# LASSO: Least-Absolute Shrinkage and Selection Operator

Ideally I can solve the problem

$$\hat{\beta} = \underset{\beta}{\arg\min} \; L(\beta) \quad s.t. \quad \|\beta\|_0 \leq t$$

SE loss

at most $t$ vars

$$= \text{building the best } t\text{-variable model}$$

(best subset selection problem)
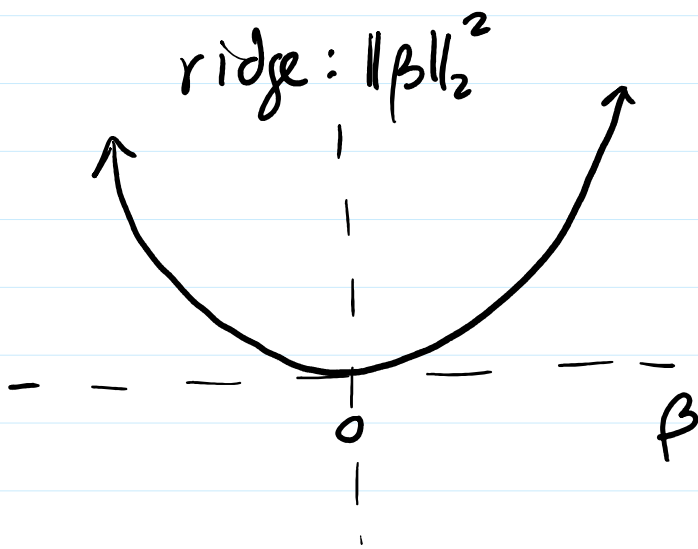
**Problem**: comp. intensive

Optimizing under $\|\cdot\|_0$ constraint is difficult b/e $\|\beta\|_0$ is neither diff'able nor convex
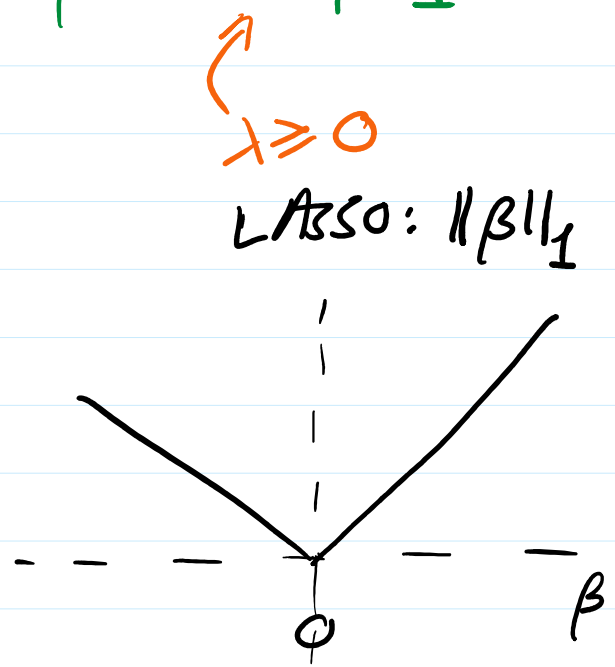
diff'able nor convex

LASSO : to make workable use the L1
norm instead — convex relaxation

(1) $\hat{\beta}^{(LASSO)} = \underset{\beta}{\arg\min} \; L(\beta) \; s.t. \; \|\beta\|_1 \leq t$

(2) $\hat{\beta}^{(LASSO)} = \underset{\beta}{\arg\min} \; L(\beta) + \lambda \|\beta\|_1$

