

## Lecture 16

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

$$\|x\|_0 = \lim_{q \rightarrow 0} \|x\|_q = \#\{x_i \neq 0\}$$

↑ how many non-zero elements of  $x$

---

LASSO : Least-Absolute Shrinkage and Selection Operator

Variable selection is like forcing some of my  $\hat{\beta}$ s to zero.

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots$$

↑ force  $\hat{\beta}_1 = 0$  then I effectively select out  $X_1$

Consider the constrained optim problem

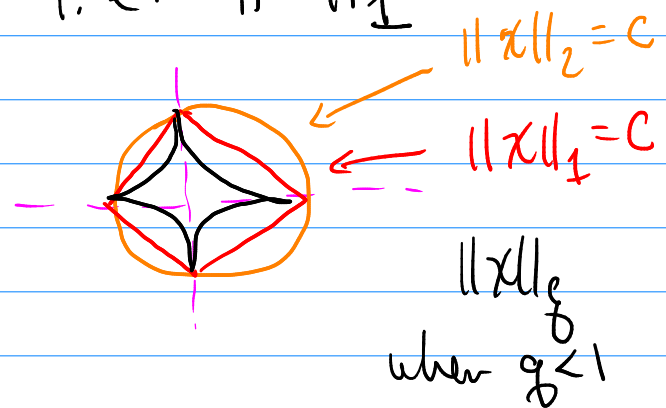
$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} L(\beta) \quad \text{s.t.} \quad \|\beta\|_0 \leq t$$

↑ build best model w/  $\leq t$  variables

Problem: generally optimizing using a  $\|\cdot\|_0$  constraint is difficult / intractable

not convex

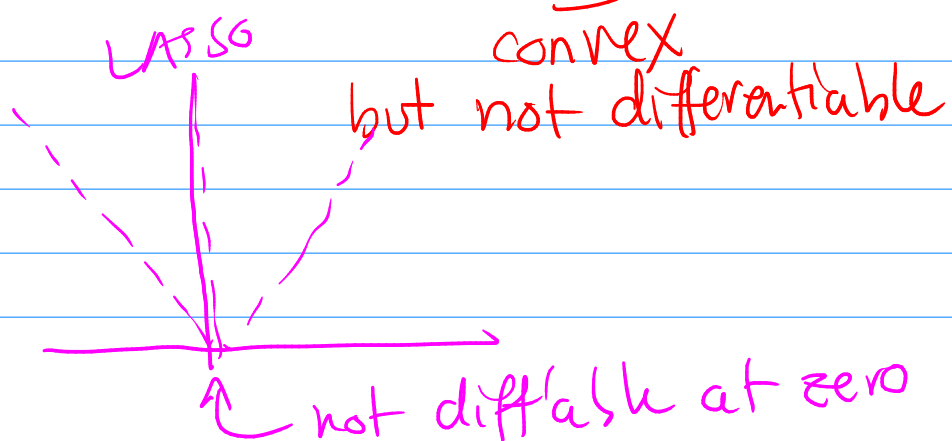
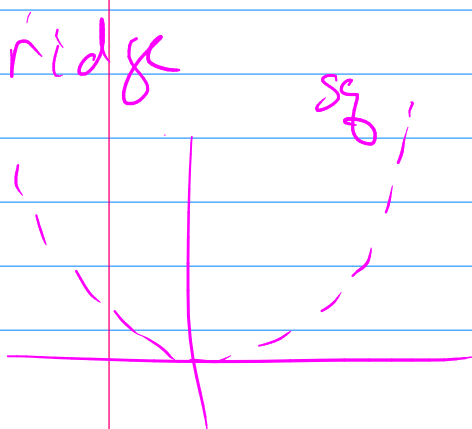
So instead, we can work w/ the convex-relaxation of  $\|\cdot\|_0$  i.e.  $\|\cdot\|_1$



LASSO:

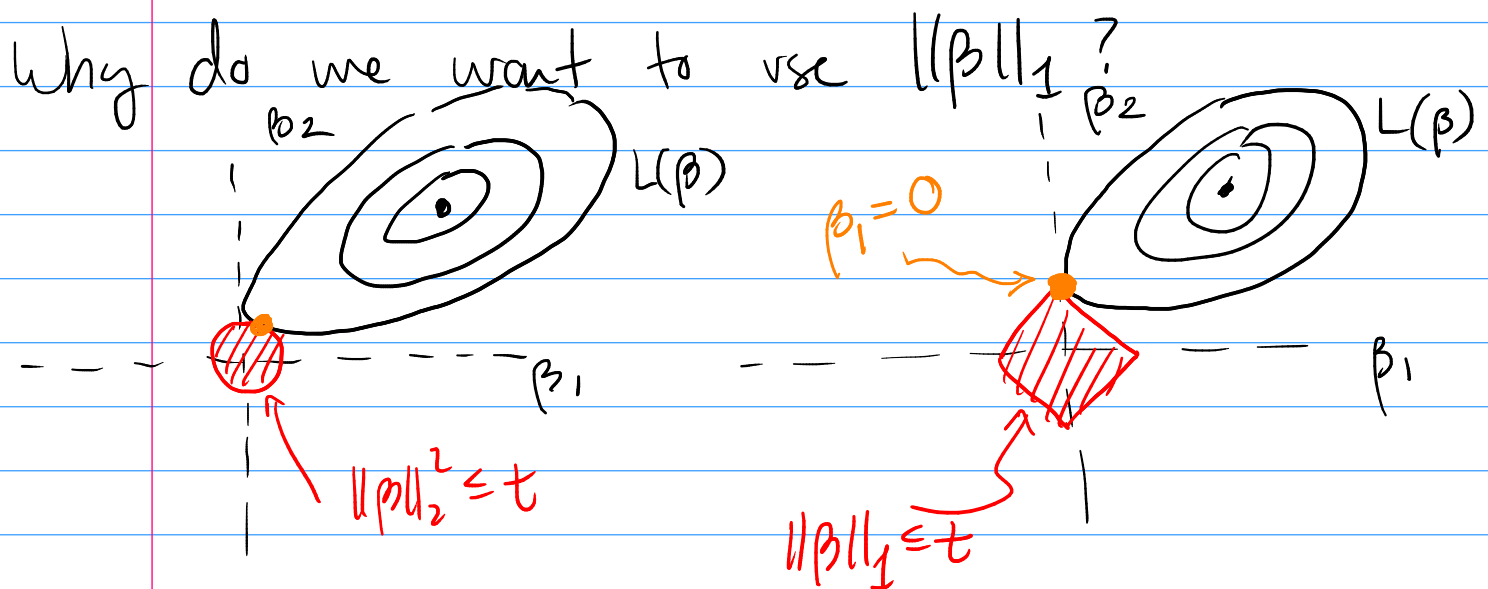
①  $\hat{\beta}^{(LASSO)} = \underset{\beta}{\operatorname{argmin}} L(\beta) \text{ s.t. } \|\beta\|_1 \leq t$

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



b/c  $\|\beta\|_1$  isn't diff'able there is  
no closed-form soln for  $\hat{\beta}^{\text{LASSO}}$

So need to use some numerical method.



Ridge

LASSO

LASSO tends to force  
Constr. optimum at points of  
Constraint i.e. zero out elements  
of  $\beta$

Comparison Assume that  $X$  is orthogonal.

