

Regularized regression

$$y = X\beta + \epsilon, \quad X \text{ } n \times p \text{ design matrix}$$

- y & X 's are correlated \Rightarrow high variance for $\hat{\beta}$'s
- \Rightarrow correlated $\hat{\beta}$'s
- \Rightarrow Difficult to interpret model

$$\rightarrow \text{Source of problem? } \hat{\beta} = (X'X)^{-1}X'y$$

↙
This inverse numerically unstable
when p large comp n and/or
 X 's are correlated

What to do? → Regularize the fit

Principal Component Regression

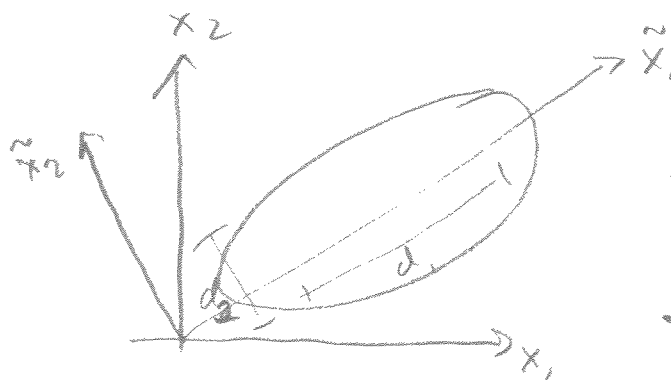
$$X = UDV' \quad \text{Singular value decomposition}$$

$$U'U = I$$

$$V'V = I$$

$$X'X = VD^2V'$$

↙
 $\text{cov}(X)$



- \tilde{X}_1 and \tilde{X}_2 are uncorrelated

- $V(\tilde{X}_2) = d_2^2 < V(\tilde{X}_1) = d_1^2$

- X_1 and X_2 are correlated

- direction v_1 (1st column in V)

- dir. maximum X -spread in this direction

- $V(\tilde{X}_1) = d_1^2$

So PC = new X-variables such that

(2)

$$V(\tilde{x}_1) > V(\tilde{x}_2) > \dots$$

If we rotate X-variables \Rightarrow \tilde{X} -variables

The new variables are uncorrelated \Rightarrow new $\hat{\tilde{\beta}}$'s are uncorrelated

$$y = \tilde{X} \tilde{\beta} + \varepsilon \Rightarrow \hat{\tilde{\beta}} = (\tilde{X}' \tilde{X})^{-1} \tilde{X}' y, \quad \tilde{X} = XV$$

$$\tilde{\beta} = V' \beta$$

\Downarrow

$$V(\hat{\tilde{\beta}}) = \sigma^2 (\tilde{X}' \tilde{X})^{-1}$$

$$= \sigma^2 (V' X' X V)^{-1} = \sigma^2 (V' V D^2 V' V)^{-1}$$

$$= \sigma^2 D^{-2}$$

\uparrow diagonal, so all $\hat{\tilde{\beta}}$'s are

uncorrelated and

$$V(\hat{\tilde{\beta}}_1) < V(\hat{\tilde{\beta}}_2) < \dots$$

$$\sigma^2 d_1^{-2} \quad \sigma^2 d_2^{-2} \quad \dots$$

In new coordinate system, the 1st variable \tilde{x}_1 has the largest spread and so its coefficient estimate $\hat{\tilde{\beta}}_1$ is the most precise.

To regularize - just use the first k components with large spread.

Pros of PC regression - simple to use

(3)

- keep only \tilde{x} -variables for which we can estimate $\tilde{\beta}$ w high precision

(con) - each new \tilde{x} -variable = linear combo of all x -variables \Rightarrow difficult to interpret.