$\lambda \geq 0$

ridge: $\|\beta\|_2^2$

LASSO: $\|\beta\|_1$



$\beta$

$\beta$
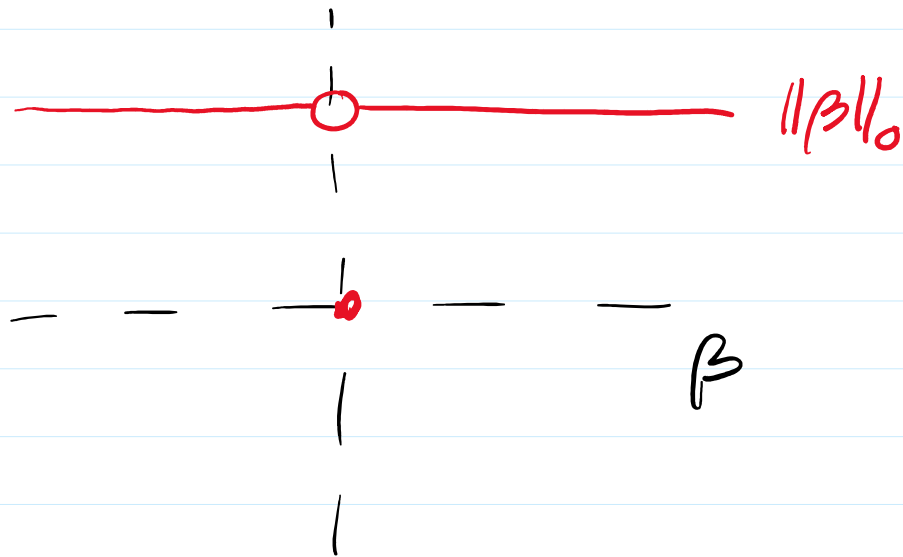
— convex
— diff'able

— convex
— not diff'able

Since $\|\cdot\|_1$ isn't diff'able, no closed form soln, need to use numerical methods.



