

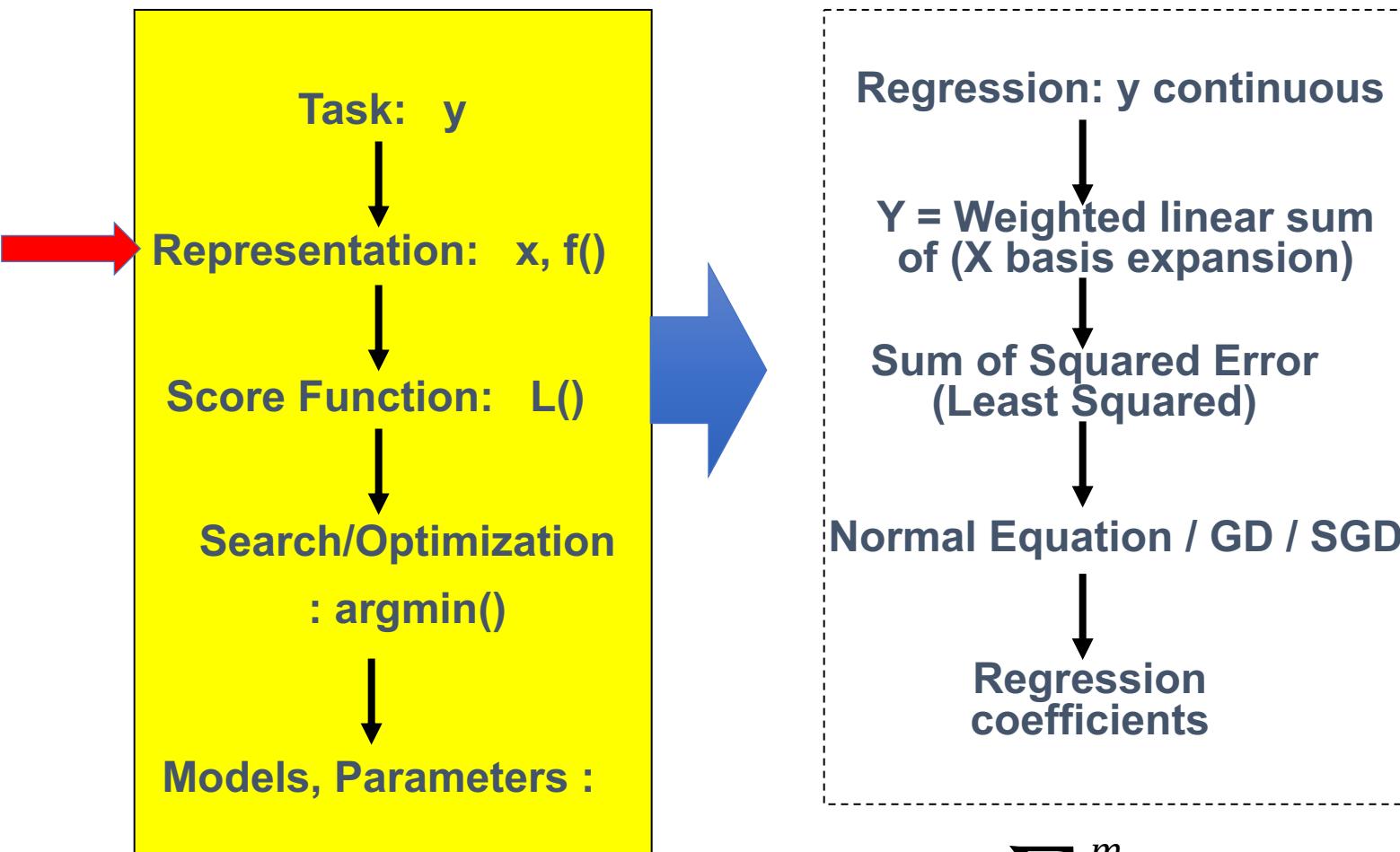
# UVA CS 6316: Machine Learning

## Lecture 6: Linear Regression Model with Regularizations

Dr. Yanjun Qi

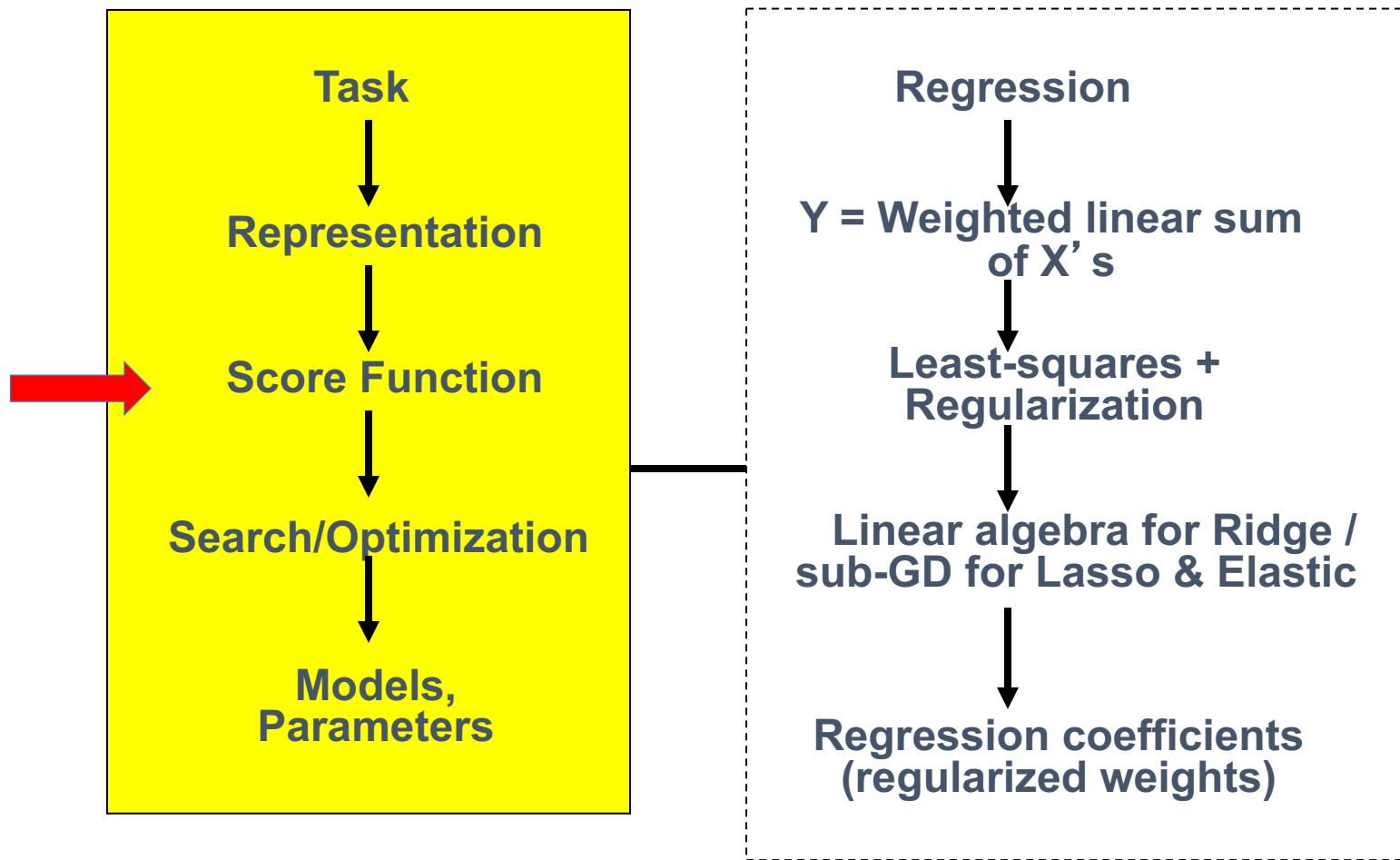
University of Virginia  
Department of Computer Science

# Last: Multivariate Linear Regression with basis Expansion



$$\hat{y} = \theta_0 + \sum_{j=1}^m \theta_j \varphi_j(x) = \varphi(x)^T \theta$$

# Today: Regularized multivariate linear regression



$$\min J(\beta) = \sum_{i=1}^n \left( Y - \hat{Y} \right)^2 + \lambda \left( \sum_{j=1}^p \beta_j^q \right)^{1/q}$$

We aim to make the learned model

- 1. Generalize Well
- 2. Computational Scalable and Efficient
- 3. Robust / Trustworthy / **Interpretable**
  - Especially for some domains, this is about trust!