Regularized regression

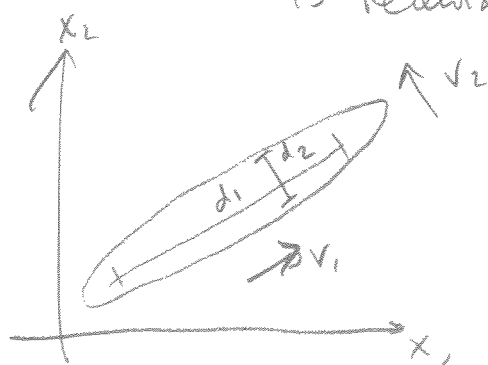
$$\text{Note } V(\hat{\beta}) = \sigma^2 (X'X)^{-1} = \sigma^2 V D^{-2} V' = \sigma^2 \sum_{k=1}^p d_k^{-2} v_k v_k'$$

spectral decomposition

$$\text{Total MSE}(\hat{\beta}) = \text{Trace}\{V(\hat{\beta})\} = \sigma^2 \text{Tr}\{(X'X)^{-1}\} = \sigma^2 \sum_{k=1}^p d_k^{-2}$$

large if any d_k is small.

When does that happen? Whenever there are x 's that are highly correlated so information in X is redundant.



Here d_2 is small

So LS estimator

④

• $E(\hat{\beta}) = \beta$ so unbiased

• $V(\hat{\beta}) = \sigma^2 (X'X)^{-1}$ — can be large if any x 's are correlated

Can we reduce variance by allowing for some bias?

Source of high variance = numerical instability $\Rightarrow X'X$ inverse

\rightarrow stabilize it $\rightarrow \hat{\beta}_R = (X'X + rI)^{-1} X'y$

\uparrow

add r to the diagonal of $X'X$ before taking inverse

artificially suppress

correlation between x 's

What are $\hat{\beta}_R$ properties?

$$\begin{aligned} \hat{\beta}_R = (X'X + rI)^{-1} X'y &\rightarrow E(\hat{\beta}_R) = (X'X + rI)^{-1} X'X\beta \\ &= (X'X + rI)^{-1} (X'X + rI - rI)\beta \\ &= \beta - r(X'X + rI)^{-1}\beta \end{aligned}$$

bias

(controlled by r
if $r=0$, no bias

depends on β
large $\beta \rightarrow$ large bias

$$V(\hat{\beta}_R) = \sigma^2 (X'X + rI)^{-1} X'X (X'X + rI)^{-1}$$

(5)

$$MSE(\hat{\beta}_R) = \text{Var} + \text{bias}^2$$

$$= (X'X + rI)^{-1} (\sigma^2 X'X + r^2 \beta\beta') (X'X + rI)^{-1}$$

$$\text{Total MSE} = \text{Trace}\{MSE(\hat{\beta}_R)\}$$

$$= \text{Trace}\{V(D^2 + rI)^{-1} V' [\sigma^2 V D^2 V' + r^2 \beta\beta'] V (D^2 + rI)^{-1} V'\}$$

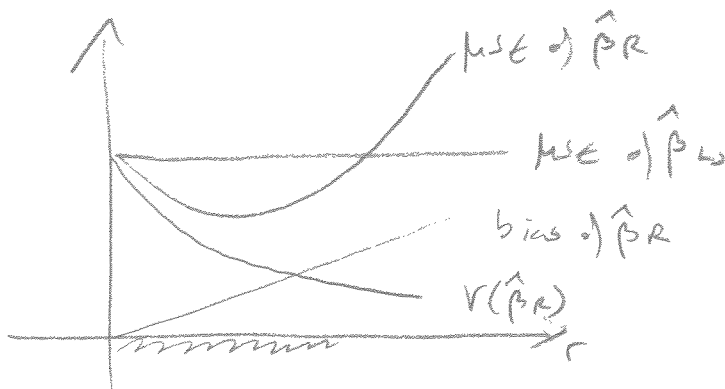
$$= \text{Trace}\{V(D^2 + rI)^{-1} [\sigma^2 D^2 + r^2 \tilde{\beta}\tilde{\beta}'] (D^2 + rI)^{-1} V'\}$$

diagonal matrix

$$= \sum_{k=1}^p \frac{\sigma^2 d_k^2 + r^2 \tilde{\beta}_k^2}{(d_k^2 + r)^2} < \sum_{k=1}^p \sigma^2 d_k^{-2} \quad \text{for some } r$$

MSE of $\hat{\beta}_R$

MSE of $\hat{\beta}_W$



r-values for which $\hat{\beta}_R$ better than $\hat{\beta}_W$

Use CV to find r in practise. — best prediction error.

