo, cp=rpart.ceo$cptable[cp.min,1])
> rpart.ceo.prune$variable.importance/sum(rpart.ceo.prune$variable.importance)
  totcomp    assets    sales    profits    tenure    age
0.52984007 0.18398873 0.10229360 0.09962347 0.05557637 0.02867777
```



Rattle 2017–Nov–16 14:56:49 okhrinya



## Generalizations

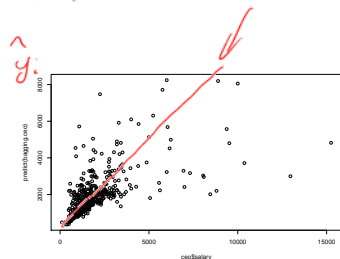
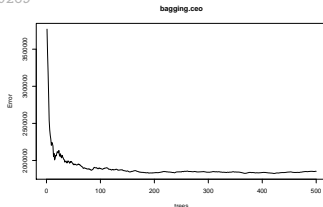
- **Bagging**: if you use for CART just a subsample, then you obtain a completely different tree.
  - Fit a CART to  $B$  random subsamples (*bootstrap*).
  - The error is measured on the remaining observations **out-of-bag**.
  - The final forecast is:

$$\hat{f}_{avr}(\mathbf{x}_0) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(\mathbf{x}_0).$$

← Love  $X\beta$

Mittelwert

```
> bagging.ceo= randomForest(salary ~ ., data=ceo, mtry=6)
> predict(bagging.ceo);
> cor(ceo$salary, predict(bagging.ceo))
[1] 0.6180269
```



## Random Forests: is a generalization of Bagging

- For each splitting you consider not all explanatory variables but just a subset of size  $M \approx \sqrt{J}$    
 *draw not only a subsample of observations but also a subset of X-variables.*
- ... this makes the trees more heterogenous and “uncorrelated” in terms of forecasts
- Each tree is grown on a bootstrap sample (as for bagging)
- The **importance** of a variable is measured by increase in (a) MSE ; (b) in node impurity over the out-of-bag sample if the variable is permuted

$$\Delta MSE_{j,b} =$$

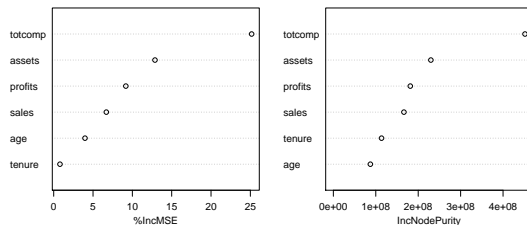
$$\frac{1}{|\bar{B}_b|} \sum_{k \in \bar{B}_b} \hat{u}^2(x_{1k}, \dots, x_{jk}) - \frac{1}{|\bar{B}_b|} \sum_{k \in \bar{B}_b} \tilde{u}_k^2(x_{1k}, \dots, x_{j-1,k}, \tilde{x}_{jk}, x_{j+1,k}, \dots, x_{Jk}),$$



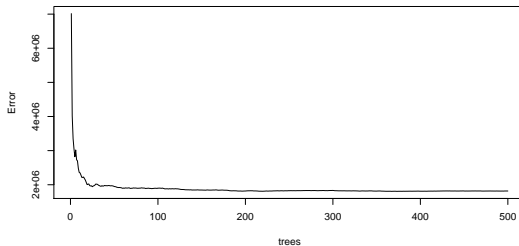
where  $\tilde{x}_j$  are the randomly permuted (reordered) observations on the  $j$ th variable and  $\bar{B}_b$  is the  $b$ th out-of-bag subsample.

```
> forest.ceo= randomForest(salary ~ ., data=ceo, importance=T)
> varImpPlot(forest.ceo)
```

forest.ceo



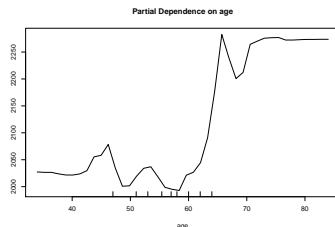
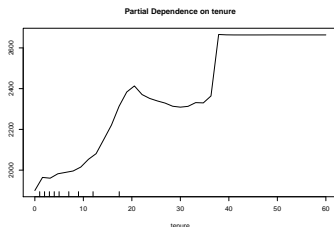
forest.ceo



**Partial dependence plots:** visualize the marginal impact of a variable/feature

$$\tilde{f}_j(x) = \frac{1}{K} \sum_{k=1}^K \hat{f}(x_{1k}, \dots, x_{j-1,k}, x, x_{j+1,k}, \dots, x_{Jk})$$