# Today

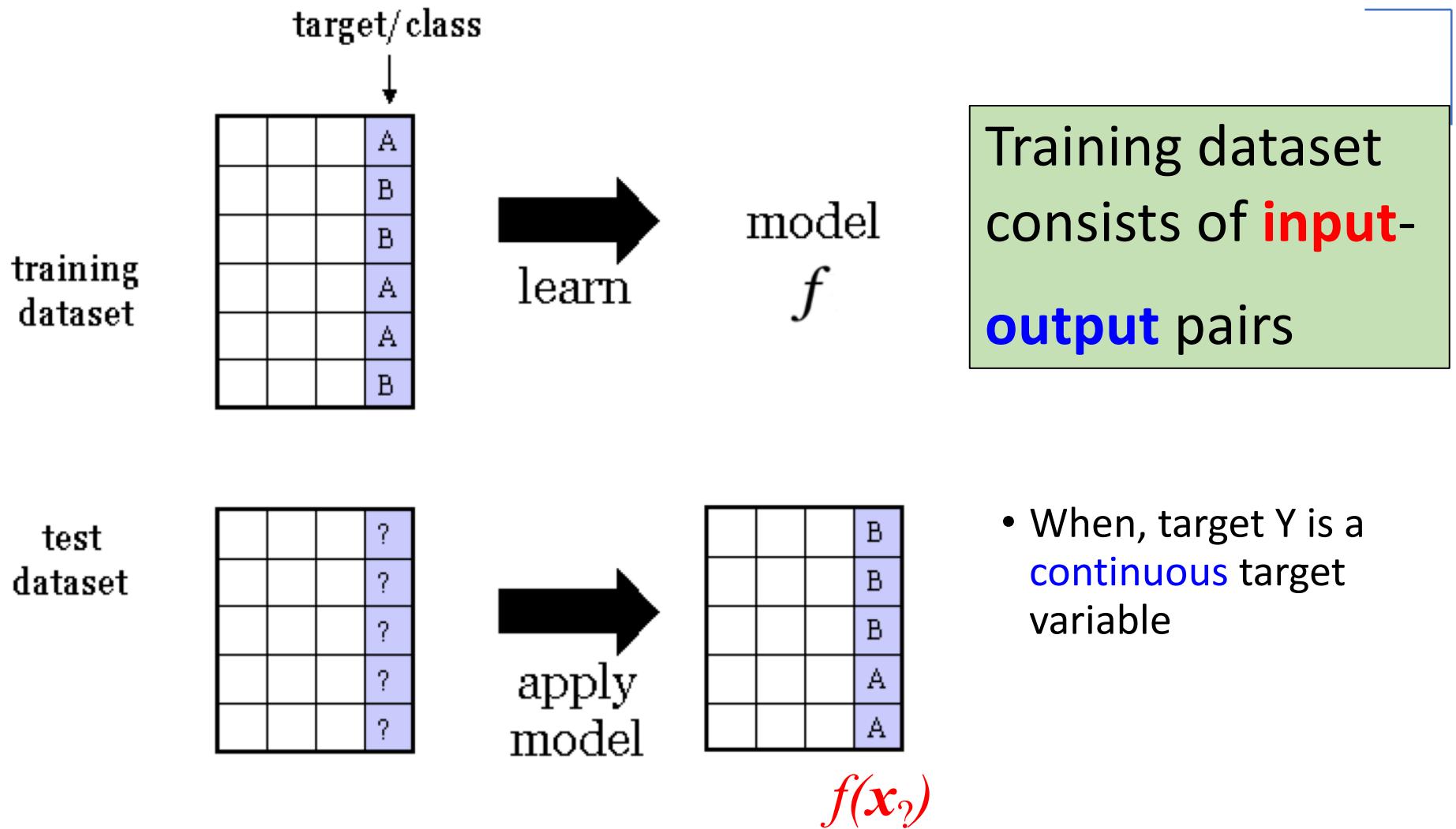


Linear Regression Model with Regularizations

→ Review: (Ordinary) Least squares: squared loss (Normal Equation)

- ✓ Ridge regression: squared loss with L2 regularization
- ✓ Lasso regression: squared loss with L1 regularization
- ✓ Elastic regression: squared loss with L1 AND L2 regularization
- ✓ How to Choose Regularization Parameter

