

## Chapter 5 - Resampling methods

### I Cross-validation

Split into training-validation-test set

leave one out: (LOOCV) → fit model  $n$  times using  $n-1$  datapoints

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}; \quad w/ \text{MSE}_i = (y_i - \hat{y}_i)^2$$

and model was fit  
w/ the  $n-1$  samples leaving  
the  $i$ th one out

- no randomness, always gives the same result  
(unlike w/ validation set, which is different even if you pick a different validation set)

- uses the whole (almost whole) dataset to fit unless if we have a validation set

- could be expensive to fit

statistic for regression:

$$CV_{(n)} = \frac{1}{n} \sum \left( \frac{y_i - \hat{y}_i}{1 - q_i} \right)^2$$

### K-fold cross-validation

- $K$  groups of the dataset, train on  $K-1$ , validate on 1

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE};$$

- less expensive to fit than LOOCV

- same variability depending how we choose the  $k$  groups

- LOOCV uses more data to fit than  $k$ -fold CV,  
so is preferred from a bias - perspective.
  - But LOOCV has a higher variance than  $k$ -fold CV  
 it is the avg of  $n$  very correlated models, while  $k$ -fold  
 has the avg of  $k$  correlated models.  
 averages of very correlated things have a higher variance
- Best is to choose  $k=5$  or  $10$ . for  $k$ -fold CV
- For classification problems, just take the misclassifications instead of MSE

## II Bootstrap

Generate new datasets w/ repeated sampling from the original dataset  
 we estimate the value of a parameter for each dataset  
 we get the mean and the standard error of the parameter, for  
 the bootstrap datasets to have an idea of the true value.

## Exercises

1)  $\min_{\alpha} \text{Var}(\alpha X + (1-\alpha)Y)$

$$= \min_{\alpha} \alpha^2 \text{Var}(X) + (1-\alpha)^2 \text{Var}(Y) + 2\alpha(1-\alpha) \text{Cov}(X, Y)$$

$$\begin{aligned} &= \min_{\alpha} \alpha^2 (\sigma_{xx}^2 + \sigma_{yy}^2 - 2\sigma_{xy}) \\ &\quad + \alpha(2\sigma_{xy} - 2\sigma_{yy}^2) \\ &\quad + \frac{\sigma_{yy}^2}{\sigma_{yy}^2} \end{aligned}$$

derivative = 0:

$$0 = 2\alpha(\sigma_{xx}^2 + \sigma_{yy}^2 - 2\sigma_{xy}) + 2\sigma_{xy} - 2\sigma_{yy}^2$$

$$\boxed{\alpha = \frac{\sigma_{yy}^2 - \sigma_{xy}}{\sigma_{xx}^2 + \sigma_{yy}^2 - 2\sigma_{xy}}}$$

2) Bootstrap sample n observations

a) p(first bootstrap observation not the jth sample)  
 $= 1 - 1/n$

b) p(second bootstrap not the jth sample)  
 $= 1 - 1/n$

c) p(jth observation is not in the bootstrap)  
 $= (1 - 1/n)^n$  (not the jth sample for any

d) p(jth observation is in the sample)  
 $= 1 - (1 - 1/n)^n$   
 $\begin{aligned} n=5 &\approx 0.67 \\ n=100 &\approx 0.63 \\ n=10000 &\approx 0.63 \end{aligned}$

In general,  $(1 - 1/n)^n \rightarrow \frac{1}{e} \approx 0.367$

3) k-fold cross-validation

You divide your data into k subsets of equal size.

For each set, you train on the rest and evaluate the error

on that set

Then you average your error on all your sets to get an estimate.

(b) advantage relative to: validation set

- you get more data to train
- result is less random (depends less on the set)

disadv:

- validating:

- more expensive to fit

adv - relative to LOOCV:

- lower variance
- less expensive to fit

disadv:

- less data to train on
- not always the same result

4)

Estimate standard deviation w/ bootstrap:

- Resample from dataset many times to get artificial datasets
- Estimate prediction for each dataset
- Calculate the std dev. of the estimates

## G. Linear model Selection and Regularization

Least squares fit is not the best for linear regression

- if  $p \approx n$  or  $p > n$  (lots of features)
  - as can result in overfitting
  - or cannot even be used
- if some variables are irrelevant - least squares won't put coeff to 0

alternatives:

- subset selection
- regularization
- dimension reduction

### I. Subset selection

#### Best subset selection

1.  $H_0$  - no predictors, just single mean.
2. for  $S=1, \dots, p$ 
  - (a) fit all  $\binom{p}{S}$  models w/  $S$  predictors
  - (b) choose best model  $H_S$  for all  $S$  (subset LSS)  
or by  $R^2$
3. Select best model among  $H_0, \dots, H_p$  using  
cross-validated prediction error,  $C_p(AIC)$ , BIC or  
adjusted  $R^2$

#### Stepwise selection

##### Forward:

1.  $H_0$  - null model
2. For  $S=0, \dots, p-1$ 
  - (a) consider all  $p-S$  models that augment the predictors in  $H_S$  w/ one additional predictor

(e) choose the one w/ the highest  $R^2$  or lowest RSS

3. Select best model from  $M_0, \dots M_p$  using cross-validated prediction error,  $C_p(AIC)$ , BIC or adjusted  $R^2$

$\Rightarrow 1 + \frac{p(p+1)}{2}$  models, not  $2^p$

- faster
- not guaranteed to find optimal model
- can be used if  $n < p$ , but only  $M_0, \dots M_{n-1}$  can be considered

backward:

1.  $M_p$  - full model

2. For  $S = p, \dots, 1$

(a) consider all  $S$  models that contain the predictors in  $M_S$  except one predictor

(b) choose the one w/ the highest  $R^2$  or lowest RSS

3. Select best model from  $M_0 \sim M_p$  using cross-validated pred. err.  
 $C_p(AIC)$ , BIC or adjusted  $R^2$

$\Rightarrow 1 + \frac{p(p+1)}{2}$  models, not  $2^p$

- faster
- not guaranteed to find optimal
- cannot be used if  $n < p$

hybrid approaches:

1. adding and/or removing in each step

choosing the best model

RSS and  $R^2$  is the lowest for the model w/ most features

this is not a good measure to use

2 approaches:  
I adjusting by taking into account the number of features

II estimating the test error w/ validation

$$g) C_p = \frac{1}{n} (RSS + 2d\hat{\sigma}^2) \quad \text{for } d \text{ predictors}$$

penalty of  $2d\hat{\sigma}^2$  for having RSS  
 $C_p$  is unbiased estimate of test RSS  
if  $\hat{\sigma}^2$  is unbiased estimate for  $\sigma^2$  (not parentheses)

choose ~~best~~ model : the one w/ lowest  $C_p$

### b) AIC

For large models fit w/ maximum likelihood

$$\text{For } Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon \quad \varepsilon \sim N(0, \hat{\sigma}^2)$$

MLE is same as least squares so we can use it.

$$AIC = \frac{1}{n\hat{\sigma}^2} (RSS + 2d\hat{\sigma}^2)$$

So for least squares,  $C_p \approx AIC$

### c) BIC

$$BIC = \frac{1}{n\hat{\sigma}^2} (RSS + \log(n)d\hat{\sigma}^2)$$

$\log n > 2$  for  $n > 7$ ,  $\Rightarrow$  heavier penalty for large models

### d) adjusted $R^2$

$$\text{adjusted } R^2 = 1 - \frac{RSS/(n-d-1)}{TSS/(n-1)} \quad (\text{instead of } 1 - \frac{RSS}{TSS})$$

For this one, we want the model w/ largest adjusted  $R^2$

Since  $TSS = \sum (y_i - \bar{y})^2$  is fix, this is equivalent to

$$\text{minimizing } \frac{RSS}{n-d-1}$$

## II Validation and CV

Choose model w/ -fold validation / CV error

- +: direct estimate for test error  
less assumptions about model  
works for models w/o explicit degrees of freedom  
a when hard to estimate  $\sigma^2$
- expensive

### One-standard-error rule

Calculate standard error for estimated MSE for all model

Select smallest where estimated test error is within one SE  
of the lowest point on the curve

meaning if models are similar, select simplest model

### 6.2 Shrinkage methods - regularizing the coeffs

#### Ridge regression $\rightarrow$ L2 penalty

$$DSS = \sum (y_i - \beta_0 - \sum \beta_j x_{ij})^2$$

Instead of minimizing DSS, we want to minimize  $DSS + \lambda \sum_j \beta_j^2$

$\lambda \geq 0$  is a tuning parameter

[we do not shrink  $\beta_0$ , as this is estimate for the mean!]