```
> partialPlot(forest.ceo, pred.data=ceo, x.var=tenure)
```



Ctree; CS.O

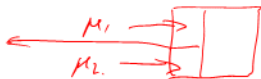
## CHAID

- An alternative approach is **CHAID (Chi-square Automatic Interaction Detectors)**: allows not only for binary splitting and is similar to ANOVA.
- Analysis is a generalization of two-sample test for the mean.
- **Idea**: let  $G$  be the number of splittings for variable  $X$ . We test if there is a significant difference between the means of  $Y$  in different regions.

CART - no tests  $\Rightarrow$  just optimization

CHAID - test.

split only  
if  $H_0: \mu_1 = \mu_2$   
is rejected.



$\Rightarrow$  we test the difference in the means for two or more rectangles



triple-splitting  
if  $H_0: \mu_1 = \mu_2 = \mu_3$   
is rejected.



## Total sum of squares

$$TSS = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{g=1}^G \sum_{i:\mathbf{x}_i \in R_g} (y_i - \bar{y})^2$$

## Within sum of squares

$$WSS = \sum_{g=1}^G \sum_{i:\mathbf{x}_i \in R_g} (y_i - \bar{y}_{R_g})^2$$

## Between sum of squares

$$BSS = TSS - WSS = \sum_{g=1}^G |R_g| (\bar{y}_{R_g} - \bar{y})^2.$$

$$H_0 : \mu_1 = \dots = \mu_G \quad \text{vs} \quad H_1 : \mu_i \neq \mu_j \text{ for at least one pair } i, j$$

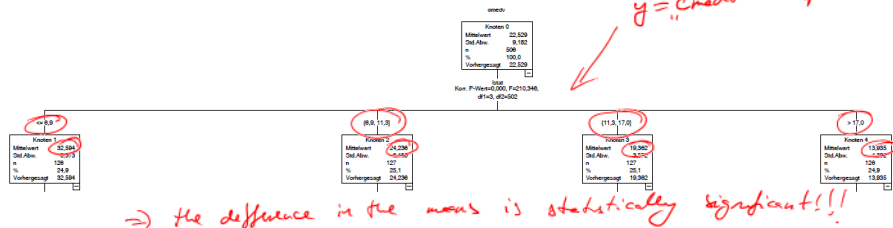
Test statistic:  $F = \frac{BSS/(G-1)}{WSS/(n-G)} \sim F_{G-1, n-G}$

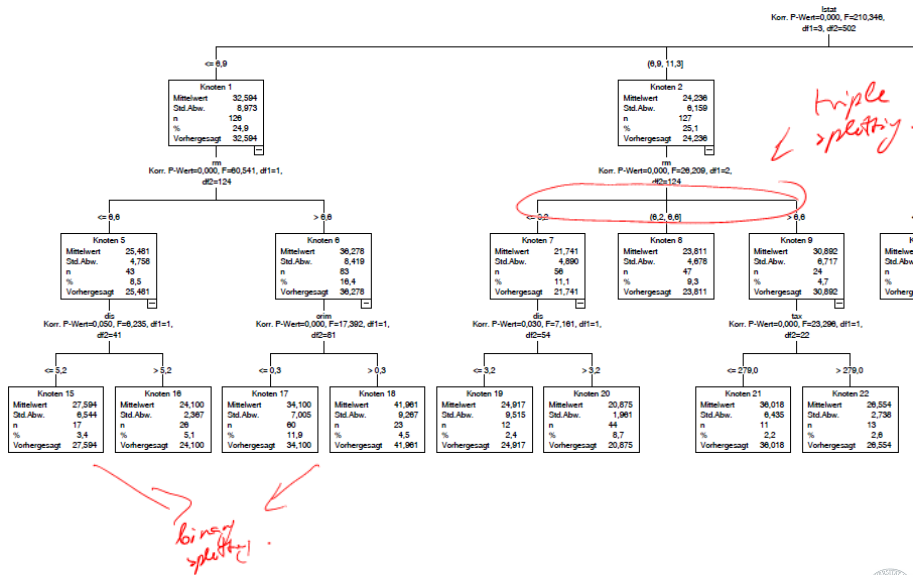
splitting  
in G rectangles  
simultaneously

F-distribution

## Idea:

- For each predictor we determine the optimal splitting, i.e. the regions with the smallest  $p$ -value of the test.
- The  $p$ -values should be corrected due to multiple testing (Bonferroni correction).
- The predictor with the smallest corrected  $p$ -value is used for splitting.







Baumtabelle												
Knoten	Mittelwert	Standardabweichung	N	Prozent	Vorhergesagter Mittelwert	Übergeordneter Knoten	Primäre unabhängige Variable					
							Variable	Sig. <sup>a</sup>	F	df1	df2	Werte aufteilen
0	22,53	9,182	506	100,0%	22,53							
1	32,59	8,973	126	24,9%	32,59	0	lstat	,000	210,346	3	502	<= 6,9
2	24,24	6,159	127	25,1%	24,24	0	lstat	,000	210,346	3	502	{6,9, 11,3}
3	19,36	3,572	127	25,1%	19,36	0	lstat	,000	210,346	3	502	{11,3, 17,0}
4	13,93	4,392	126	24,9%	13,93	0	lstat	,000	210,346	3	502	> 17,0
5	25,48	4,758	43	8,5%	25,48	1	rm	,000	60,541	1	124	<= 6,6