# SUPERVISED Regression



# Review: Normal equation for LR

- Write the cost function in matrix form:

$$\begin{aligned} J(\beta) &= \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \beta - y_i)^2 \\ &= \frac{1}{2} (\mathbf{X}\beta - \bar{\mathbf{y}})^T (\mathbf{X}\beta - \bar{\mathbf{y}}) \\ &= \frac{1}{2} (\beta^T \mathbf{X}^T \mathbf{X}\beta - \beta^T \mathbf{X}^T \bar{\mathbf{y}} - \bar{\mathbf{y}}^T \mathbf{X}\beta + \bar{\mathbf{y}}^T \bar{\mathbf{y}}) \end{aligned}$$

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1^T & \cdots \\ \cdots & \mathbf{x}_2^T & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \mathbf{x}_n^T & \cdots \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

To minimize  $J(\theta)$ , take derivative and set to zero:

$$\Rightarrow \boxed{\mathbf{X}^T \mathbf{X}\beta = \mathbf{X}^T \bar{\mathbf{y}}}$$

The normal equations

$$\beta^* = \boxed{\downarrow} \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \bar{\mathbf{y}}$$

Assume  
that  $\mathbf{X}^T \mathbf{X}$  is  
invertible

# Comments on the normal equation

What if  $X$  has less than full column rank?

→ Not Invertible

$$\text{rank}(\Sigma_{n \times p}) = \min(n, p)$$

when  $p > n$

$$\text{rank}(\Sigma) < p$$

~~$$(\Sigma^T \Sigma)^{-1}$$~~

$$\text{rank} \left( \begin{matrix} \Sigma^T & \Sigma \\ p \times n & n \times p \end{matrix} \right) \leq \min(r(\Sigma^T), r(\Sigma))$$

$< p$

For any matrix  $A \in \mathbb{R}^{m \times n}$ , it turns out that the column rank of  $A$  is equal to the row rank of  $A$  (though we will not prove this), and so both quantities are referred to collectively as the **rank** of  $A$ , denoted as  $\text{rank}(A)$ . The following are some basic properties of the rank:

- For  $A \in \mathbb{R}^{m \times n}$ ,  $\boxed{\text{rank}(A) \leq \min(m, n)}$ . If  $\text{rank}(A) = \min(m, n)$ , then  $A$  is said to be **full rank**. (2)
- For  $A \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(A) = \text{rank}(A^T)$ .
- For  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times p}$ ,  $\boxed{\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))}$ . (1)
- For  $A, B \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$ .

Page 11 of  
Handout L2

$$\underbrace{X^T X}_{p \times p}$$

$$\text{rank}(X^T X) \leq \text{rank}(X) \leq \min(n, p)$$

When  $n < p$

$$\text{rank}(X^T X) < p$$

$\Downarrow$  singular / not invertible

# Today



Linear Regression Model with Regularizations

- ✓ Review: (Ordinary) Least squares: squared loss (Normal Equation)
- Ridge regression: squared loss with L2 regularization
- ✓ Lasso regression: squared loss with L1 regularization
- ✓ Elastic regression: squared loss with L1 AND L2 regularization
- ✓ How to Choose Regularization Parameter

# Review: Vector norms

A norm of a vector  $\|x\|$  is informally a measure of the “length” of the vector.