Another way of looking at this;

6

Shrinkage estimator

What does adding r to diagonal of $X'X$ do?

if X 's are uncorrelated st. $X'X = D^2$

$$\Rightarrow (X'X + rI) = (D^2 + rI) = \begin{pmatrix} d_1^2 + r & & \\ & d_2^2 + r & \\ & & \ddots \\ & & & d_p^2 + r \end{pmatrix}$$

$(X'X + rI)^{-1} X'y =$ each $\hat{\beta}_j$ shrunk by a factor r

General ; $\hat{\beta}_S = \frac{\hat{\beta}_{OLS}}{1+c} \rightarrow$ bias $\frac{c\beta}{1+c}$
 $c > 0$

\rightarrow variance $\frac{V(\hat{\beta}_{OLS})}{(1+c)^2} < V(\hat{\beta}_{OLS})$

if X 's are not uncorrelated

$$X'X + rI = V D^2 V' + rI = V (D^2 + rI) V'$$

affects mostly the

PC directions where

d_k^2 is small.

Idea behind shrinkage is that by not allowing $\hat{\beta}$'s to get too big we cover ourselves from risk of poor estimates having a huge impact on prediction - just make everything smaller = bias[↑] some and variance ↓ a lot.

We can arrive at this solution by doing

penalized regression

$$LS: \min_{\beta} (y - X\beta)'(y - X\beta) = \|y - X\beta\|^2 = \sum_{i=1}^n (y_i - x_i' \beta)^2$$

$$Penalized LS: \min_{\beta} (y - X\beta)'(y - X\beta)$$

subject to $\|\beta\|^2 \leq \tau$

same thing as $\beta' \beta \leq \tau$

$$- \dots - \sum_{j=1}^p \beta_j^2 \leq \tau$$

So, we want to fit the data to a model using LS

but we only accept solutions where the total

magnitude of β^2 's is less than τ - either all small - or one big $\leq \tau$ rest 0 and anything in between.

How to solve constrained optimization problems?

$$\min_{\beta} (y - X\beta)' (y - X\beta)$$

⇒ Lagrange multipliers

subject to $\beta' \beta \leq \tau$

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

→ solve this for different values of λ and find the smallest λ so that solution satisfies the constraint $\beta' \beta \leq \tau$

if $\lambda = 0 \Rightarrow$ LS solution \Rightarrow if $\hat{\beta}_{LS}' \hat{\beta}_{LS} \leq \tau$ done!

if not

increase $\lambda \Rightarrow$ if $\hat{\beta}' \hat{\beta} \leq \tau$ done!

if not
increase λ

Solution gives λ :

(9)