$\lambda$  penalizing large  $\beta$  values: as long,  $\|\beta\|_2^2$  decreases always

Normal least squares is scale equivariant - multiplying  $x_i$  by c  
leads to the same model.

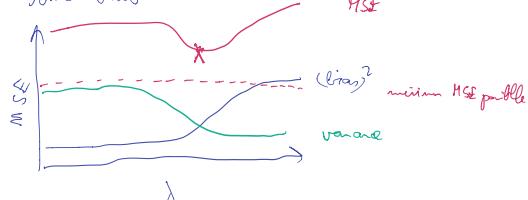
THIS IS NOT TRUE FOR RIDGE REGRESSION!!!

$$\text{STANDARDIZE: } \hat{x}_{ij} = \frac{x_{ij}}{\sqrt{\sum (x_{ij} - \bar{x})^2}} \quad (\text{s.t. dev of } 1)$$

$$\text{Since } E(y_i - \hat{f}(x_i))^2 = \text{Var}(\hat{f}(x_i)) + [\text{Bias}(\hat{f}(x_i))]^2 + \text{Var}(\epsilon)$$

This helps us reducing the variance while introducing

some bias



disadvantage: doesn't set my coeffs to 0  
makes it harder to interpret the model  
always uses all of the predictors

Lasso — tries to overcome the problem of ridge  
by forcing some params to 0

using  $L_1$  penalty

$$\text{minimize } RSS + \lambda \sum_{j=1}^p |\beta_j|$$

variable selection  
sparse models  
= not using all variables

Equivalent formulations:

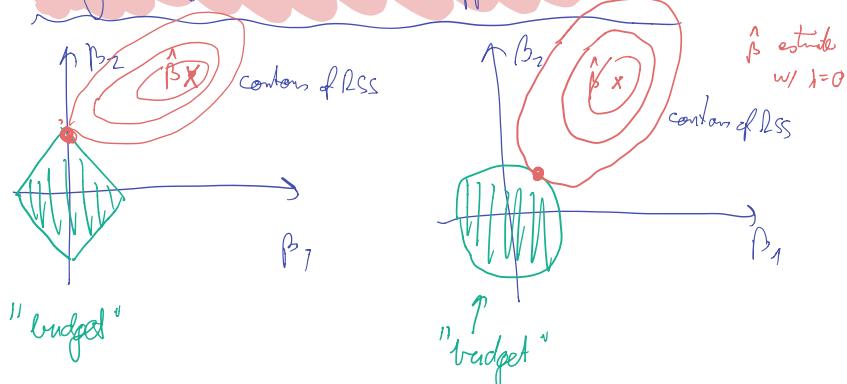
$$\min_{\beta} RSS \quad \text{given} \quad \sum_{j=1}^p |\beta_j| \leq s = \begin{cases} \text{Ridge} & m=1 \\ \text{Lasso} & m=2 \end{cases}$$

$m=1$  for lasso  
 $m=2$  for ridge

each  $\lambda$  is equivalent to some value  $s$   
(meaning they result in same coeffs)

"budget for  $\|\beta\|_m$ " — as long as least squares is  
within the budget, it will result in that

Why does Lasso set coefficients to 0?



because of the sharp edge, Lasso gets an intersection w/ lowest RSS at the corner more often

## Ridge vs Lasso

Ridge is better when all variables related

Lasso is better when only few are actually related

→ can't know in advance, try both

## Bayesian Interpretation

$\beta$  has a prior  $p(\beta)$

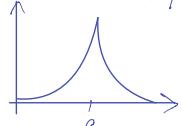
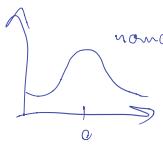
$$P(\beta|x,y) \propto f(y|x,\beta) p(\beta|x) = f(y|x,\beta) p(\beta)$$

Assume  $y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$

$$\text{and } p(\beta) = \prod_{j=1}^p p(\beta_j) \sim \text{prior}$$

• If  $\beta \sim \text{Normal}(0, R\mathbf{I})$

• If  $\beta \sim \text{dgamma}(0, b(x))$



→ more likely to set to 0 as distribution expected there

## Selecting the Tuning Parameter

Choose a grid for  $\lambda$  values ( $\alpha^{-1}$ )

using cross-validation, compute one w/ smallest error, then refit model using all values