$$\|x\|_q = \left( \sum_{i=1}^n |x_i|^q \right)^{1/q} \quad q=1, 2, \dots$$

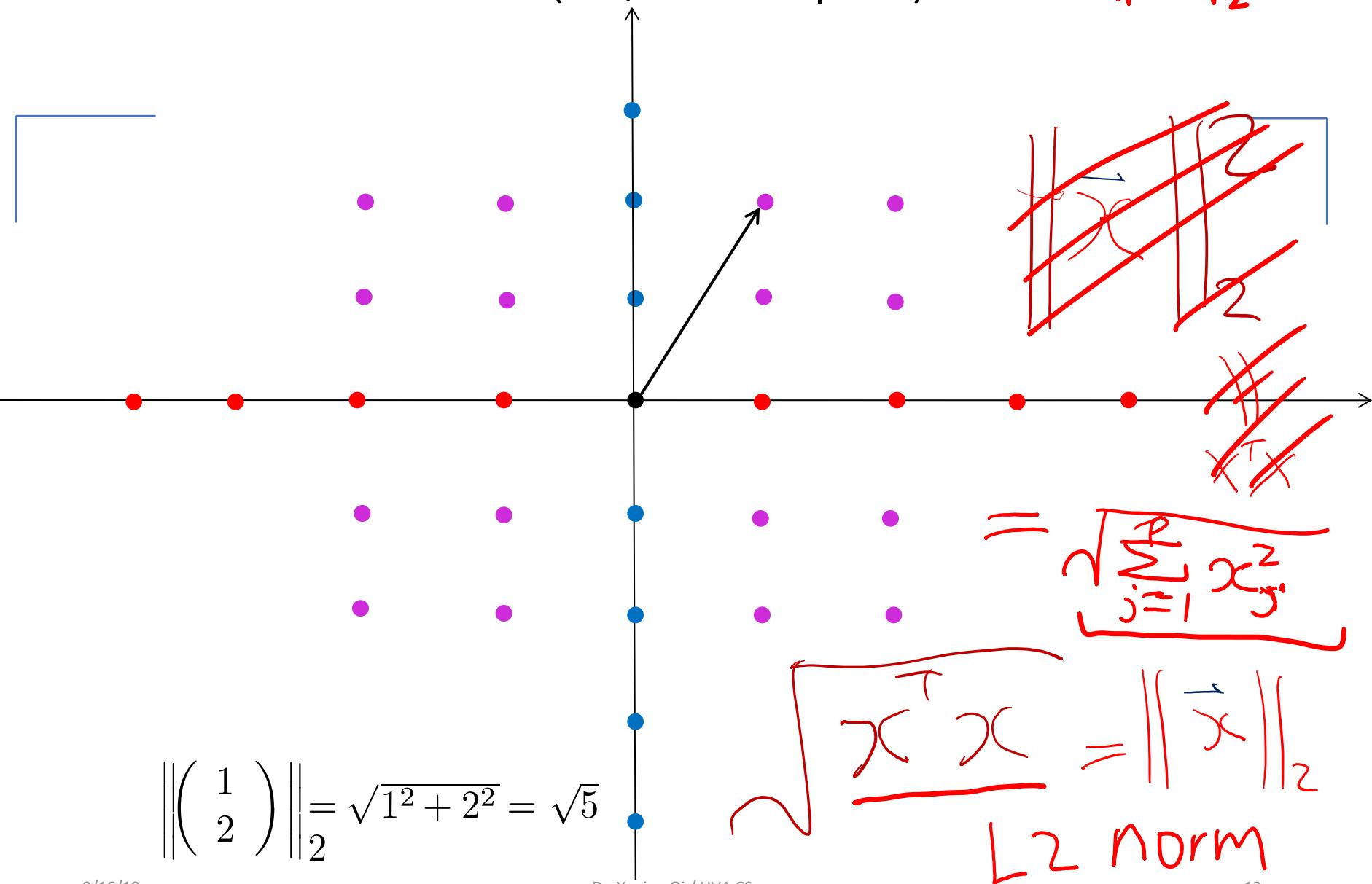
- Common norms:  $L_1$ ,  $L_2$  (Euclidean)

$$\|x\|_1 = \sum_{i=1}^n |x_i| \quad \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

- $L_{\infty}$

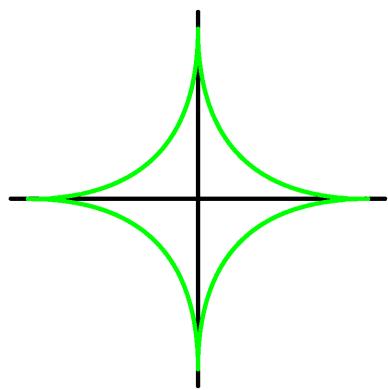
$$\|x\|_\infty = \max_i |x_i|$$

Review: Vector Norm (L2, when p=2)  $\vec{x}^T \vec{x} = \|\vec{x}\|_2^2$

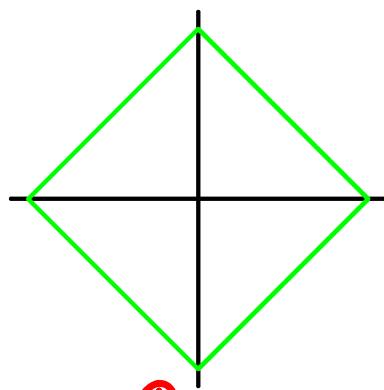


# $q$ Norms

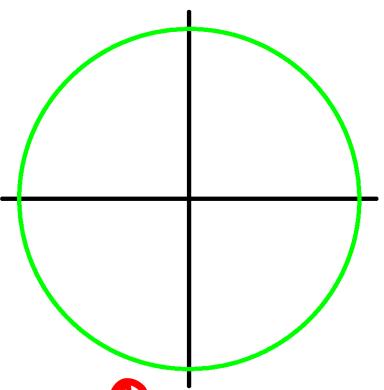
$$\|x\|_q = \left( \sum_{i=1}^n |x_i|^q \right)^{1/q}$$