6	36,28	8,419	83	16,4%	36,28	1	rm	,000	60,541	1	124	> 6,6
7	21,74	4,890	56	11,1%	21,74	2	rm	,000	26,209	2	124	<= 6,2
8	23,81	4,678	47	9,3%	23,81	2	rm	,000	26,209	2	124	{6,2, 6,6}
9	30,89	6,717	24	4,7%	30,89	2	rm	,000	26,209	2	124	> 6,6
10	22,59	3,393	17	3,4%	22,59	3	tax	,000	18,192	1	125	<= 279,0
11	18,86	3,345	110	21,7%	18,86	3	tax	,000	18,192	1	125	> 279,0
12	19,31	2,881	14	2,8%	19,31	4	nox	,000	36,153	2	123	<= ,5
13	15,72	3,623	44	8,7%	15,72	4	nox	,000	36,153	2	123	{,5, ,6}
14	11,67	3,554	68	13,4%	11,67	4	nox	,000	36,153	2	123	> ,6
15	27,59	6,544	17	3,4%	27,59	5	dis	,050	6,235	1	41	<= 5,2
16	24,10	2,367	26	5,1%	24,10	5	dis	,050	6,235	1	41	> 5,2
17	34,10	7,005	60	11,9%	34,10	6	crim	,000	17,392	1	81	<= ,3
18	41,96	9,267	23	4,5%	41,96	6	crim	,000	17,392	1	81	> ,3
19	24,92	9,515	12	2,4%	24,92	7	dis	,030	7,161	1	54	<= 3,2
20	20,88	1,961	44	8,7%	20,88	7	dis	,030	7,161	1	54	> 3,2
21	36,02	6,435	11	2,2%	36,02	9	tax	,000	23,296	1	22	<= 279,0
22	26,55	2,738	13	2,6%	26,55	9	tax	,000	23,296	1	22	> 279,0
23	19,51	2,847	75	14,8%	19,51	11	nox	,007	9,596	1	108	<= ,6
24	17,47	3,911	35	6,9%	17,47	11	nox	,007	9,596	1	108	> ,6
25	17,24	3,698	24	4,7%	17,24	13	rad	,022	11,566	1	42	{2,0; 5,0; 6,0; 24
26	13,90	2,596	20	4,0%	13,90	13	rad	,022	11,566	1	42	4,0
27	10,85	3,087	56	11,1%	10,85	14	dis	,000	22,769	1	66	<= 2,1
28	15,53	3,094	12	2,4%	15,53	14	dis	,000	22,769	1	66	> 2,1

Linear model: linear in  $\beta$ 's

$$y_k = \beta_0 + \beta_1 x_{k1} + \dots + \beta_J x_{kJ} + u_k$$

## Nonlinear regression

The general form of a nonlinear regression is:

*h - multivariate*

$$y_k = h(\mathbf{x}_k, \boldsymbol{\beta}) + u_k,$$

*→ nonlinear function of  $\beta$ 's*

where  $h(\cdot, \cdot)$  is some unknown function of the regressors and parameters.

$$\log y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

→ •  $y = e^{\beta_0} e^{\beta_1 x_1} e^{\beta_2 x_2} e^u$  - can be linearized

→ •  $y = \beta_0 + \beta_1 e^{\beta_2 x_1} + u$  - cannot be linearized

→ •  $y = \beta_0 + \beta_1 x_1^\gamma + u$  - cannot be linearized

*→ OLS cannot be applied directly.*

A popular special case of the non-linear regression is the single-index model

$$y_k = h(\mathbf{x}'_k \boldsymbol{\beta}) + u_k$$

*link function.* *← =  $h(\beta_0 + \beta_1 x_{k1} + \dots + \beta_J x_{kJ}) + u_k$*  *here h - univariate function.*

thus  $h$  is a function of a linear combination of the regressors.