$\|\beta\|_0$

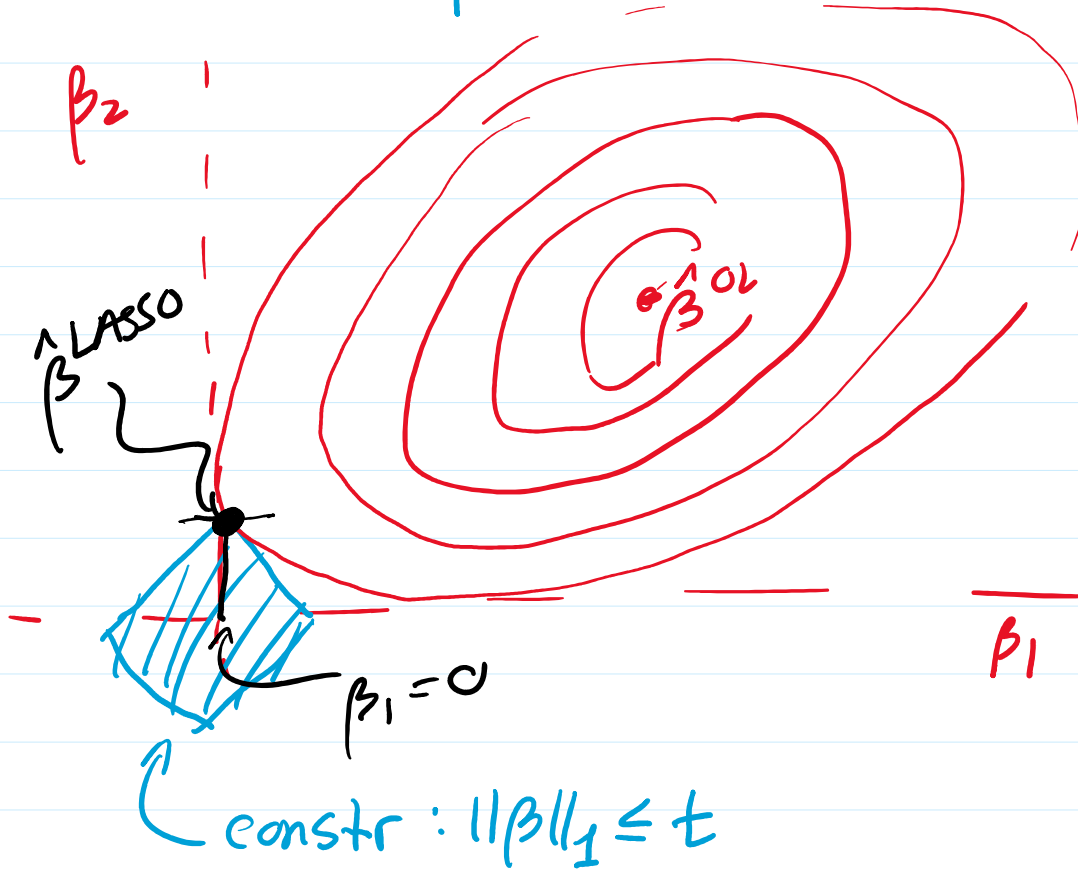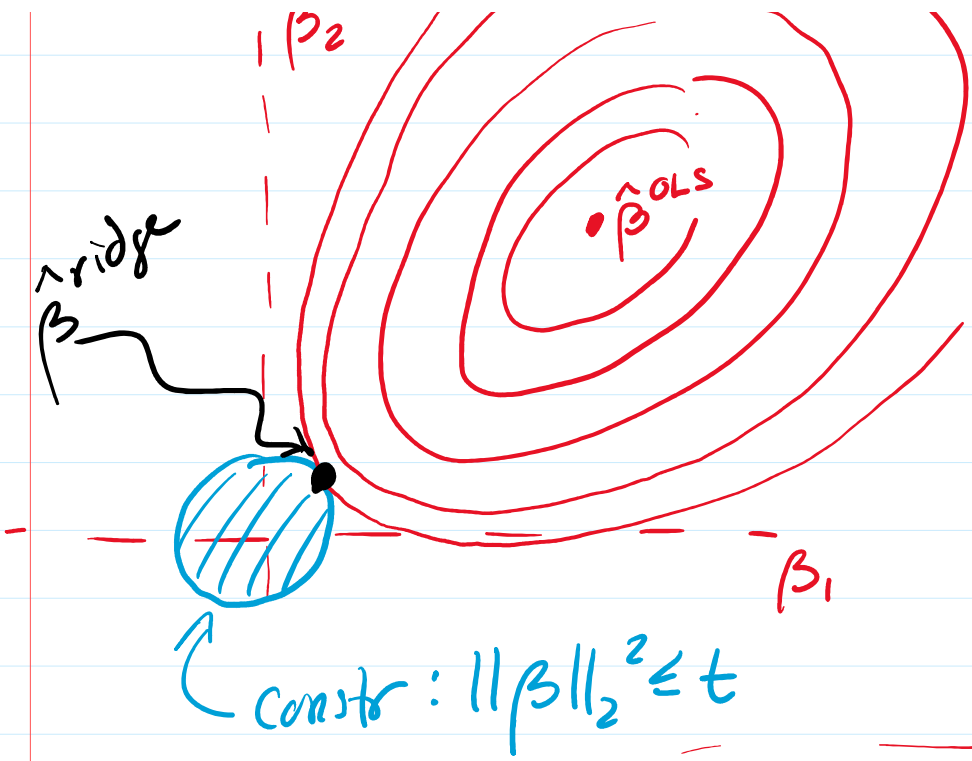$\beta$

## Why use LASSO?

$\hat{\beta}^{(LASSO)}$ will exactly zero-out some elements of the coef vector for large enough $\lambda$ — ridge doesn't do this

## Why?

## Ridge:

$|\beta_2$

$\beta_2$

$\hat{\beta}^{OLS}$

$\hat{\beta}^{ridge}$

$\beta_1$

constr : $||\beta||_2^2 \leq t$

$\beta_2$

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

$\hat{\beta}^{OL}$

$\beta_1 = 0$

$\beta_1$

constr : $||\beta||_1 \leq t$

Comparison:  assume $X$ is orthogonal

# ① Variable Selection (Hard-thresholding)

$$\hat{\beta}_j^{HT} = \begin{cases} \hat{\beta}_j^{OLS} & \text{if } |\hat{\beta}_j^{OLS}| \geq t \\ 0 & \text{else} \end{cases}$$

## ② Ridge:

$$\hat{\beta}_j^{ridge} = \frac{\hat{\beta}_j^{OLS}}{1 + \lambda} \qquad \left(\begin{array}{c} \text{proportional} \\ \text{shrinkage} \end{array}\right)$$

