

Last time: Ridge RegressionOLS:

$$\hat{\beta}^{(ols)} = \underset{\beta}{\operatorname{argmin}} L(\beta)$$

$$L(\beta) = \|Y - X\beta\|^2$$

$$= (X^T X)^{-1} X^T Y$$

Ridge: Shrinkage Estimator

$$(1) \quad \hat{\beta}^{(ridge)} = \underset{\beta}{\operatorname{argmin}} L(\beta) + \lambda \|\beta\|_2^2$$

$\lambda \geq 0$

$$(2) \quad \hat{\beta}^{(ridge)} = \underset{\beta}{\operatorname{argmin}} L(\beta) \quad \text{s.t.} \quad \|\beta\|_2^2 \leq t$$

corresp. t

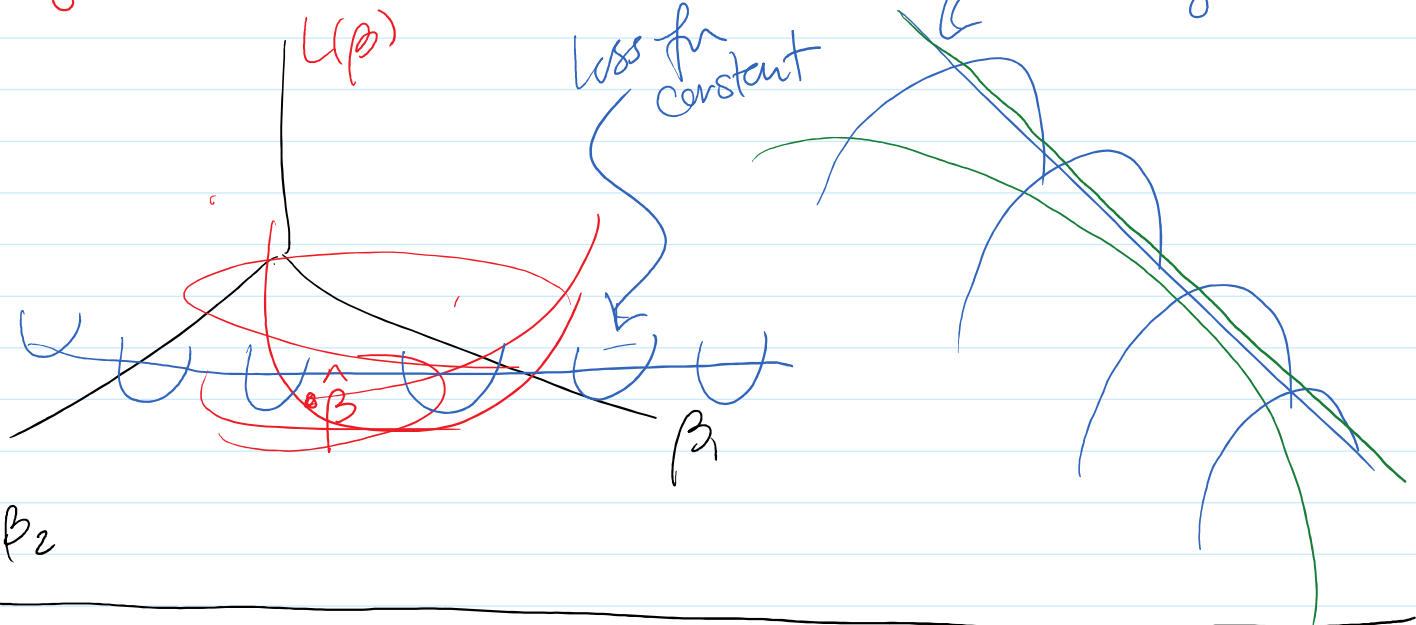
$$(3) \quad \hat{\beta}^{(ridge)} = (X^T X + \lambda I)^{-1} X^T Y$$

Why called ridge regression

Imagine $K(X^T X) = 0$
 $\lambda = 0$

ridge

Imagine $K(X, X) = \omega$



Why do variable selection

→ highly correlated vars (ridge)