## 6.3 - Dimension Reduction Methods

$Z_1, \dots, Z_M$  are linear combinations of  $x_1, \dots, x_p$  w/  $M < p$

$$\sum_m^M = \sum_{j=1}^p \phi_{jm} x_j \quad \text{for some } \phi_{1m}, \phi_{pm}, m=1, \dots, M$$

$$y_i = \theta_0 + \sum_{m=1}^M \theta_m z_{im} + \epsilon_i \quad i=1, \dots, n \quad \text{using least squares}$$

reduces  $p+1$  coeffs  $\beta_0, \dots, \beta_p$  to  $\theta_0, \dots, \theta_M$

having to fit  $M+1 < p+1$  coeffs:

$$\sum_m^M \theta_m z_{im} = \sum_{m=1}^M \theta_m \sum_{j=1}^p \phi_{jm} x_j = \sum_{j=1}^p \sum_{m=1}^M \theta_m \phi_{jm} x_j = \sum_{j=1}^p \beta_j x_j$$

$$\Rightarrow \beta_j = \sum_{m=1}^M \theta_m \phi_{jm} \quad (*)$$

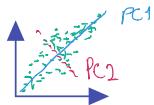
So this is a special case of linear regression where coeffs have to be  $\theta$ . This is a constraint, so increases bias and decreases variance if  $M < p$ . If  $M=p$  and all  $Z_m$  are truly indep, this is not a constraint.

How to get  $Z_i$ ? (or  $\phi_{jm}$ )

## Principal component regression (PCR)

principal components are chosen s.t.

- they are a linear combination of the data
- w/  $\sum_{j=1}^p \phi_j^2 = 1$
- s.t. they minimize the variance given  
that they are uncorrelated / perpendicular to the previous components



We assume that variables in which  $x_1, \dots, x_p$  vary the most are directly associated w/  $y$ .

Fitting least squares to  $z_1, \dots, z_m$  is better, because we have most of the signal, and we overfit less due to fewer params

→ NOT A FEATURE SELECTION METHOD,  
IT USES ALL

closely related to ridge regression

- STANDARDIZE BEFORE PCR  
(values in same units)
- Choose number of components w/ cross-validation

## Partial Least Squares (PLS)



- identify  $z_1, \dots, z_m$ , fit w least squares

that = choose  $z_1, \dots, z_m$  in a reversed way

using  $Y$ : those that are related to the response

directions explaining both response and predictors

1, Standardize predictors

2, Compute  $z_1$  by setting  $\beta_{j1}$  to the coeff from a single linear regression of  $Y$  onto  $X_j$ :

→ this places highest weight on variables that are most correlated to the response

3, We regress each variable on  $z_1$ , and take residuals  
(removing info that is not explained by  $z_1$ )

4, We get  $z_2, \dots, z_m$  as  $z_i$ , using the residuals

Complicated, mostly not better than PCR

## High-dimensional data

most methods have been developed for  $n \gg p$ ,

having large  $p$  is relatively new

→ high-dimensional data

least squares goes wrong if  $p \geq n$ , because we have many sets of coeffs that fit the data perfectly (like line through point)

→ overfit

- Cp, AIC and BIC don't help in High dimensional setting

(estimating  $\hat{\sigma}^2$  is problematic)

- adjusted  $R^2$  also not, as can easily get adjusted  $R^2 = 1$

adding noise features makes the model worse

→ (overfitting)

- if we identify a set of variables that predict the data well, it's not true it's the only set of variables that do so

→ (multicollinearity)

- if  $p > n$ , don't use ~~say~~ p-values,  $R^2$   
(they could be arbitrarily small due to overfitting)

→ use test set error / CV error

MSE /  $R^2$  on independent test set is valid

but not on training set

## Exercises

1) Best subset Forward stepwise Backward stepwise

we get  $p+1$  models for  $0, \dots, p$  predictors

for  $q$  predictors:

smallest training RSS:

X

smallest test RSS:

not true . least subset selection has the most flexibility,  
but other methods might be better ---