## Assumptions

- as before +
- $E(\mathbf{u}|\mathbf{X}) = \mathbf{0}$  is replaced with  $E(u_i|h(\mathbf{x}_i, \boldsymbol{\beta})) = 0$ : if  $u$  is uncorrelated with  $\mathbf{x}$  it still may be correlated with some function of  $\mathbf{x}$ . In general  $E(\mathbf{u}|\mathbf{X}) = \mathbf{0}$  is not needed.
- Identifiability of the model parameters: the model is identifiable if there is no a non-zero parameter  $\boldsymbol{\beta}_0$ , such that  $h(\mathbf{x}_i, \boldsymbol{\beta}_0) = h(\mathbf{x}_i, \boldsymbol{\beta})$  for all  $\mathbf{x}_i$ .

**Note:** in the linear regression it is sufficient to assume  $\text{rank}(\mathbf{X}'\mathbf{X}) = J + 1$ . Here it is not enough.

$$y = \frac{2\beta_0 + 2\beta_1 x_1}{2\beta_2 + 2\beta_3 x_2} + u.$$

← the parameters cannot be identified uniquely!

⇒ impose some restriction:  $\beta_0 = 1$       $y = \frac{1 + \beta_1 x_1}{\beta_2 + \beta_3 x_2} + u.$

$$y = \beta_0 + \beta_1 \cdot e^{\beta_2 + \beta_3 x_2} + u = \beta_0 + \underbrace{\beta_1 \cdot e^{\beta_2}}_{\beta_4} \cdot e^{\beta_3 x_2} + u.$$



**Estimation:** the LS estimation can be used, but the asymptotic theory follows in a straightforward way from the quasi (!) maximum-likelihood estimation.

Assuming Gaussian residuals it holds:

$$\begin{aligned}
 \mathcal{L}(y_k | \mathbf{x}_k, \boldsymbol{\beta}, u_k) &= \frac{1}{K} \sum_{k=1}^K \ln f(y_k | \mathbf{x}_k, \boldsymbol{\beta}, u_k) \\
 &= \frac{1}{K} \sum_{k=1}^K \ln \left\{ \frac{1}{\sqrt{2\sigma^2}} \exp \left( -\frac{1}{2\sigma^2} (y_k - h(\mathbf{x}_k, \boldsymbol{\beta}))^2 \right) \right\} \\
 &= \frac{1}{K} \sum_{k=1}^K \left[ -\ln \sqrt{2\sigma^2} - \frac{1}{2\sigma^2} (y_k - h(\mathbf{x}_k, \boldsymbol{\beta}))^2 \right] \rightarrow \max_{\text{w.r.t. } \boldsymbol{\beta}}
 \end{aligned}$$

*density function of  $y^i$*   
*Gaussian density*  
*our nonlinear component*

Thus the first order conditions for  $\boldsymbol{\beta}$  are

$$\frac{\partial \mathcal{L}(y_k | \mathbf{x}_k, \boldsymbol{\beta}, u_k)}{\partial \boldsymbol{\beta}} = \sum_{k=1}^K (y_k - h(\mathbf{x}_k, \boldsymbol{\beta})) \frac{\partial h(\mathbf{x}_k, \boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \mathbf{0}.$$

*system of nonlinear equations*  
*⇒ no explicit solution.*

⇒ mostly a highly nonlinear system of equations solved numerically.