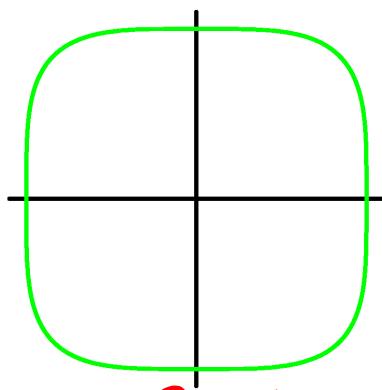
$q=0.5$



$q=1$   
diamond  
contour



$q=2$   
circle  
contour



$q=4$

# Ridge Regression / L2 Regularization

$$\hat{\beta}_{OLS} = \beta^* = (X^T X)^{-1} X^T \bar{y}$$



- If not **invertible**, a classical solution is to add a small positive element to diagonal

$$\lambda > 0$$

$$\beta^* = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

By convention, the bias/intercept term is typically not regularized.

Here we assume data has been centered ... therefore no bias term

# Extra: Positive Definite Matrix

- A symmetric matrix  $A \in \mathbb{S}^n$  is ***positive definite*** (PD) if for all non-zero vectors  $x \in \mathbb{R}^n$ ,  $x^T Ax > 0$ . This is usually denoted  $A \succ 0$  (or just  $A > 0$ ), and often times the set of all positive definite matrices is denoted  $\mathbb{S}_{++}^n$ .
- A symmetric matrix  $A \in \mathbb{S}^n$  is ***positive semidefinite*** (PSD) if for all vectors  $x^T Ax \geq 0$ . This is written  $A \succeq 0$  (or just  $A \geq 0$ ), and the set of all positive semidefinite matrices is often denoted  $\mathbb{S}_+^n$ .

One important property of positive definite matrices is that

- They are always full rank, and hence, invertible.
- Extra: See Proof at Page 17-18 of Linear-Algebra Handout

positive definite (PD)

$$\forall a \neq 0 \quad \underbrace{a^T (X^T \Sigma + \lambda I) a}_{\text{}} > 0$$

$$= a^T X^T \Sigma a + \lambda a^T a$$

$$= \|\Sigma a\|_2^2 + \lambda \|a\|_2^2 > 0$$

$$\beta^* = \left( X^T X + \lambda I \right)^{-1} X^T \bar{y}$$

Extra: Positive Definite Matrix

$$\forall \vec{a} \neq 0, \quad \vec{a}^T A \vec{a} \geq 0 \Rightarrow A \succeq 0$$

$$\textcircled{1} \quad \vec{a}^T X^T X \vec{a} = \underbrace{(X\vec{a})^T (X\vec{a})}_{\substack{N \times P \\ P \times 1}} = \|X\vec{a}\|_2^2 \geq 0$$

[for any non-zero vector  $\vec{a} \in \mathbb{R}^P$ ]



$$X^T X \text{ PSD}$$

$$\textcircled{2} \quad \vec{a}^T \underbrace{(X^T X + \lambda I)}_{\substack{\text{PD} \rightarrow \text{invertible}}} \vec{a} = \vec{a}^T X^T X \vec{a} + \lambda \vec{a}^T I \vec{a} = \|X\vec{a}\|_2^2 + \lambda \|\vec{a}\|_2^2 \geq 0$$

$\lambda > 0, \vec{a} \neq 0$

# Ridge Regression / Squared Loss+L2

$$\beta^* = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

- As the solution from

HW2

$$\hat{\beta}^{ridge} = \operatorname{argmin}_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

to minimize, take derivative and set to zero

$$\sum \beta \rightarrow \hat{y}$$

$n \times p$   $p \times 1$        $n \times 1$

# Ridge Regression / Squared Loss+L2

$$\beta^* = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

- As the solution from

$$\hat{\beta}^{ridge} = \operatorname{argmin}_{\beta} \sum_{j=1}^n (y_j - \beta^T \tilde{x}_j)^2$$