(i) True

(ii) True

(iii) False

(iv) False

(v) False

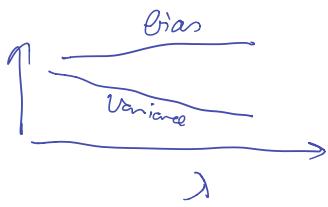
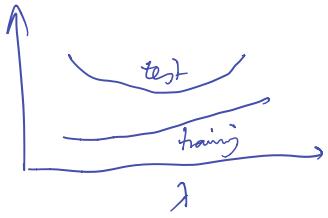
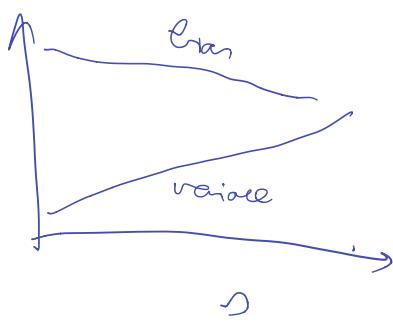
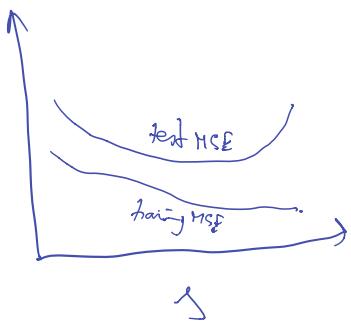
2) lasso relative to least-squares: Ridge relative to least-squares non-linear relative to a.s.  
(iii) (iii) (ii)

3) training RSS : iv 4, iii

1, 1, 0 or ...

v

test loss : 11  
 Variance : 33  
 Bias : 44  
 irreducible error : 55



5)  $n=2$        $\beta_0 = 0$

$$\begin{aligned} X_{11} &= X_{12} = X_1 & y_1 + y_2 &= 0 \\ X_{21} &= X_{22} = X_2 & X_{11} + X_{21} &= 0 \\ && X_{12} + X_{22} &= 0 \end{aligned}$$

$$X_{ij}, \begin{matrix} i: 1 \rightarrow n \\ j: 1 \rightarrow p \end{matrix}$$

a) Ridge:

$$\begin{aligned} \text{min } & (y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_{12})^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2) \\ & + (y_2 - \hat{\beta}_1 x_2 - \hat{\beta}_2 x_{22})^2 \end{aligned}$$

$$\Rightarrow \text{min } (y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_{12})^2 + (y_2 - \hat{\beta}_1 x_2 - \hat{\beta}_2 x_{22})^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2)$$

derivative :

$$\begin{aligned}\frac{\partial}{\partial \hat{\beta}_1} &: -f(y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2) x_1 \\ &- f(y_2 - \hat{\beta}_1 x_2 - \hat{\beta}_2 x_1) x_2 \\ &+ \lambda \hat{\beta}_1\end{aligned}$$

$$\lambda \hat{\beta}_1 = y_1 x_1 + y_2 x_2 - \hat{\beta}_1 x_1^2 - \hat{\beta}_2 x_1 x_2 - \hat{\beta}_1 x_2^2 - \hat{\beta}_2 x_1 x_2$$

$$\textcircled{I} \quad \hat{\beta}_1(\lambda + x_2^2 + x_1^2) = y_1 x_1 + y_2 x_2 - \hat{\beta}_2(x_1^2 + x_2^2)$$

$$\textcircled{II} \quad \hat{\beta}_2(\lambda + x_2^2 + x_1^2) = y_1 x_2 + y_2 x_1 - \hat{\beta}_1(x_1^2 + x_2^2)$$

$$\textcircled{I}-\textcircled{II} \quad \lambda(\hat{\beta}_1 - \hat{\beta}_2) = 0 \quad \Rightarrow \quad \hat{\beta}_1 = \hat{\beta}_2$$

c) Lasso:

$$\min. (y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2)^2 + (y_2 - \hat{\beta}_1 x_2 - \hat{\beta}_2 x_1)^2 + \lambda(|\hat{\beta}_1| + |\hat{\beta}_2|)$$

$$\text{or} \quad \min. (y_1 - \hat{\beta}_1 x_1 - \hat{\beta}_2 x_2)^2 + (y_2 - \hat{\beta}_1 x_2 - \hat{\beta}_2 x_1)^2 \quad \text{given } |\hat{\beta}_1| + |\hat{\beta}_2| \leq s$$