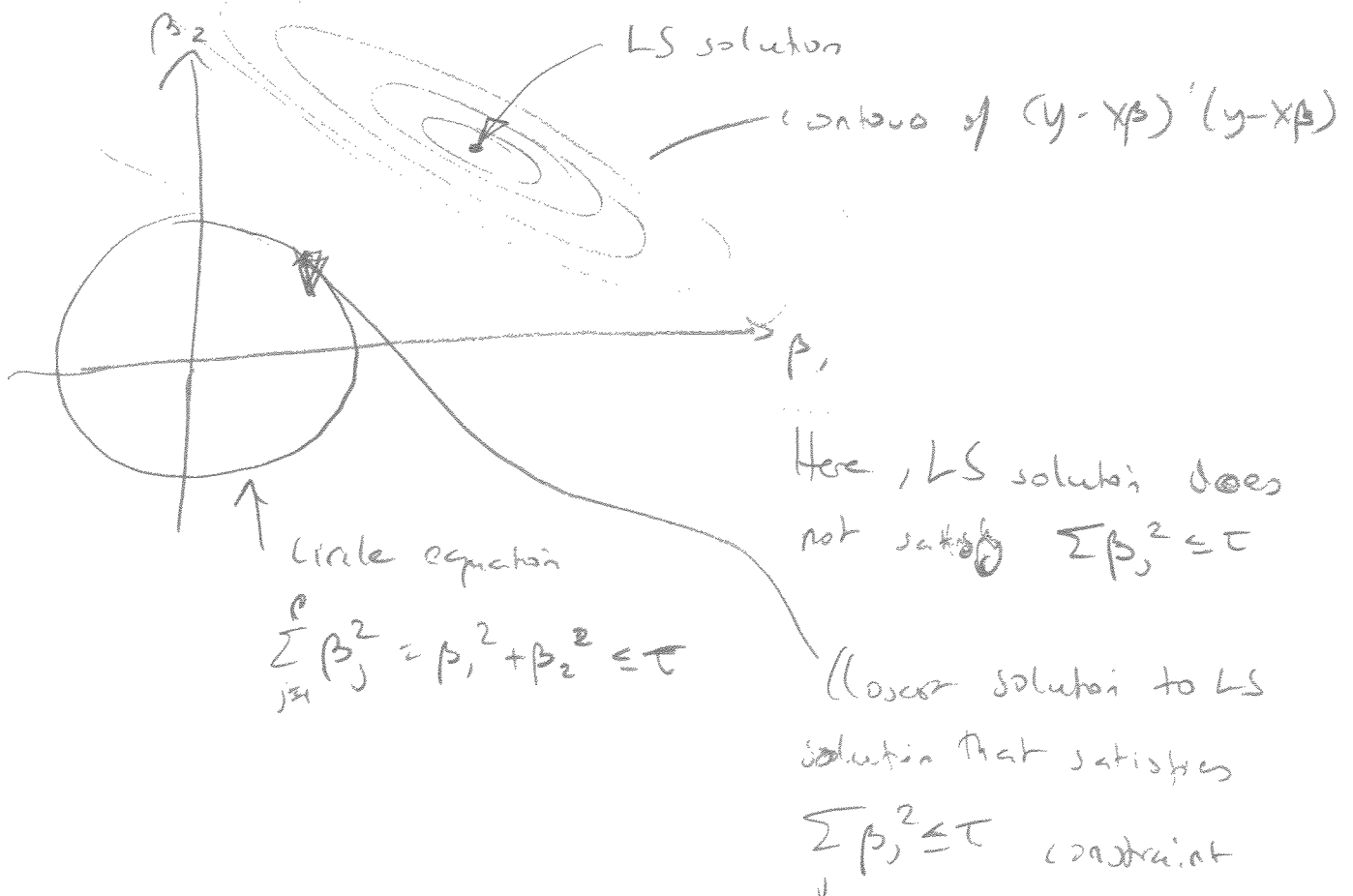
$$\min_{\beta} (y - X\beta)'(y - X\beta) + \lambda \beta' \beta$$