→ interpretability (look at this) (LASSO)

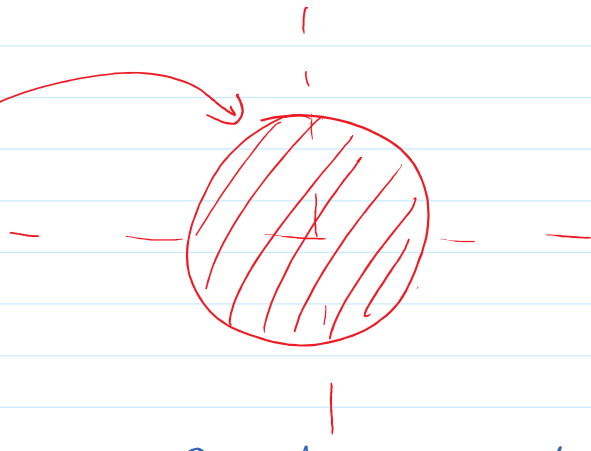
Aside: Norms

Ex. euclidean norm

$$x \in \mathbb{R}^p \text{ then } \|x\| = \sqrt{\sum_{i=1}^p x_i^2}$$

Consider

$$\{y \mid \|y\| \leq 1\}$$



can generalize: let $q \geq 0$ define the q -norm as

$$\|x\|_q = \left(\sum_{i=1}^p |x_i|^q \right)^{1/q}$$

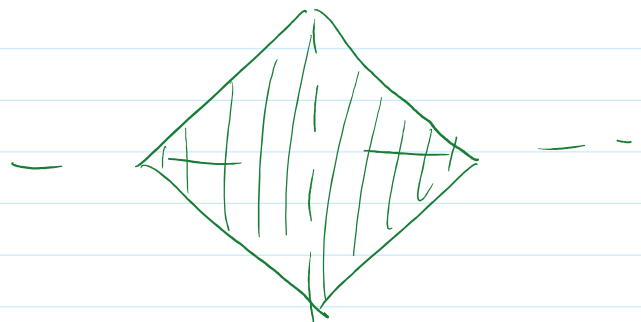
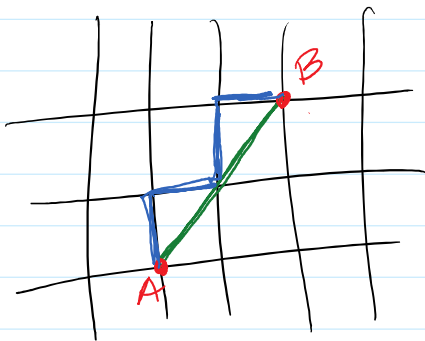
when $q=2$ then $\|x\|_2$ is euclidean norm

$q=1$ then $\|x\|_1$ is called the

Manhattan norm (1-norm)