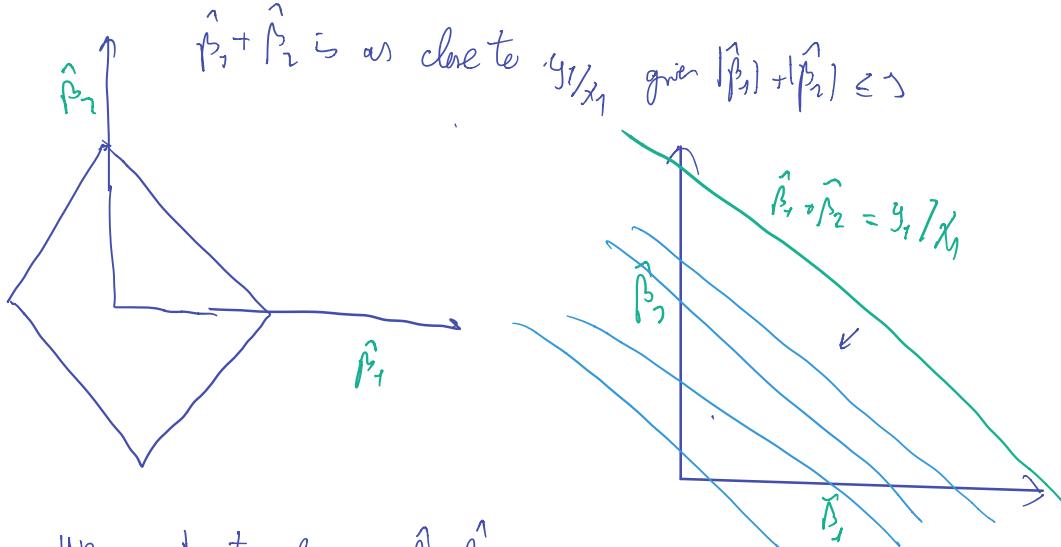
since  $x_1 = -x_2$ , this is the same as  
 $y_1 = -y_2$

$$\min. [y_1 - (\hat{\beta}_1 + \hat{\beta}_2)x_1]^2 \quad \text{given } |\hat{\beta}_1| + |\hat{\beta}_2| \leq s$$

$$\text{so min } |y_1 - (\hat{\beta}_1 + \hat{\beta}_2)x_1|$$

This is true when  $\hat{\beta}_1 + \hat{\beta}_2 = \frac{y_1}{x_1}$

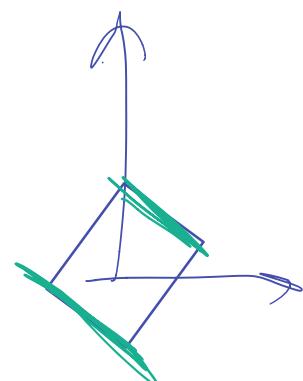
So we optimize when



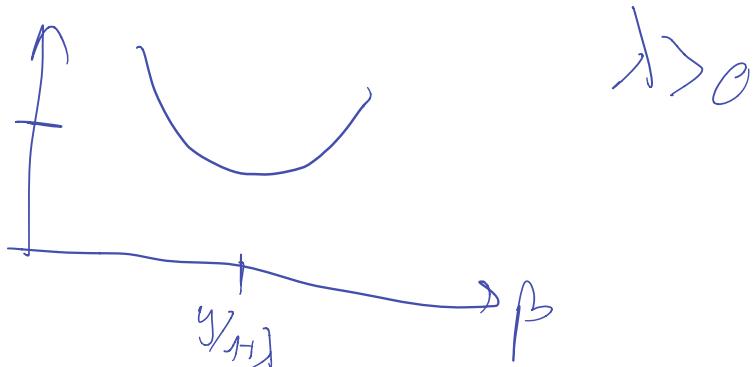
We want to have  $\hat{\beta}_1 + \hat{\beta}_2$  as big as possible but to intersect the diamond, so the solution is the first intersection w/ the diamond.