$$\frac{\partial}{\partial \beta} \quad \quad \quad = 0$$

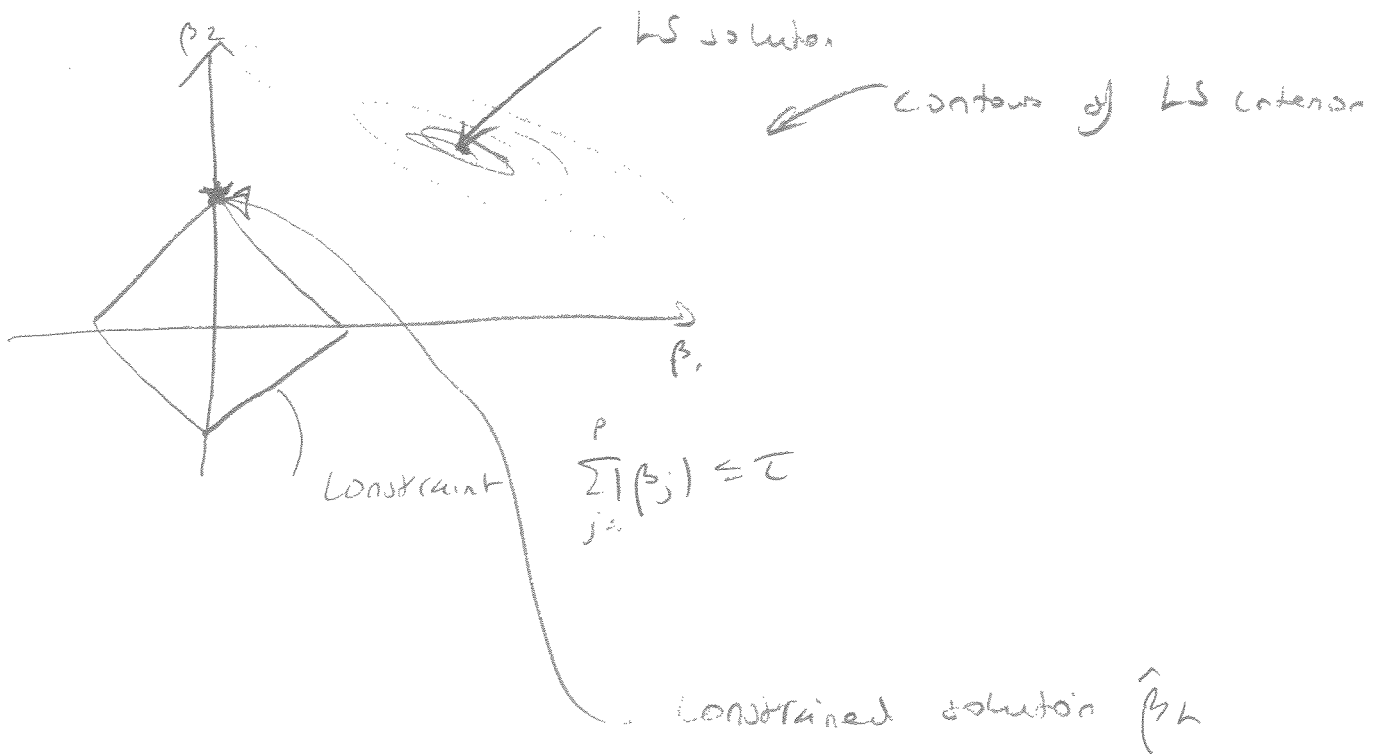
$$\Rightarrow -2X'(y - X\beta) + 2\lambda\beta = 0$$

$$\Rightarrow (X'X + \lambda I)\beta = X'y \Rightarrow \hat{\beta} = (X'X + \lambda I)^{-1} X'y$$
$$= \hat{\beta}_R \quad \checkmark$$

Another way of looking at it;



Recently — explore the use of other constraint shapes (10)

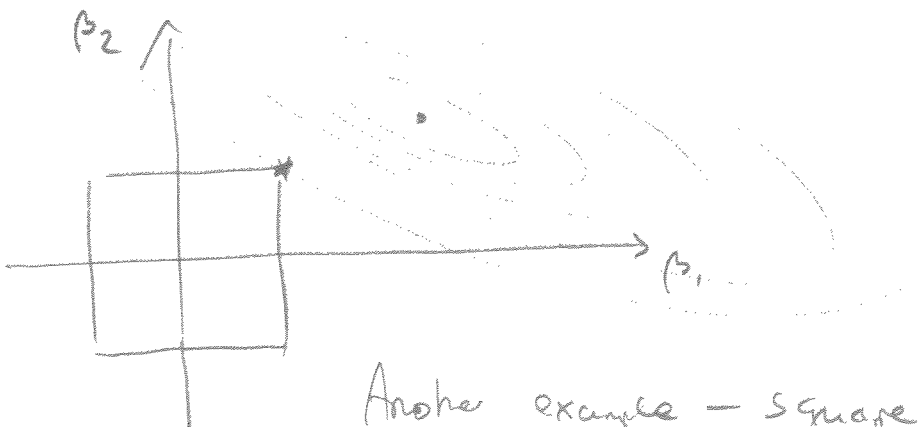


Using this diamond-shaped constraint is called LASSO

It does automatic model selection.

