

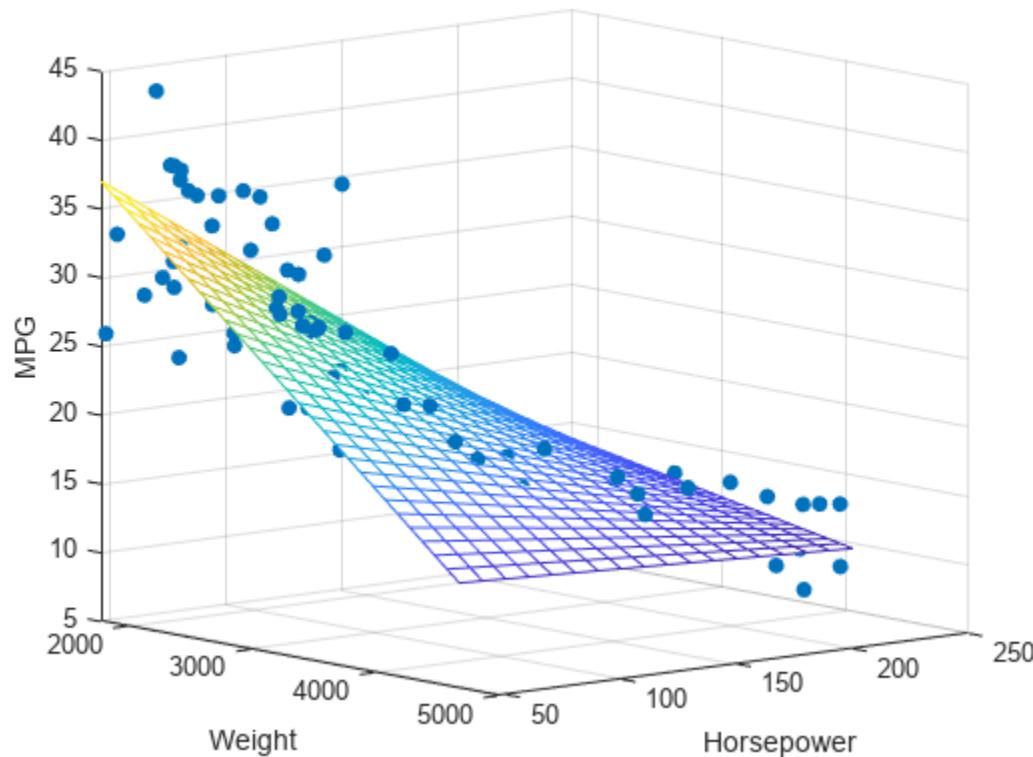
Linear Regression

Announcements

- HW 1 due today at 11:59 PM
- HW 2 posted today
- Python tutorial Friday evening – see announcement on Canvas

Regression

- Predict a continuous response variable y from a multi-dimensional feature vector \mathbf{x}



Boston Housing Dataset

x

y

- CRIM - per capita crime rate by town
- ZN - proportion of residential land zoned for lots over 25,000 sq.ft.
- INDUS - proportion of non-retail business acres per town.
- CHAS - Charles River dummy variable (1 if tract bounds river; 0 otherwise)
- NOX - nitric oxides concentration (parts per 10 million)
- RM - average number of rooms per dwelling
- AGE - proportion of owner-occupied units built prior to 1940
- DIS - weighted distances to five Boston employment centres
- RAD - index of accessibility to radial highways
- TAX - full-value property-tax rate per \$10,000
- PTRATIO - pupil-teacher ratio by town
- B - $1000(Bk - 0.63)^2$ where Bk is the proportion of blacks by town
- LSTAT - % lower status of the population
- MEDV - Median value of owner-occupied homes in \$1000's