*⇒ in contrast to  $\boldsymbol{\beta} = (X'X)^{-1} X'Y$  in the linear model.*

**Consequences:** since the resulting  $\hat{\beta}$  is a non-linear function of the residuals  $u_k$

- ... the unbiasedness can not be proven in simple fashion;
  - ... the variance of  $\hat{\beta}$  is not easy to derive;
  - ... the exact distribution of  $\hat{\beta}$  is not Gaussian;
  - ... all the inferences, like tests, are valid only asymptotically.
- *t-test; F-test*

**but** the ML estimators are consistent and efficient (they possess the smallest variance among all consistent and asymptotically normal estimators)

Where all this comes from?

- Taylor expansion of  $f(x)$  in neighborhood of  $x_0$

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 + \dots$$

- Exact Taylor expansion of  $f(x)$  in neighborhood of  $x_0$  (mean-value theorem)

$$f(x) = f(x_0) + f'(x_+)(x - x_0),$$

where  $x_+$  lies between  $x$  and  $x_0$ .

- 

$$\left. \frac{\partial \mathcal{L}(y_k | \mathbf{x}_k, \boldsymbol{\beta}_k, u_k)}{\partial \boldsymbol{\beta}} \right|_{\hat{\boldsymbol{\beta}}} = \left. \frac{\partial \mathcal{L}(y_k | \mathbf{x}_k, \boldsymbol{\beta}_k, u_k)}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}} + \left. \frac{\partial^2 \mathcal{L}(y_k | \mathbf{x}_k, \boldsymbol{\beta}_k, u_k)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \right|_{\boldsymbol{\beta}_+} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})$$

- 

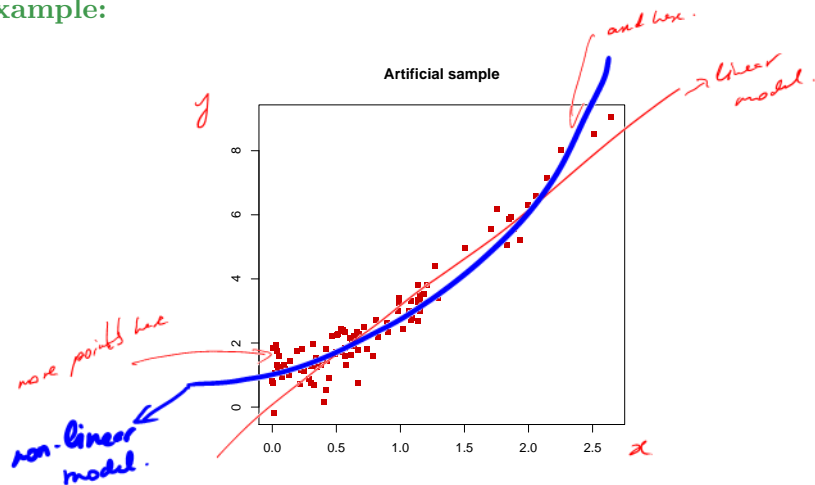
$$\sqrt{K}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = - \left( \overbrace{\left. \frac{\partial^2 \mathcal{L}(y_k | \mathbf{x}_k, \boldsymbol{\beta}_k, u_k)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \right|_{\boldsymbol{\beta}_+}}^{\mathbf{A}} \right)^{-1} \sqrt{K} \overbrace{\left. \frac{\partial \mathcal{L}(y_k | \mathbf{x}_k, \boldsymbol{\beta}_k, u_k)}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}}}^{\mathbf{B}}$$

- numerical estimates of the parameters

$$\sqrt{K}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \overset{\text{approx}}{\sim} N(\mathbf{0}, \mathbf{A}^{-1} \mathbf{B} \mathbf{A}^{-1})$$

*from the central limit theorem*  
*messy due to nonlinearity*

## Example:



- Model 1 :  $y = \beta_0 + \beta_1 x + u$

- Model 2 :  $y = \beta_0 + \beta_1 x^{\beta_2} + u$

cannot be linearized  
power regression.