HW2

$$\hat{\beta}^{ridge} = \operatorname{argmin}_{\beta} (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

to minimize, take derivative and set to zero

By convention, the bias/intercept term is typically not regularized.  
Here we assume data has been centered ... therefore no bias term

$$\sum \beta \rightarrow \hat{y}$$

$n \times p$   $p \times 1$        $n \times 1$

# Ridge Regression / Squared Loss+L2

$$\beta^* = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

- As the solution from

$$\sum_{j=1}^n (y_j - \beta^T \tilde{x}_j)^2$$

HW2

$$\hat{\beta}^{ridge} = \operatorname{argmin}(y - X\beta)^T(y - X\beta) + \lambda \beta^T \beta$$

to minimize, take derivative and set to zero

- Equivalently  $\hat{\beta}^{ridge} = \operatorname{argmin}(y - X\beta)^T(y - X\beta)$

subject to  $\sum_{j=\{1..p\}} \beta_j^2 \leq s^2$

circle  
with radius S

By convention, the bias/intercept term is typically not regularized.

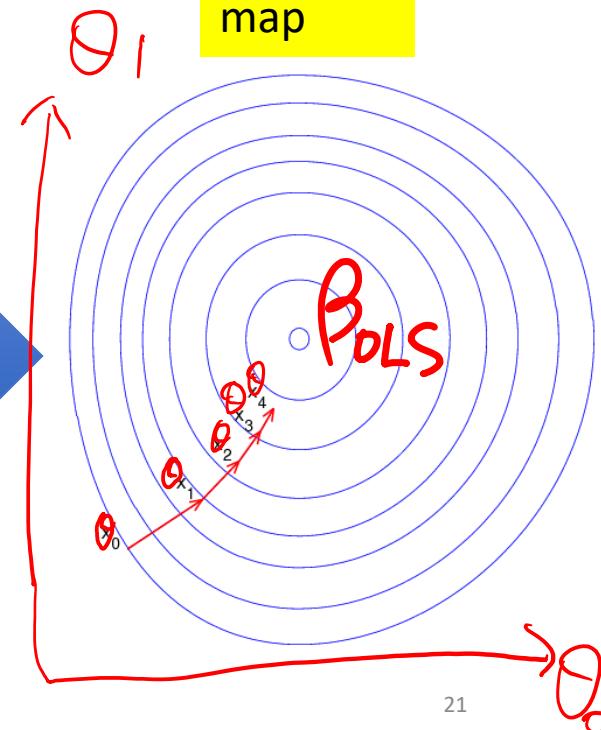
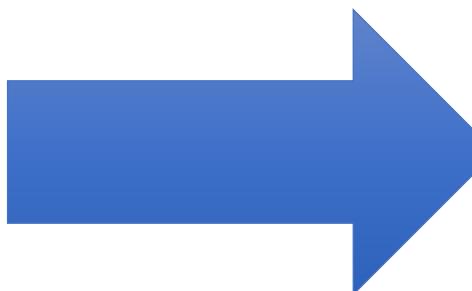
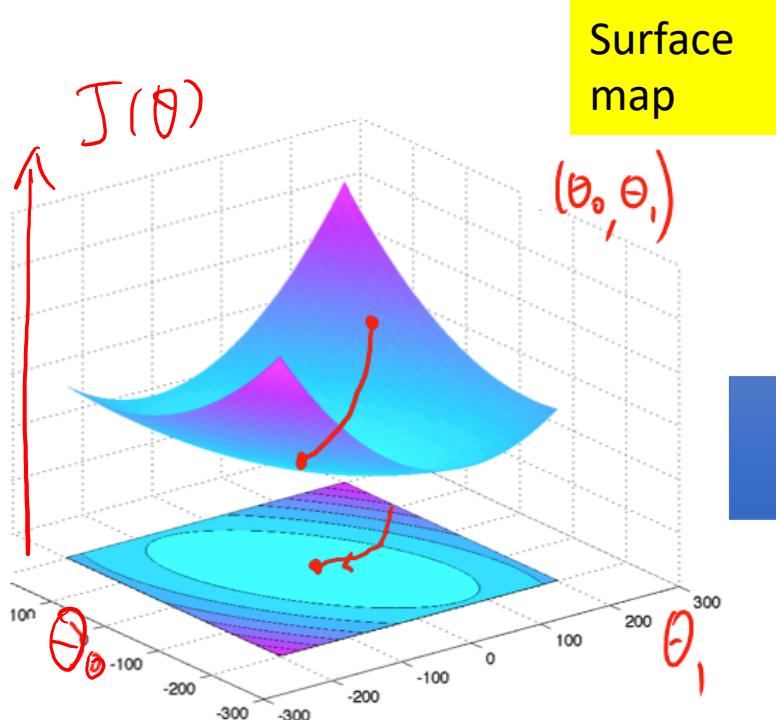
Here we assume data has been centered ... therefore no bias term