MSE

- Probabilistic setting: We have jointly distributed variables (\mathbf{X}, Y) where

$$\mathbf{X} \in \mathbb{R}^d, \quad Y \in \mathbb{R}$$

and the goal is to predict Y from \mathbf{X} using a function

$$f : \mathbb{R}^d \rightarrow \mathbb{R}.$$

- Performance measure:

$$\Pr(f(\mathbf{x}) \neq y)$$

mean squared error

$$R(f) = \mathbb{E}_{x,y} [(y - f(x))^2]$$

Linear Regression

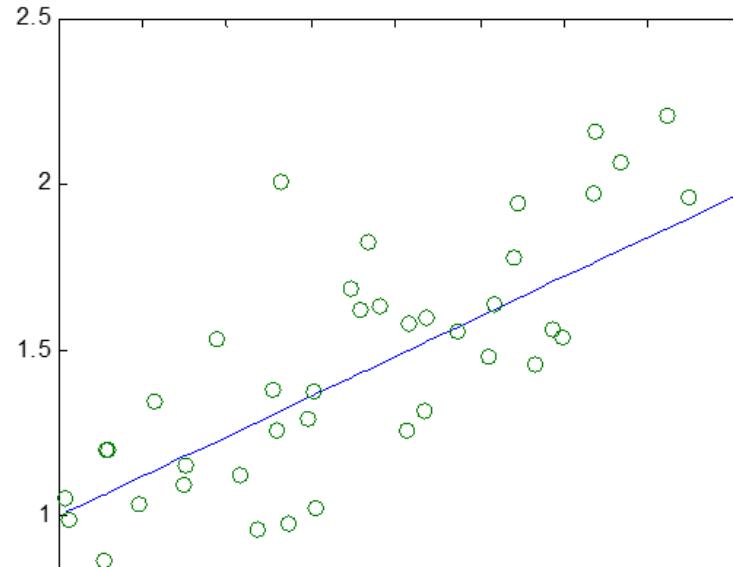
- In practice we don't have access to the joint distribution, but we do have training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$.
- Choose f to minimize the *empirical MSE*

$$\min_f \quad \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 \quad \stackrel{\text{LLN}}{\approx} \mathbb{E}_{x,y} [(y - f(x))^2]$$

- To formulate a tractable optimization problem, we need to restrict f to belong to a *regression model*, i.e., a class of candidates for f .
- Today we'll focus on the *linear model*

$$f(x) = \mathbf{w}^T \mathbf{x} + b$$

where $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$.



Least Squares Linear Regression

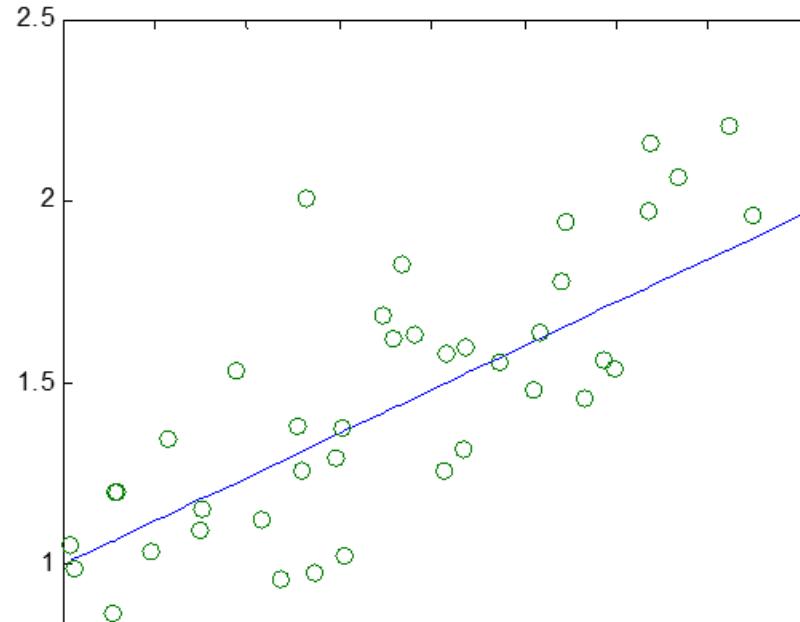
- Least squares linear regression solves

$$\min_{w, b} \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i - b)^2$$

The method is also known as
ordinary least squares.

- For greater generality, we can add a regularization term

$$\min_{w, b} \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i - b)^2 + \lambda \|w\|^2$$



This method is known as *ridge regression*, and the term $\lambda \|w\|^2$ is called the *ridge penalty*. $\lambda \geq 0$ is the *regularization parameter*.