$$\{y \mid \|y\|_1 \leq 1\}$$

$$\|x\|_1 = \sum_{i=1}^p |x_i|$$

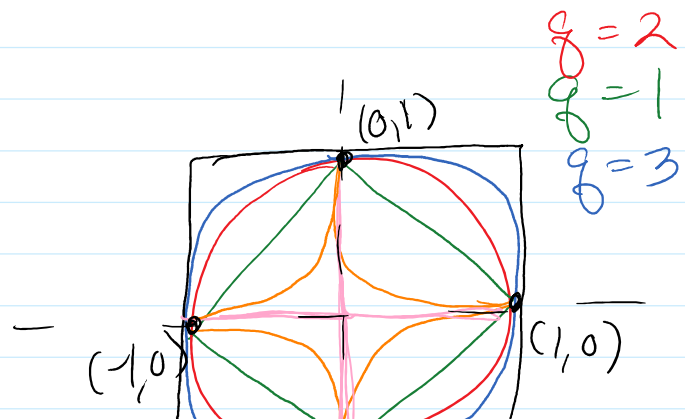


$$d_1(A, B) = \|A - B\|_1$$

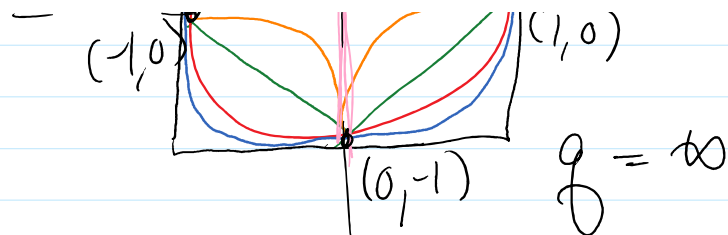
$$d_2(A, B) = \|A - B\|_2$$

$$\|x\|_{\infty} = \lim_{q \rightarrow \infty} \|x\|_q$$

$$= \max |x_i|$$



$$= \max_i |x_i|$$



$$\|x\|_0 = \lim_{g \rightarrow 0} \|x\|_g = \#\{x_i \neq 0\}$$

$$g = 1/2$$

↑
of elems
of x non-zero

$$\|(0,1)\|_0 = 1 \quad \|(1/2, 1/2)\| = 2 \text{ etc.}$$

$$\|(0,0)\| = 0$$

LASSO: Least-Absolute Shrinkage and Selection Operator

Variable selection is like zeroing out entries of β

$$\hat{Y} = X \hat{\beta} = X_1 \hat{\beta}_1 + X_2 \hat{\beta}_2 + \dots$$

↑ force to zero
drop X_2 from
model

Want:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} L(\beta) \quad \text{s.t.} \quad \|\beta\|_0 \leq t$$

of non-zero elements of $\beta \leq t$

Unfortunately, this $\|\cdot\|_0$ not convex
this is intractable problem.

Can't do this directly.

Relax the problem (convex relaxation)
using $q=1$

LASSO:

$$(1) \quad \hat{\beta}^{(\text{LASSO})} = \underset{\beta}{\operatorname{argmin}} L(\beta) \quad \text{s.t.} \quad \|\beta\|_1 \leq t$$

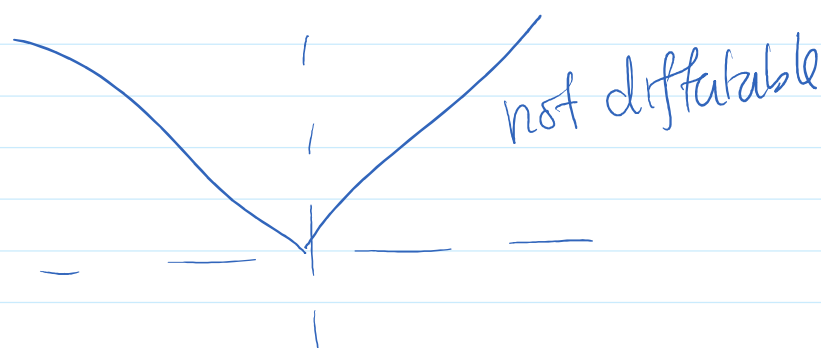
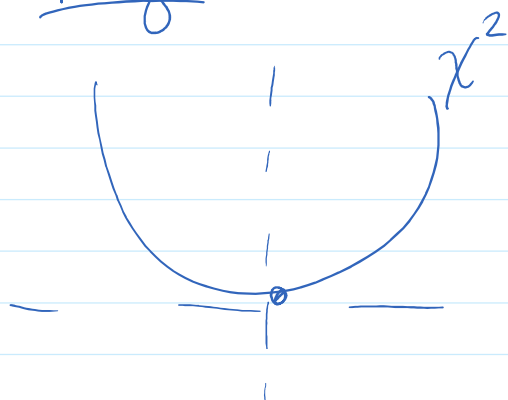
$$(2) \quad \hat{\beta}^{(\text{LASSO})} = \underset{\beta}{\operatorname{argmin}} L(\beta) + \lambda \|\beta\|_1$$

(3) amazing fact: $\|\beta\|_1$ continuous, convex,
not differentiable
So no analytic solution