① Variable Selection (Hard-thresholding)

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

② Ridge:

$$\hat{\beta}^{ridge} = \frac{\hat{\beta}^{OLS}}{1 + \lambda}$$

proportional shrinkage

③ LASSO:

$$\hat{\beta}_i^{LASSO}$$

soft thresholding

$$\left\{ \begin{array}{ll} \text{sign}(\hat{\beta}_i^{OLS}) \left[ |\hat{\beta}_i^{OLS}| - \lambda \right] & \text{if } |\hat{\beta}_i^{OLS}| - \lambda > 0 \\ 0 & \text{if } |\hat{\beta}_i^{OLS}| - \lambda \leq 0 \end{array} \right.$$

---

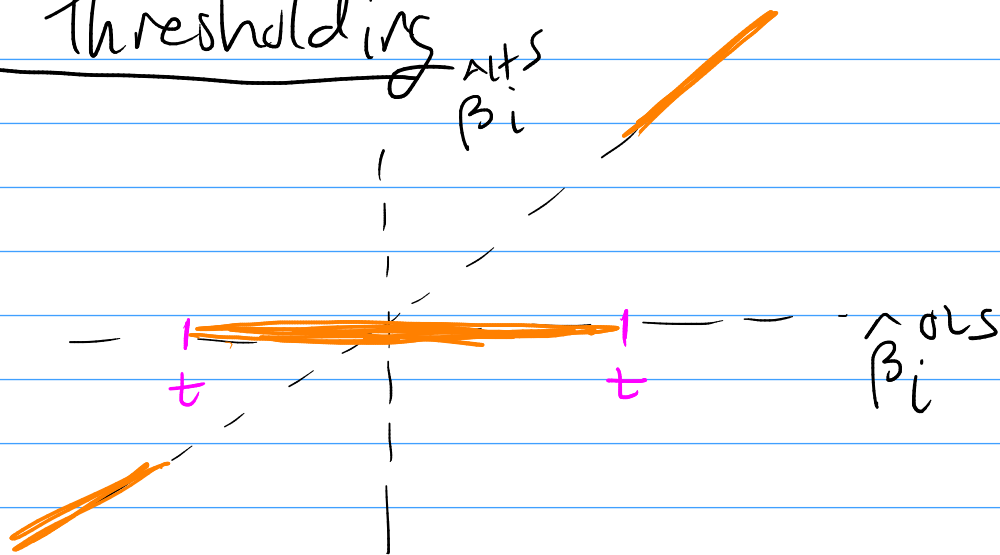
Elastic Net :  $L1 + L2$

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

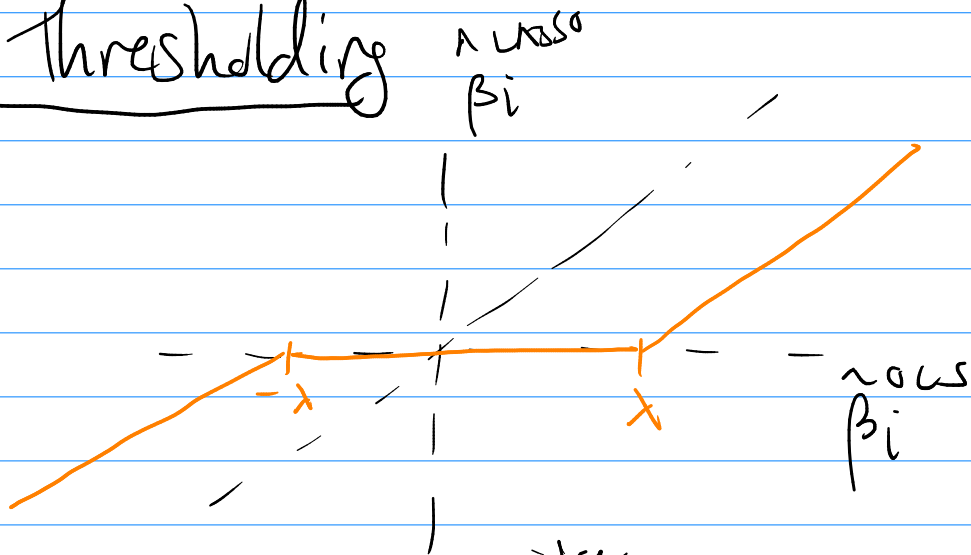
$\alpha = 0 \Rightarrow \text{ridge}$   
 $\alpha = 1 \Rightarrow \text{LASSO}$