OLS Derivation

$$\theta = \begin{bmatrix} b \\ \omega \end{bmatrix} \in \mathbb{R}^{d+1}$$

Assume $\lambda = 0$ (ordinary least squares). Then

$$= \frac{1}{n} \|y - X\theta\|^2$$

$$= \frac{1}{n} (y - X\theta)^T (y - X\theta)$$

$$= \frac{1}{n} (y^T y - 2y^T X\theta + \theta^T X^T X\theta)$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \cdots & x_{nd} \end{bmatrix}$$

$$x_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}$$

$$z^T X^T X z = \|Xz\|^2 \geq 0$$

$$(AB)^T = B^T A^T \quad \Rightarrow \quad (Xz)^T (Xz)$$

Poll

Warning: Overloaded notation ahead.

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ \vdots & & & \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix}$$

Assume \mathbf{X} has full column rank. Then the solution of

$$\min_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$$

is

- (A) $\hat{\boldsymbol{\theta}} = \mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{y}$
- (B) $\hat{\boldsymbol{\theta}} = \mathbf{X}^{-1}\mathbf{y}$
- (C) $\hat{\boldsymbol{\theta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$
- (D) $\hat{\boldsymbol{\theta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$ ✓

OLS Derivation

Fact: X has full col. rank $\Leftrightarrow X^T X$ is invertible

$$\min_{\theta} \quad y^T y - 2y^T X \theta + \theta^T X^T X \theta$$

$$\nabla_{\theta} (\downarrow) = -2X^T y + 2X^T X \theta = 0$$

$$\Rightarrow \hat{\theta} = \underbrace{(X^T X)^{-1} X^T y}_{\text{pseudo inverse of } X} \quad \text{is a critical pt.}$$

Hessian $= 2X^T X$ \Rightarrow strictly convex $\Rightarrow \hat{\theta}$ is the $\underset{\lambda}{\text{global min}}$
is PD unique

Poll

What happens in ordinary least squares when $d \geq n$? Select the best answer.

(A) \mathbf{X} does not have full column rank ✓

$$\begin{matrix} \uparrow\downarrow \\ d+1 > n \end{matrix}$$

(B) $\mathbf{X}^T \mathbf{X}$ is not invertible ✓

(C) The OLS problem has infinitely many minimizers ✓

(D) All of the above ✓

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix}$$

Ridge Regression Derivation

- Optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2 + \lambda \|\mathbf{w}\|^2$$

- First step: eliminate b

$$\min_{\mathbf{w}} \left[\min_b \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2 + \lambda \|\mathbf{w}\|^2 \right]$$

$$\frac{\partial}{\partial b} \left(\cdot \right) = -\frac{2}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - b) = 0$$

$$\Rightarrow \hat{b} = \hat{b}(\mathbf{w}) = \bar{y} - \mathbf{w}^T \bar{\mathbf{x}}, \quad \bar{y} = \frac{1}{n} \sum y_i \in \mathbb{R}$$

$$\bar{\mathbf{x}} = \frac{1}{n} \sum \mathbf{x}_i \in \mathbb{R}^d$$

Ridge Regression Derivation

$$\min_w \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i - (\bar{y} - w^T \bar{x}))^2 + \lambda \|w\|^2$$

$$\min_w \frac{1}{n} \sum_i ((y_i - \bar{y}) - w^T (x_i - \bar{x}))^2 + \lambda \|w\|^2$$

$$\min_w \frac{1}{n} \sum_i (\tilde{y}_i - w^T \tilde{x}_i)^2 + \lambda \|w\|^2$$

" "

$$\frac{1}{n} \|\tilde{y} - \tilde{X}w\|^2 + \lambda \|w\|^2$$

$$\begin{aligned}\tilde{y}_i &= y_i - \bar{y} \\ \tilde{x}_i &= x_i - \bar{x}\end{aligned}$$

$$\tilde{y} = \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_n \end{bmatrix}, \tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix}$$

Ridge Regression Derivation

- The regularized least-squares (i.e., ridge regression) objective function can be written (after eliminating b)

$$J(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{r}^T \mathbf{w} + c,$$

where $\mathbf{A} = 2(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + n\lambda \mathbf{I})$, $\mathbf{r} = -2\tilde{\mathbf{X}}^T \tilde{\mathbf{y}}$, and $c = \tilde{\mathbf{y}}^T \tilde{\mathbf{y}}$ where

$$\tilde{\mathbf{y}} = \begin{bmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_n \end{bmatrix} \quad \tilde{\mathbf{X}} = \begin{bmatrix} \tilde{x}_{11} & \cdots & \tilde{x}_{1d} \\ \vdots & \ddots & \vdots \\ \tilde{x}_{n1} & \cdots & \tilde{x}_{nd} \end{bmatrix},$$

$\tilde{y}_i = y_i - \bar{y}$ and $\tilde{x}_i = \mathbf{x}_i - \bar{\mathbf{x}}$.