## ③ LASSO:

$$\hat{\beta}_j^{LASSO} = \begin{cases} \hat{\beta}^{OLS} - \lambda, & \hat{\beta}^{OLS} \geq \lambda \\ \hat{\beta}^{OLS} + \lambda, & \hat{\beta}^{OLS} \leq -\lambda \\ 0, & |\hat{\beta}^{OLS}| \leq \lambda \end{cases}$$

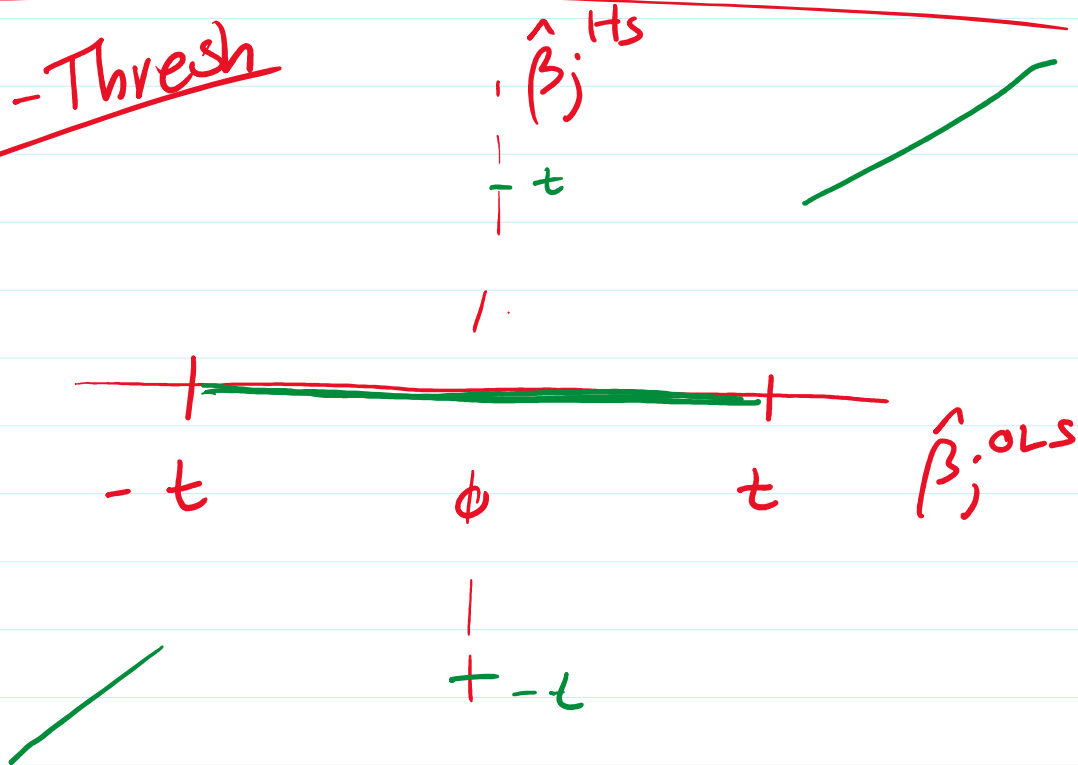$$= \text{Sign}(\hat{\beta}_j^{OLS})\left(|\hat{\beta}_j^{OLS}| - \lambda\right)_+$$
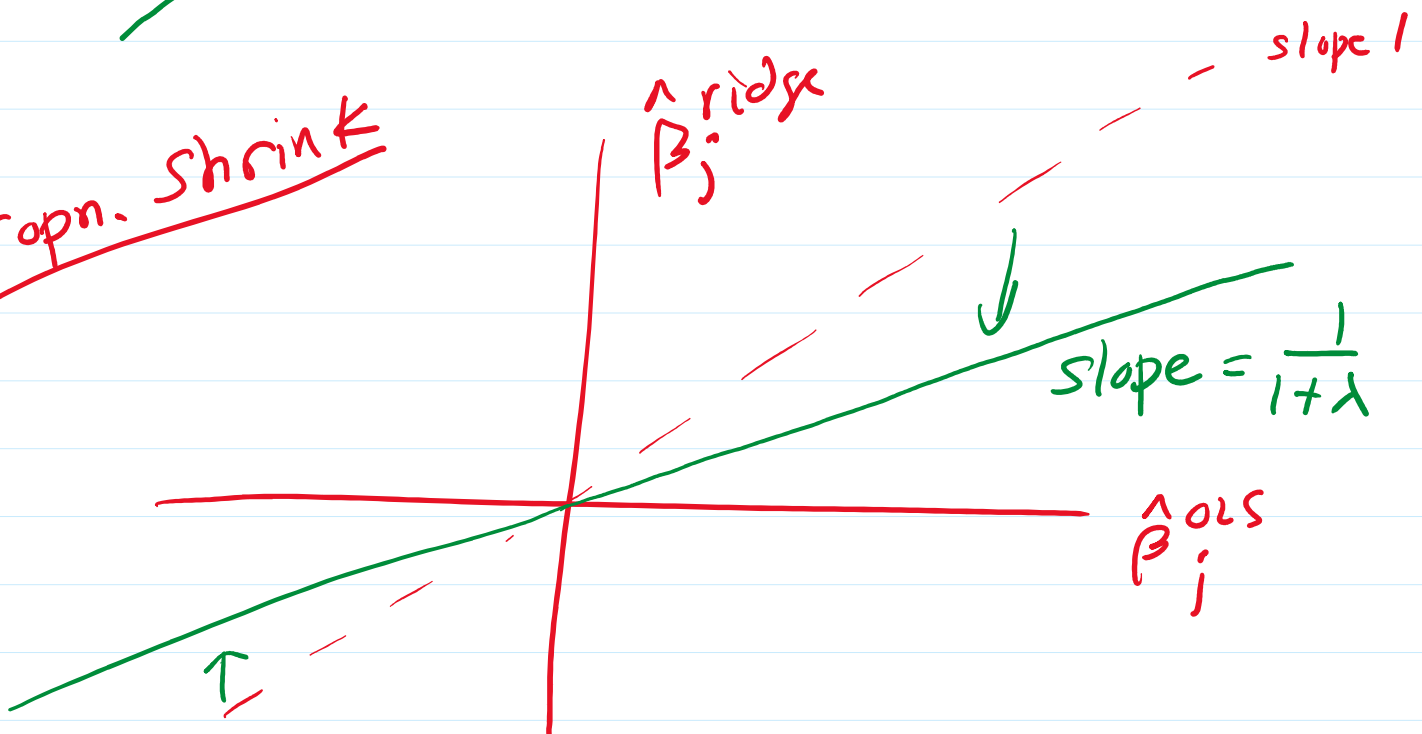
$$= \text{sign}(\beta_j)(|\beta_j| - \lambda)_+$$