If  $\frac{y_1}{x_1} > 0$ , top right edge,

if  $\frac{y_1}{x_1} < 0$ , left bottom edge

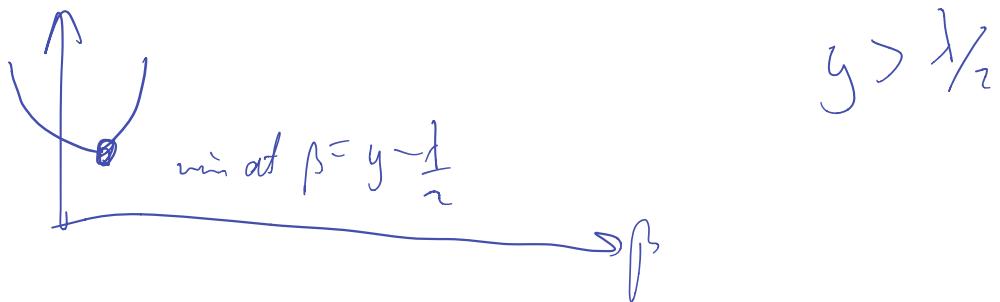


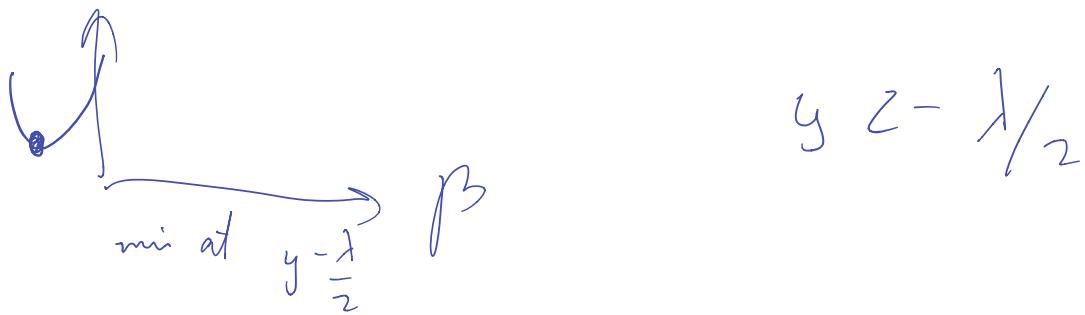
$$6) \text{ a)} (y - \beta)^2 + \lambda |f|^2 \quad \begin{matrix} P=1 \\ n=1 \end{matrix} \text{ edge}$$



$$\beta^2(\lambda+1) - 2\beta y + y^2$$

$$(b) (y - \beta)^2 + \lambda |\beta|$$





$$\text{min at } -2(y-\beta) + \lambda(-1 \cdot v + 1)$$

assume  $\beta > 0$      $\lambda - 2y + 2\beta = 0$

$$\text{min at } \beta = y - \frac{\lambda}{2}.$$

so min at  $\beta = y - \frac{\lambda}{2}$  when  $y > \frac{\lambda}{2}$

assume  $\beta < 0$      $\lambda - 2y + 2\beta = 0$

$$\beta = y + \frac{\lambda}{2} \rightarrow y < -\frac{\lambda}{2}$$

min at  $\beta = y + \frac{\lambda}{2}$  when  $y < -\frac{\lambda}{2}$

what if  $|y| < \frac{\lambda}{2}$ ?

?

-

A

$$\begin{cases} \beta^2 - (2y + \lambda)\beta + y^2 & \text{if } \beta \leq 0 \\ \beta^2 + (\lambda - 2y)\beta + y^2 & \text{if } \beta > 0 \end{cases}$$

min at  $2\beta = 2y + \lambda \rightarrow \beta = \frac{y + \lambda}{2}$

But if  $|y| < \frac{\lambda}{2}$  and  $\lambda > 0$ ,  $\beta > 0$ , so

this has  $\nexists$  min at 0

same for  $\beta > 0$ : min at  $y - \frac{\lambda}{2} < 0$ , so

this has min at 0

$\rightarrow$  min at 0 for  $(y - \beta)^2 + \lambda|\beta|$

$$\mathcal{F} \quad y_i = \beta_0 + \sum_{j=1}^p x_{ij}\beta_j + \varepsilon_i \quad \varepsilon_i \sim N(0, \sigma^2) \text{ iid}$$

likelihood:

$$P(y_i | x_{ij}, \beta_j) = \prod_j \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - (\beta_0 + \sum_j x_{ij}\beta_j))^2}{2\sigma^2}}$$

b) Assume  $\beta_1, \dots, \beta_n$  are iid  $\text{exp}(0, \theta)$

$$p(\beta) = \frac{1}{2\sigma} \exp \frac{-|\beta|}{\sigma}$$

posterior for  $\beta$ :

$$p(\beta | y_i) = \frac{p(y | \beta) p(\beta)}{p(y)} \propto p(y | \beta) p(\beta)$$

$$\propto \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^n \left( \frac{1}{2\sigma} \right) \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n \left[ Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij}) \right]^2 - \frac{|\beta|}{\sigma} \right)$$