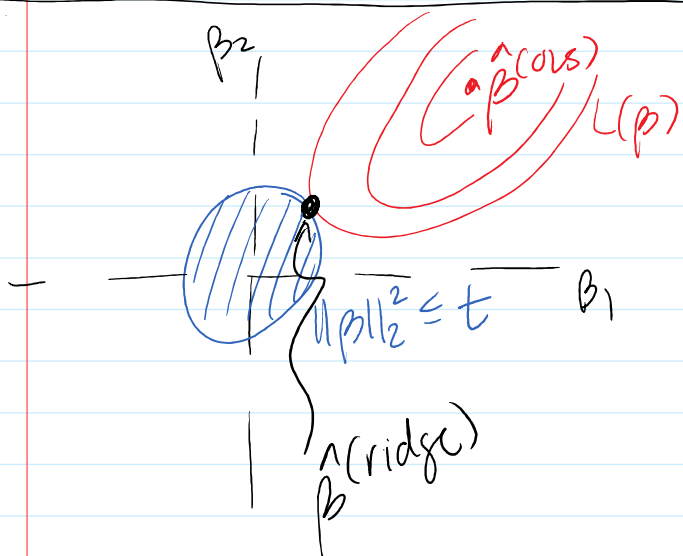
think of:

ridge: x^2

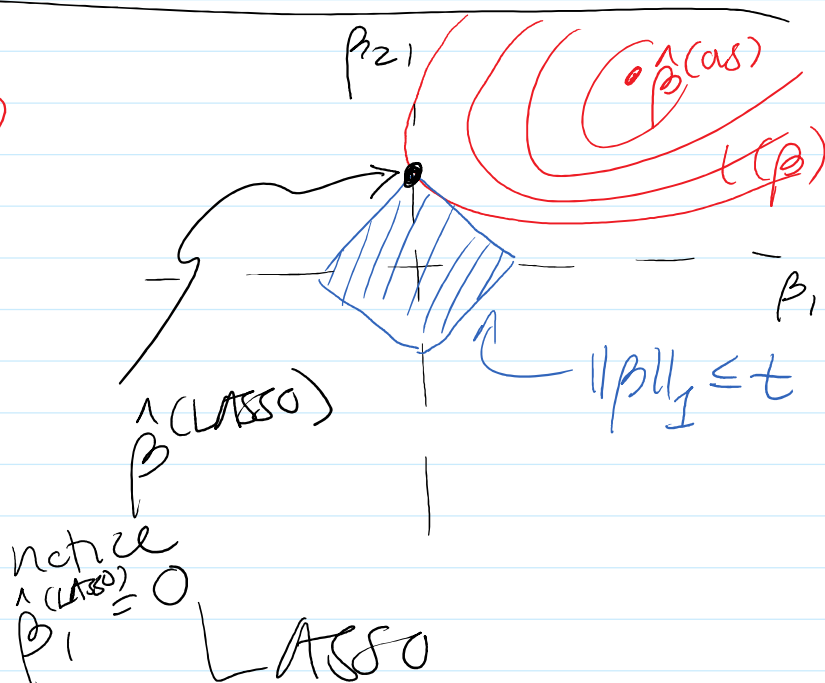
v. LASSO, $|x|$



punchline: no simple formula for $\hat{\beta}^{(LASSO)}$
need to numerically solve w/ convex optim.



Ridge



Ridge

$\hat{\beta}_i^{(LASSO)} \equiv 0$ LASSO

Comparison: Assume X is orthogonal
(vars uncorrelated)

① Subset selection (Hard thresholding)

$$\hat{\beta}_i^{(HS)} = \hat{\beta}_i^{(OLS)} \mathbb{1}(|\hat{\beta}_i^{(OLS)}| \geq t) = \begin{cases} \hat{\beta}_i & \text{if large enough} \\ 0 & \text{else} \end{cases}$$

② Ridge:

$$\hat{\beta}^{(ridge)} = \frac{\hat{\beta}_i^{(OLS)}}{1 + \lambda} \quad (\text{proportional shrinkage})$$