$$(\cdot)_+ = \max(\cdot, 0)$$
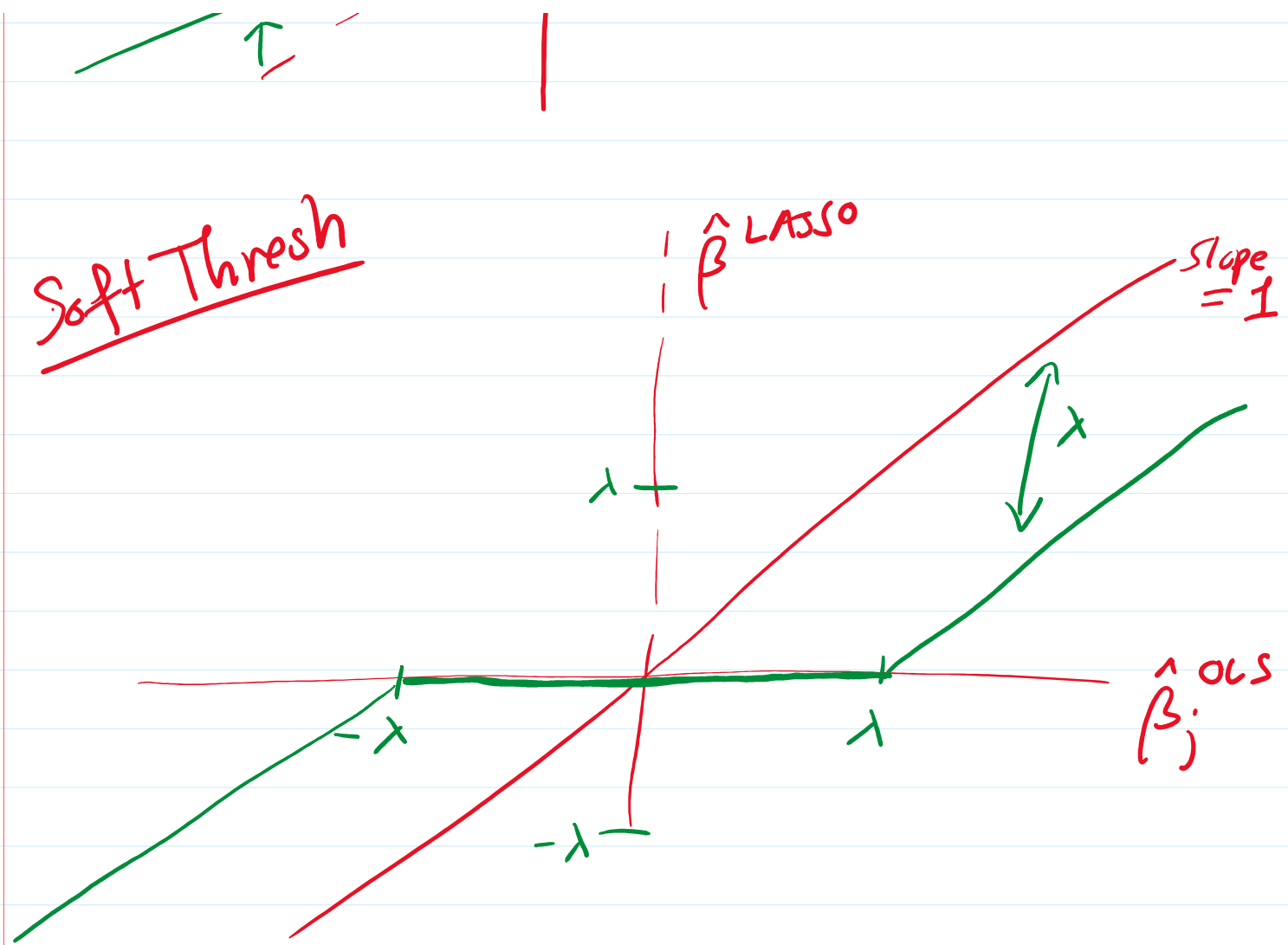
# Soft – Thresholding.