```

> z1 = lm(y ~ x)
> z2 = nls(y ~ b0 + b1 * x^b2, start=list(b0=0, b1=1, b2=2))
> summary(z1)

```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.44547	0.09694	4.595	1.29e-05 ***
x	<u>2.78805</u>	0.09787	28.488	<u>&lt; 2e-16 ***</u>

Residual standard error: 0.604 on 98 degrees of freedom

Multiple R-squared: 0.8923, Adjusted R-squared: 0.8912

F-statistic: 811.5 on 1 and 98 DF, p-value: &lt; 2.2e-16

&gt; summary(z2)

Formula: y ~ b0 + b1 \* x^b2

Parameters:

	Estimate	Std. Error	t value	Pr(> t )
b0	1.08702	0.09345	11.63	<2e-16 ***
b1	1.77926	0.13041	13.64	<2e-16 ***
b2	1.56810	0.08382	18.71	<2e-16 ***

Residual standard error: 0.4678 on 97 degrees of freedom

Number of iterations to convergence: 5

Achieved convergence tolerance: 2.905e-06

True model:  $y = 1 + 2x^{1.5} + u, u \sim N(0, 0.5^2)$ .

I used to simulate the data

close to the true values!

 $\beta_0 \quad \beta_1 \quad \beta_2$ 

from linear model  
 $\Rightarrow$  from here linear model is a good model, with significant  $\beta_1$  and high  $R^2$ !  
 $\Rightarrow$  look at residuals!!!

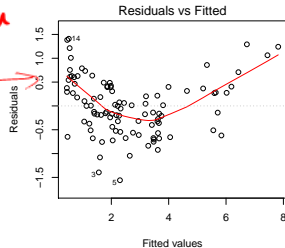
Note: that  $\beta_1 = 2.78805$  from the linear model has no relationship to  $\beta_1$  and  $\beta_2$  in the non-linear model.



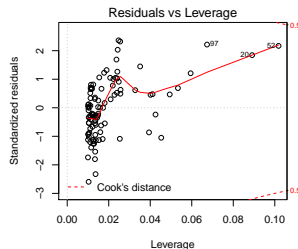
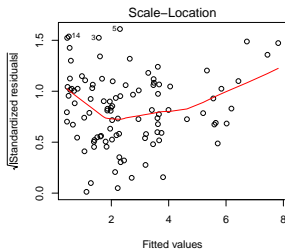
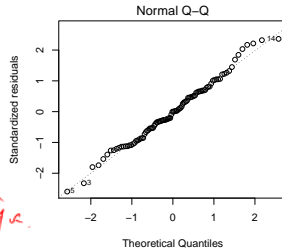
$$\text{lm}(y \sim x)$$

but the residuals are not really random, but have a pattern!

u



y



LR: big data: first estimate the model  $\hat{\beta}_0, \dots, \hat{\beta}_J \Rightarrow$  only after this the model selection  
 $K$  - small;  $J$  - big  $\Rightarrow$  low precision  $\Rightarrow$  many  $\beta$ 's insignificant  $\Rightarrow$  unreliable selection.

## Lasso regression

Next: model selection and estimation simultaneously

### Problems:

- **accuracy:** In  $K$  is much larger than  $J$ , then the variances are small and the inferences are precise. Low number of observations per parameter implies general high variability.
- **interpretability:** In large data sets there always irrelevant variables which make the economic interpretability difficult.
- **sparsity:** only a subset of the explanatory variables is relevant economically and statistically.

**Solution:** stepwise variable selection procedures based on statistical properties of the estimators or **lasso regression**

**Idea:** minimize the sum of squared residuals with constraints on the parameters.  $\Rightarrow$  model selection: which  $\beta$  are set to zero and which not!

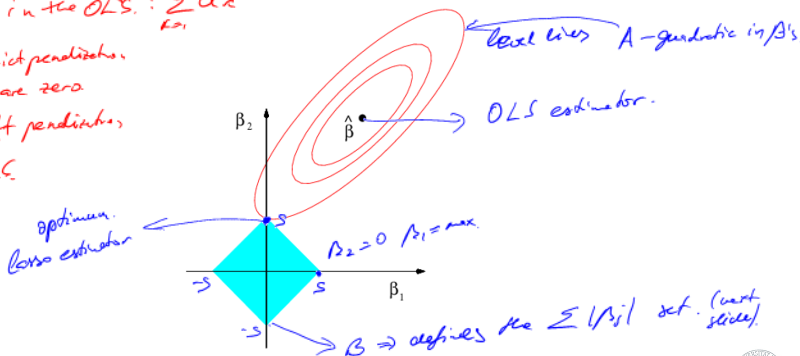