The mode for  $\beta$  is maximizing the posterior

Take log and drop irrelevant terms:

$$\max -\frac{1}{2\sigma^2} \sum \left[ Y_i - (\beta_0 + \sum \beta_j X_{ij}) \right]^2 - \frac{|\beta|}{\sigma}$$

$$= \min \sum \left[ Y_i - (\beta_0 + \sum \beta_j X_{ij}) \right]^2 + \frac{|\beta|}{\sigma} \cdot 2\sigma^2$$

$$= \min DSS + \frac{2\sigma^2}{\sigma} \sum |\beta|$$

Take  $\lambda = \frac{2\sigma^2}{\sigma}$ , this is the Lasso solution.

$$\text{d) } \beta_1 - \beta_p \sim N(0, c)$$

posterior for  $\beta$ :

The posterior distributed according to Normal distribution with mean 0 and variance  $c$  is:

$$f(\beta | X, Y) \propto f(Y | X, \beta)p(\beta | X) = f(Y | X, \beta)p(\beta)$$

Our probability distribution function then becomes:

$$p(\beta) = \prod_{i=1}^p p(\beta_i) = \prod_{i=1}^p \frac{1}{\sqrt{2c\pi}} \exp\left(-\frac{\beta_i^2}{2c}\right) = \left(\frac{1}{\sqrt{2c\pi}}\right)^p \exp\left(-\frac{1}{2c} \sum_{i=1}^p \beta_i^2\right)$$

Substituting our values from (a) and our density function gives us:

$$\begin{aligned} f(Y | X, \beta)p(\beta) &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2\right) \left(\frac{1}{\sqrt{2c\pi}}\right)^p \exp\left(-\frac{1}{2c} \sum_{i=1}^p \beta_i^2\right) \\ &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \left(\frac{1}{\sqrt{2c\pi}}\right)^p \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2 - \frac{1}{2c} \sum_{i=1}^p \beta_i^2\right) \end{aligned}$$

e

Like from part c, showing that the Ridge Regression estimate for  $\beta$  is the mode and mean under this posterior distribution is the same thing as showing that the most likely value for  $\beta$  is given by the lasso solution with a certain  $\lambda$ .

We can do this by taking our likelihood and posterior and showing that it can be reduced to the canonical Ridge Regression Equation 6.5 from the book.

Let's start by simplifying it by taking the logarithm of both sides:

Once again, we can take the logarithm of both sides to simplify it:

$$\begin{aligned} \log f(Y | X, \beta)p(\beta) &= \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \left(\frac{1}{\sqrt{2c\pi}}\right)^p \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2 - \frac{1}{2c} \sum_{i=1}^p \beta_i^2\right) \\ &= \log \left[ \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \left(\frac{1}{\sqrt{2c\pi}}\right)^p \right] - \left( \frac{1}{2\sigma^2} \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2 + \frac{1}{2c} \sum_{i=1}^p \beta_i^2 \right) \end{aligned}$$

We want to maximize the posterior, this means:

$$\arg \max_{\beta} f(\beta | X, Y) = \arg \max_{\beta} \log \left[ \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \left(\frac{1}{\sqrt{2c\pi}}\right)^p \right] - \left( \frac{1}{2\sigma^2} \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2 + \frac{1}{2c} \sum_{i=1}^p \beta_i^2 \right)$$

Since we are taking the difference of two values, the maximum of this value is the equivalent to taking the difference of the second value in terms of  $\beta$ . This results in:

$$\begin{aligned} &= \arg \min_{\beta} \left( \frac{1}{2\sigma^2} \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2 + \frac{1}{2c} \sum_{i=1}^p \beta_i^2 \right) \\ &= \arg \min_{\beta} \left( \frac{1}{2\sigma^2} \right) \left( \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2 + \frac{\sigma^2}{c} \sum_{i=1}^p \beta_i^2 \right) \end{aligned}$$

By letting  $\lambda = \sigma^2/c$ , we end up with:

$$\begin{aligned} &= \arg \min_{\beta} \left( \frac{1}{2\sigma^2} \right) \left( \sum_{i=1}^n \left[Y_i - (\beta_0 + \sum_{j=1}^p \beta_j X_{ij})\right]^2 + \lambda \sum_{i=1}^p \beta_i^2 \right) \\ &= \arg \min_{\beta} \text{RSS} + \lambda \sum_{i=1}^p \beta_i^2 \end{aligned}$$

Since we know our posterior is normal, mean equals median