If $\lambda > 0$, then \mathbf{A} is PD:

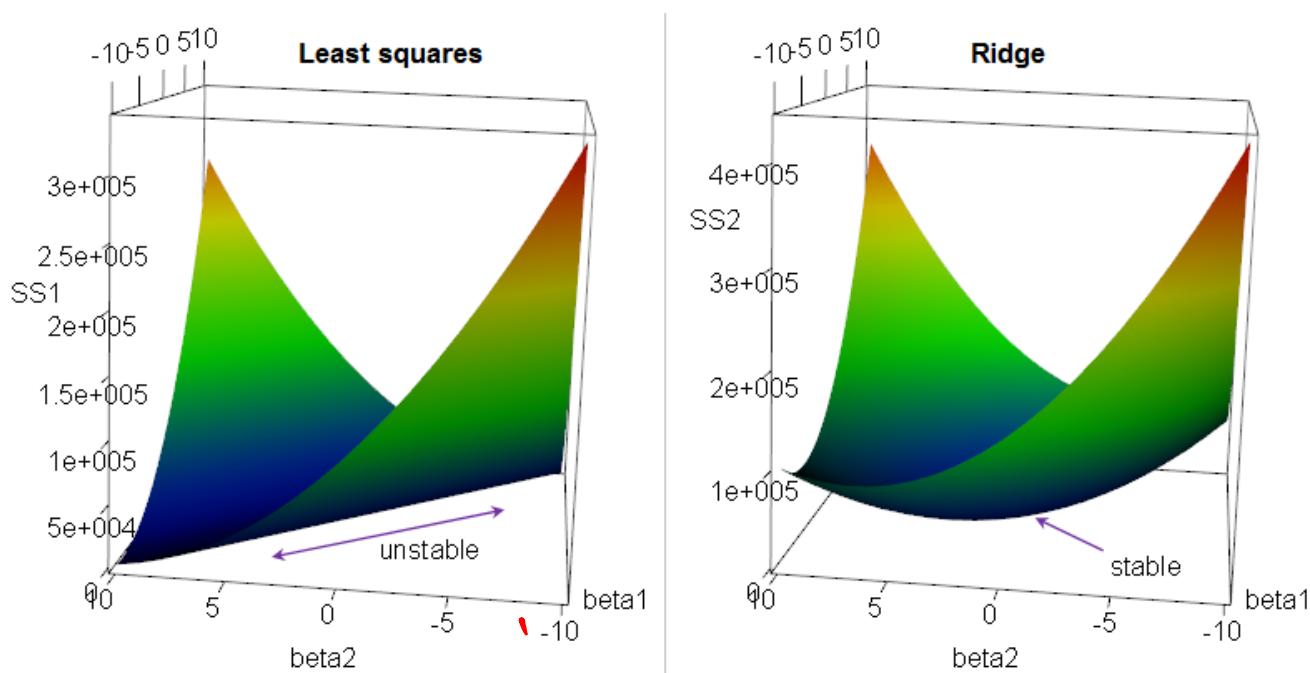
- Solution:

$$\hat{\mathbf{w}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \lambda n \mathbf{I})^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}$$

- Then $\hat{b} = \bar{y} - \hat{\mathbf{w}}^T \bar{\mathbf{x}}$.
$$\begin{aligned} \text{If } z \neq 0, z^T \mathbf{A} z &= z^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + n\lambda \mathbf{I}) z \\ &= \|\tilde{\mathbf{X}} z\|^2 + n\lambda \|z\|^2 > 0 \end{aligned}$$

OLS vs RR

- Ridge penalty: makes objective strictly convex, solution unique



The Lasso

- ℓ_p norm

$$\|w\|_p = \left(\sum_{j=1}^d |w_j|^p \right)^{1/p}$$

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

$p=2$: $\|w\|_2$ = Euclidean norm

$$p=1: \|w\|_1 = \sum_{j=1}^d |w_j|$$

- The lasso is an alternative to ridge regression:

$$\frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i - b)^2 + \lambda \|w\|_1$$

- Motivation?