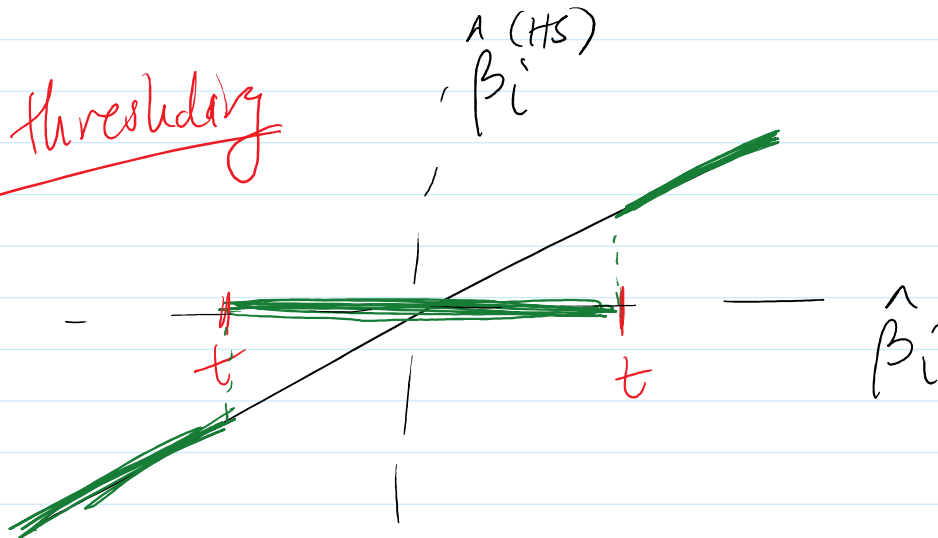
③ LASSO:

soft threshold

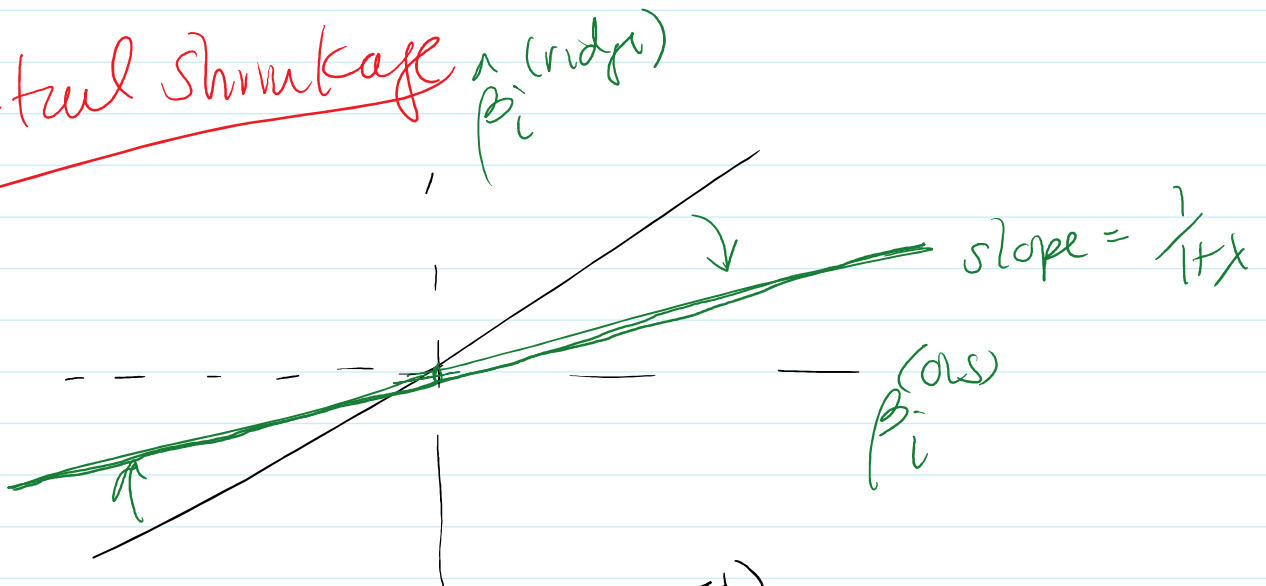
$$\begin{aligned} \hat{\beta}_i^{(LASSO)} &= \text{sign}(\hat{\beta}_i^{(OLS)}) \max(|\hat{\beta}_i^{(OLS)}| - \lambda, 0) \\ &= \begin{cases} \hat{\beta}_i^{(OLS)} - \lambda & \text{if this is } \geq 0 \end{cases} \end{aligned}$$