Hard – Thresh

$\hat{\beta}_j^{Hs}$

$- t$

$-t \qquad \phi \qquad t \qquad \hat{\beta}_j^{OLS}$

$+ -t$

Propn. Shrink

$\hat{\beta}_j^{ridge}$

slope $1$

slope $= \dfrac{1}{1+\lambda}$

$\hat{\beta}_j^{OLS}$

## Soft Thresh



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

slope = 1

$\lambda$

$-\lambda$

$-\lambda$

$\lambda$

$\hat{\beta}_j^{OLS}$

## Elastic Net

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

$\alpha \in [0,1]$

= tradeoff between
~~and~~ L2 penalty

= tradeoff between
L1 and L2 penalty

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

Can generally fit penalized methods

$$\hat{f} = \arg\min_{f} L(f)$$

penalize:

$$\hat{f} = \arg\min_{f} L(f) + \lambda J(f)$$

$$J(f) = \text{measure complexity of } f$$

Ex. Penalized logistic regr

$$\hat{\beta} = \arg\min NLL(\beta) + \lambda \|\beta\|_2^2$$

$$\hat{\beta} = \underset{\beta}{\arg\min} \; NLL(\beta) + \lambda \|\beta\|_2$$

$$\prime\prime \qquad \prime\prime\prime \qquad + \lambda \|\beta\|_1$$