Leads to sparse w (many entries are 0)
 \Rightarrow feature selection

The Lasso

- The contours of the sum of squared errors are elliptical:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i - b)^2 &= \| \tilde{y} - \tilde{X}w \|_2^2 \\ &= (\tilde{y} - \tilde{X}w)^\top (\tilde{y} - \tilde{X}w) = \tilde{y}^\top \tilde{y} - 2\tilde{y}^\top \tilde{X}w + w^\top \tilde{X}^\top \tilde{X}w \\ &= \frac{1}{2} w^\top A w + r^\top w + c \\ &= \frac{1}{2} (w - w^*)^\top A (w - w^*) + c' \end{aligned}$$

$w^* = -A^{-1}r$ $c' = c - r^\top A^{-1}r$

$\xrightarrow{\text{elliptical contours}}$

The Lasso

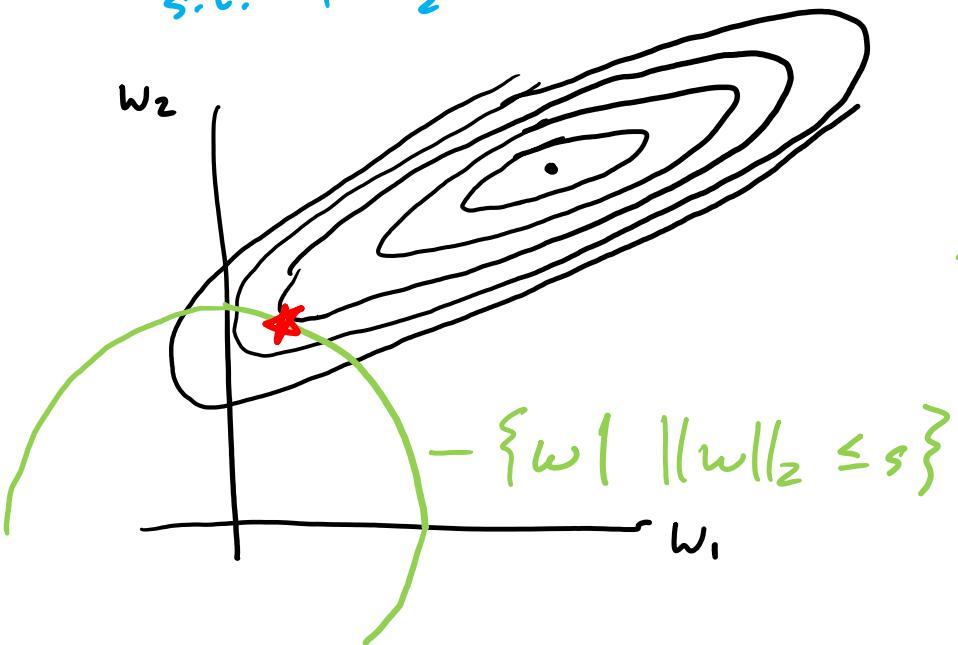
- Why does the lasso yield a sparse solution?