$0 \leq \alpha \leq 1$   
tradeoff between  
 $L1 / L2$  penalty

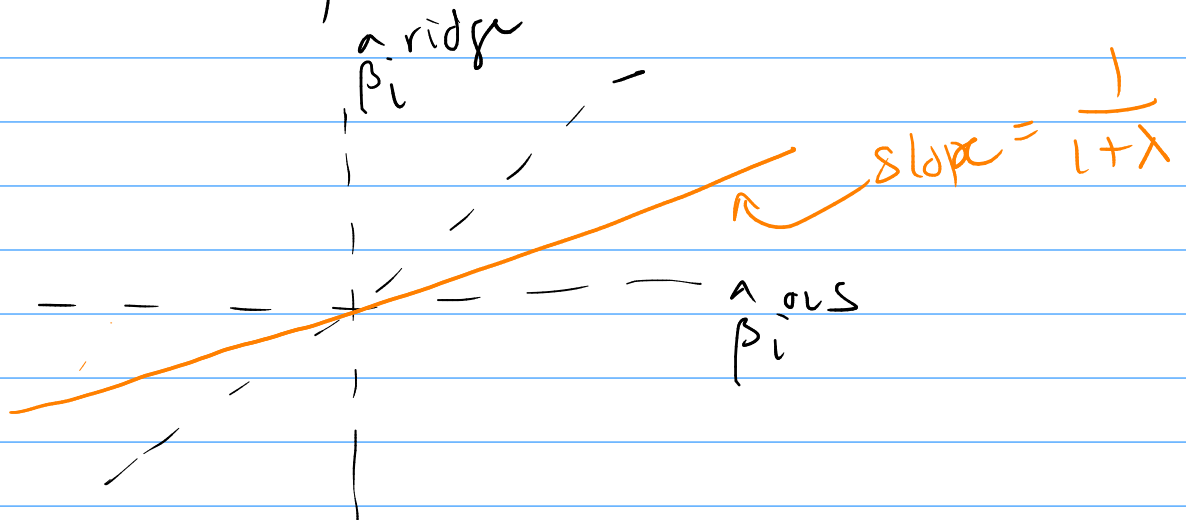
# Hard thresholding



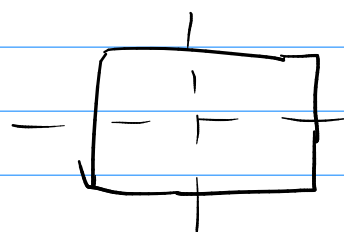
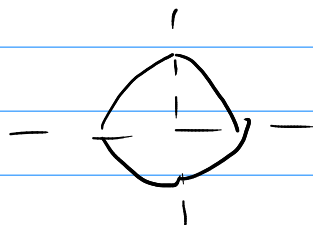
# Soft thresholding



# Ridge



$$q = 1.5$$



Can do this w/ all sorts of other methods.

Ex. Logistic Regression

$$y_n | \underline{X}_n = \underline{x}_n \sim \text{Bern}(p(\underline{x}_n))$$

$$p(\underline{x}_n) = \text{logistic}(X\beta)$$

like regr.

(neg.)  
 $l(\beta) = \text{log-likelihood fn}$

$$\hat{\beta} = \underset{\beta}{\text{argmin}} l(\beta) + \lambda \|\beta\|^2$$

or  $+ \lambda \|\beta\|_1$

---

More generally,

$$\hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} L(f)$$

could consider a penalized version

$$\hat{f} = \underset{f \in \mathcal{F}}{\text{argmin}} L(f) + \lambda J(f)$$

data term

penalty term

penalty fn for  $f$   
that measures complexity