The objective function of the OLS procedure is replaced with

$$\sum_{k=1}^K (y_k - \beta_0 - \sum_{j=1}^J \beta_j x_{kj})^2 + \lambda \sum_{j=1}^J |\beta_j| \rightarrow \min, \text{ w.r.t } \beta_j\text{'s}$$

*Handwritten notes:*

- $\lambda$  (blue arrow)  $\rightarrow$  control parameter
- $\beta_j$  (blue arrow)  $\rightarrow$  to set risk of  $> 0 < 0$
- $\sum_{k=1}^K (y_k - \beta_0 - \sum_{j=1}^J \beta_j x_{kj})^2$  (red bracket)  $\rightarrow$  as in the OLS:  $\sum_{k=1}^K u_k^2$
- $\sum_{j=1}^J |\beta_j|$  (red bracket)  $\rightarrow$  forces some of  $\beta$  to become zero.

$\lambda$ -large  $\Rightarrow$  strict penalization,  
 $\Rightarrow$  more  $\beta$ 's are zero  
 $\lambda$ -small  $\Rightarrow$  soft penalization,  
 $\lambda = 0 \Rightarrow$  OLS



Note:  $f(x) \rightarrow \min, g(x) \leq a. \Rightarrow f(x) + \lambda(g(x) - a) \rightarrow \min, \lambda$ -Lagrange.

- The problem is equivalent to the following problem, i.e. for each  $\lambda$  there exists  $s$  such that both problems lead to the same lasso-coefficients.

$$OLS \Rightarrow \sum_{k=1}^K (y_k - \beta_0 - \sum_{j=1}^J \beta_j x_{kj})^2 \rightarrow \min, \text{ w.r.t. } \beta_j\text{'s}$$

$$\text{s.t. } \sum_{j=1}^J |\beta_j| \leq s.$$

constrained OLS  
 $s \leftrightarrow \lambda$

- Minimizing the objective is not trivial and there many specific numerical methods developed for this purpose.
- Selecting a good value for  $\lambda$  is crucial. The optimal value is chosen by cross-validation.

## Special case

Assume an individual constant for each observation:

$$\sum_{k=1}^K (y_k - \beta_k)^2$$

*replace our xP by a individual constant, for every yx.*

with the OLS solution  $\hat{\beta}_k = y_k$ .

With lasso we obtain:

$$\sum_{k=1}^K (y_k - \beta_k)^2 + \lambda \sum_{k=1}^K |\beta_k| \rightarrow \min$$

*loss part*

$$2(y_k - \beta_k) \pm \lambda = 0$$

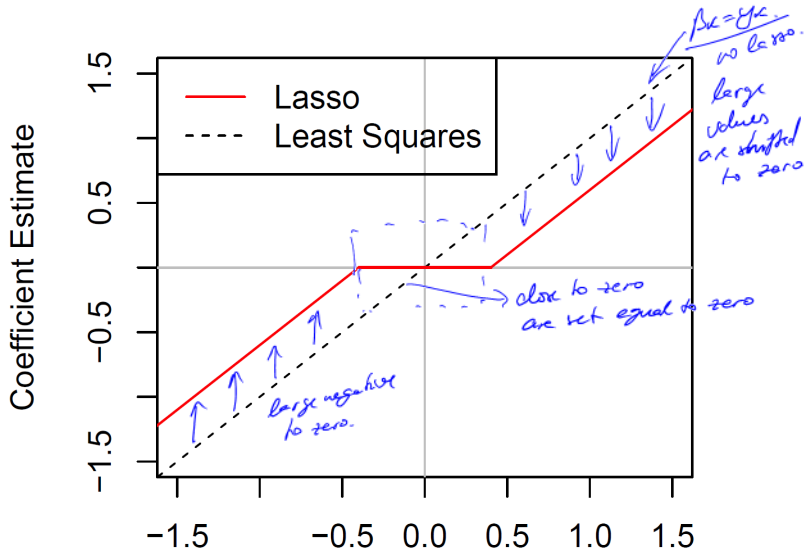
$$\beta_k = y_k \mp \lambda/2$$

with the solution

*lasso estimator of the constant.*

$$\hat{\beta}_k^{(lasso)} = \begin{cases} y_k - \lambda/2, & \text{if } y_k \geq \lambda/2 \\ y_k + \lambda/2, & \text{if } y_k \leq -\lambda/2 \\ 0, & \text{if } |y_k| \leq \lambda/2 \end{cases}$$

*→ yk is large positive*  
*→ yk is large negative*  
*→ yk is close to zero*



$\Rightarrow$  consequence of the penalty  $\lambda \sum |\beta_j|$

## Example:

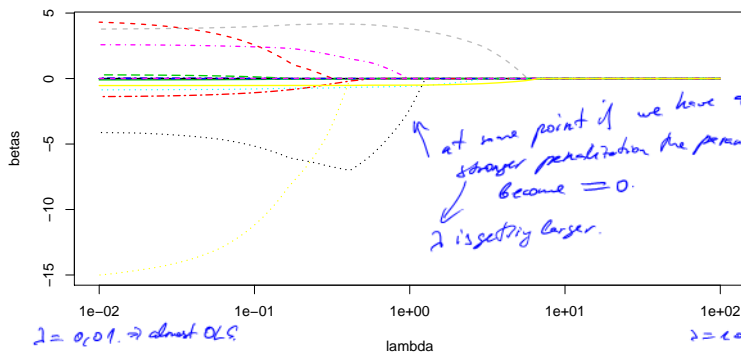
Trevor Hastie (Statistical learning book!)

lasso.