$$\min_w \|\tilde{y} - \tilde{X}w\|_2^2 + n\lambda \|w\|_2^2$$

\Downarrow

$$\min_w \|\tilde{y} - \tilde{X}w\|_2^2$$

s.t. $\|w\|_2 \leq s$

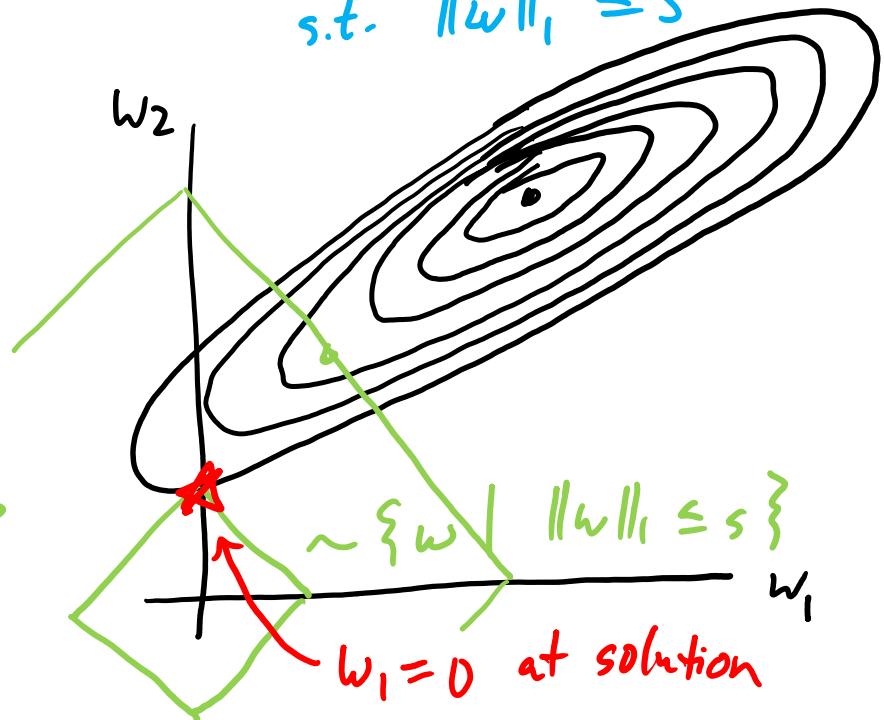


$$\min_w \|\tilde{y} - \tilde{X}w\|_2^2 + n\lambda \|w\|_1$$

\Downarrow

$$\min_w \|\tilde{y} - \tilde{X}w\|_2^2$$

s.t. $\|w\|_1 \leq s$



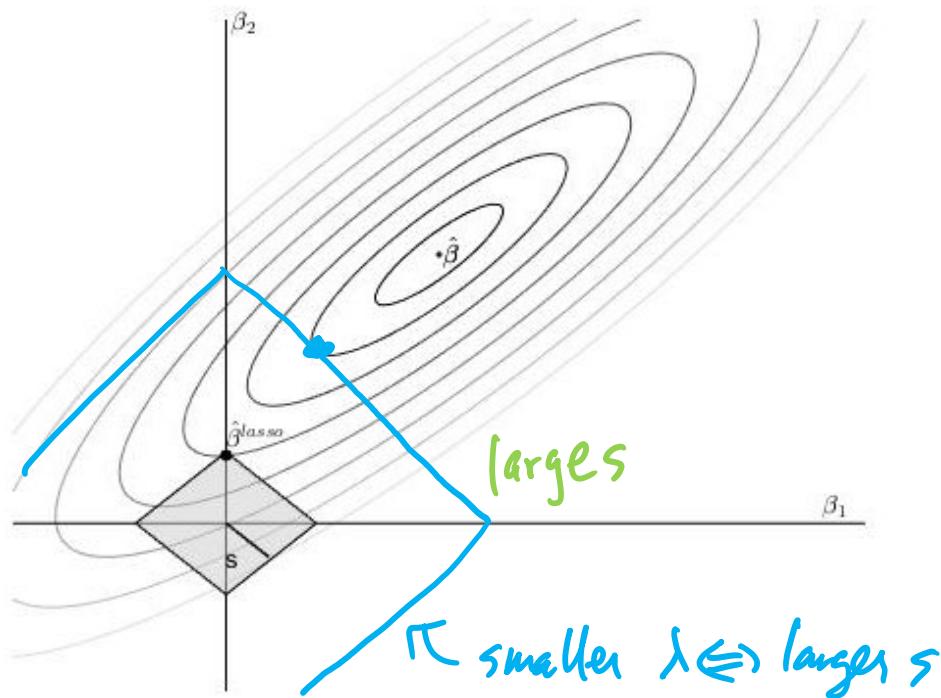
$$\text{P} \underset{\mathbf{w}, b}{\text{OII}} \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i - b)^2 + \lambda \|\mathbf{w}\|_1$$

- As λ increases, the sparsity of the minimizer \mathbf{w}

- (A) Increases (fewer nonzero entries)
 (B) Decreases (more nonzero entries)

1) Larger λ makes $\|\mathbf{w}\|_1$ smaller which leads to more 0 entries

2) Smaller λ places more emphasis on the SSE and hence a larger s .



Concluding Thoughts

- Least squares: exact solution possible, but may be infeasible for large d
- Ridge penalty: makes objective strictly convex, solution unique
- Lecture notes: robust regression
- Linear regression is the basis for several nonlinear regression methods
- Two lectures from now: computation and optimization for RR/lasso