$$= \begin{cases} \beta^{(OLS)} - \lambda & \text{if this is } \geq 0 \\ 0 & \text{else} \end{cases} \quad \text{for } \beta^{(OLS)} > 0$$

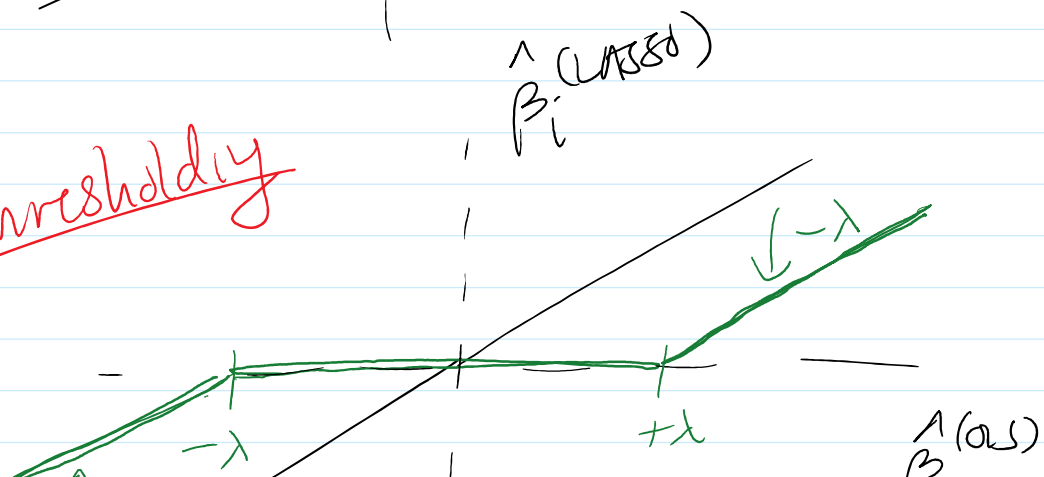
Hard thresholding

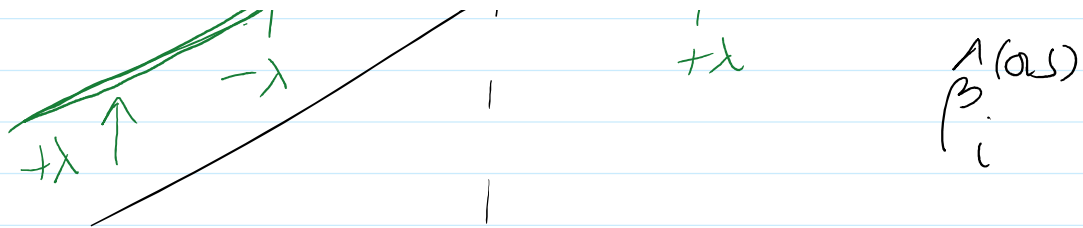


Proportional shrinkage



Soft Thresholding





① Can combine ridge/Lasso

Elastic Net

$$0 \leq \alpha \leq 1$$

$$\hat{\beta}^{(EN)} = \underset{\beta}{\operatorname{argmin}} \left[L(\beta) + \lambda \left[(1-\alpha) \|\beta\|_2^2 + \alpha \|\beta\|_1 \right] \right]$$

$\alpha=0 \Rightarrow \text{ridge}$

$\alpha=1 \Rightarrow \text{Lasso}$

② Can do penalization w/ any method.

→ penalize my Loss

→ Ex, penalized Logistic regression

→ ...