```

> grid = 10^seq(2,-2, length=100)
> lasso = glmnet(X, y.boston, alpha=1, lambda=grid);
> plot(lasso$beta)

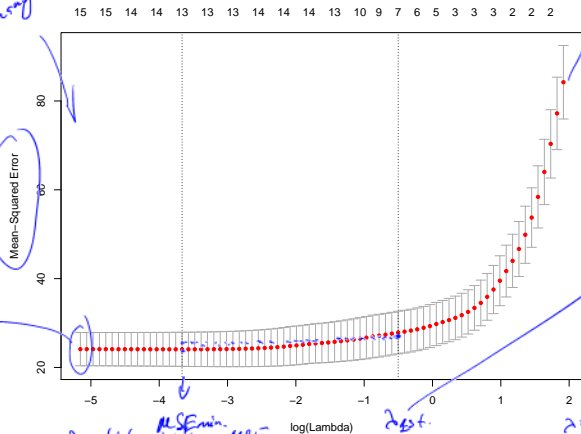
```



```
> cv.lasso = cv.glmnet(X, y.boston, alpha=1);
> plot(cv.lasso)
> cv.lasso$lambda.min
[1] 0.0255856
```

many nonzero betas

few nonzero betas.



large MSE,  
because many  
 $\beta = 0 \Rightarrow$   
few regressors  
in the model  
 $\Rightarrow$  high MSE  
lambda which  
is by 1.5  
larger than  
 $\lambda_{min.}$   
 $\Rightarrow$  to avoid  
overfitting.

we apply CV  
for every  $\lambda$   
 $\Rightarrow$  interval is computed  
using MSE's  
for CV's

$\lambda$  which minimises MSE.

$\lambda_{+}$

$\lambda = 100$



```
> lasso.coef = predict(lasso, type="coefficients", s=cv.lasso$lambda.min);
```

```
> lasso.coef
```

```
16 x 1 sparse Matrix of class "dgCMatrix"
```

```
1
```

```
(Intercept) -4.392079e+02
```

```
lon -4.250433e+00
```

```
lat 4.017435e+00
```

```
crim -9.553123e-02
```

```
zn 4.149031e-02
```

```
indus .
```

$\Rightarrow \text{set} = 0$

```
chas1 2.557345e+00
```

```
nox -1.432740e+01
```

```
rm 3.816713e+00
```

```
age .
```

$\Rightarrow \text{set} = 0$

```
dis -1.329954e+00
```

```
rad 2.572257e-01
```

```
tax -1.071061e-02
```

```
ptratio -8.520927e-01
```

```
b 8.974756e-03
```

```
lstat -5.326668e-01
```

$\lambda$  which minimizes MSE

others are smaller, than OLS  
(in most of the cases)