Since constraint region has sharp angles on the axis of β -space, we encourage solutions where some $\hat{\beta} = 0$

Above example $\hat{\beta}_1 = 0$ and $\hat{\beta}_2 \neq 0$.



Another example — square constraint

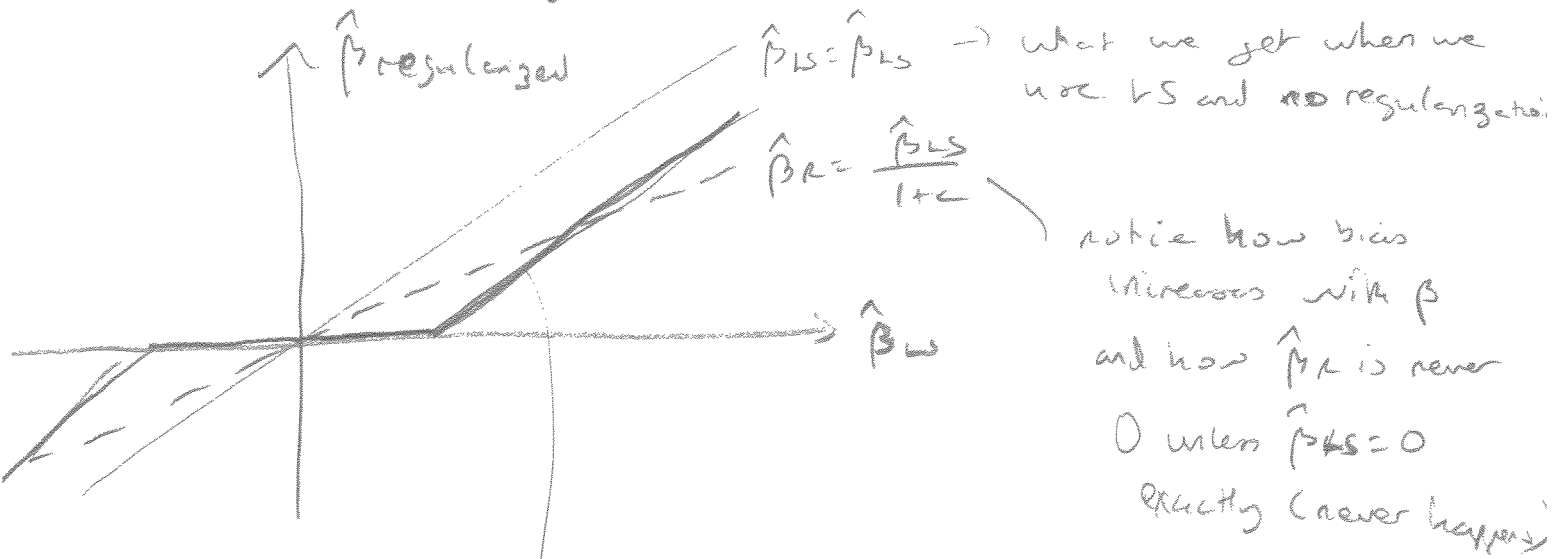
$$\max(\beta_1, \beta_2) \leq \tau$$

\Rightarrow encourages $\hat{\beta}_1 = \hat{\beta}_2$

Recently - lots of work in statistics to come up with automated model selection using interesting constraint regions like these.

(11)

Another look at the regularized estimates



$$\hat{\beta}_L = (|\hat{\beta}_{LS}| - \delta)_+ \text{sign}(\hat{\beta}_{LS})$$

$$= \begin{cases} \hat{\beta}_{LS} - \delta & \text{i) } \hat{\beta}_{LS} > \delta \\ 0 & \text{ii) } |\hat{\beta}_{LS}| \leq \delta \\ \hat{\beta}_{LS} + \delta & \text{i) } \hat{\beta}_{LS} < -\delta \end{cases}$$

= soft thresholding of $\hat{\beta}_{LS}$

- notice bias is constant