# Review

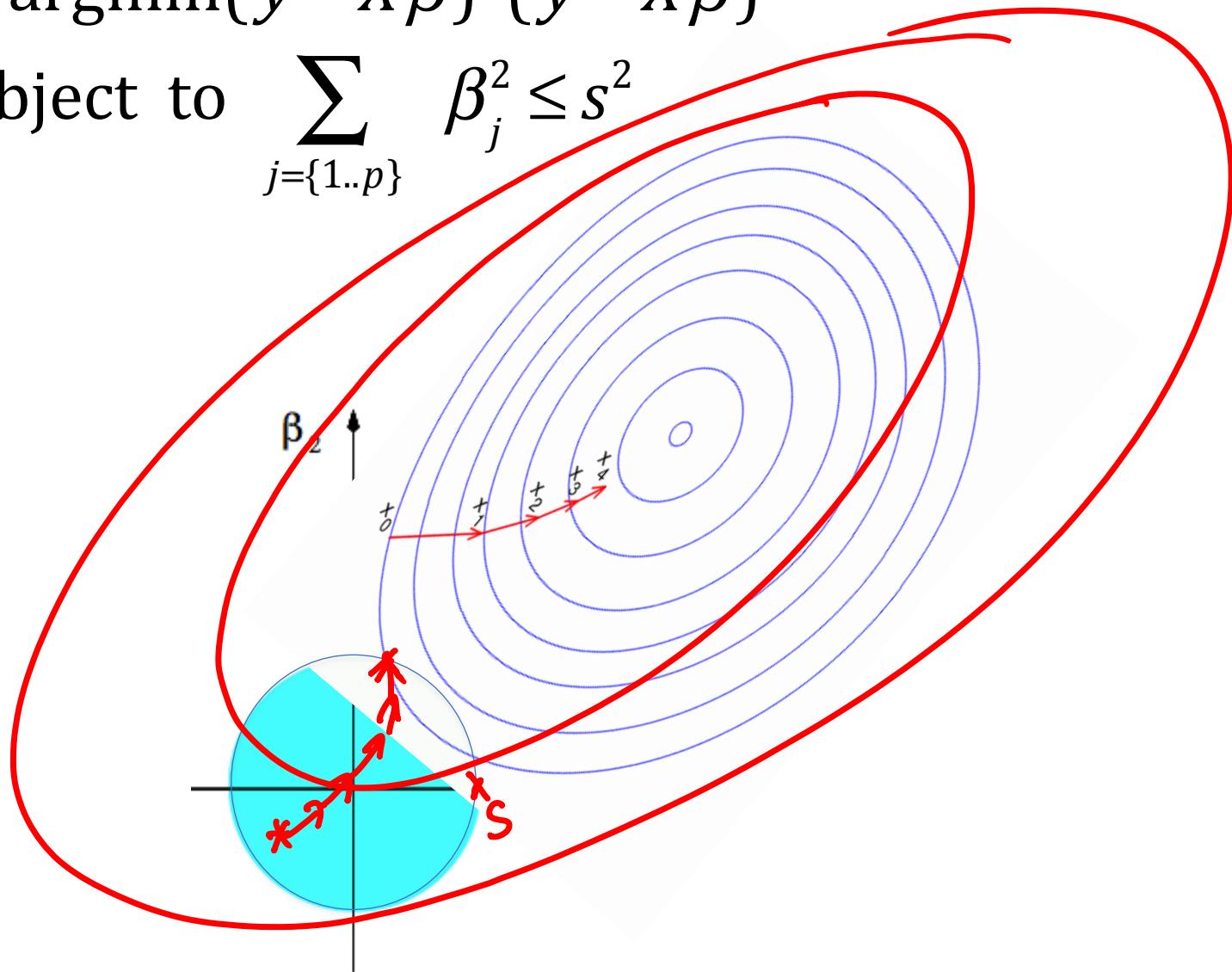


Contour map



$$\hat{\beta}^{ridge} = \operatorname{argmin}(y - X\beta)^T(y - X\beta)$$

subject to  $\sum_{j=\{1..p\}} \beta_j^2 \leq s^2$

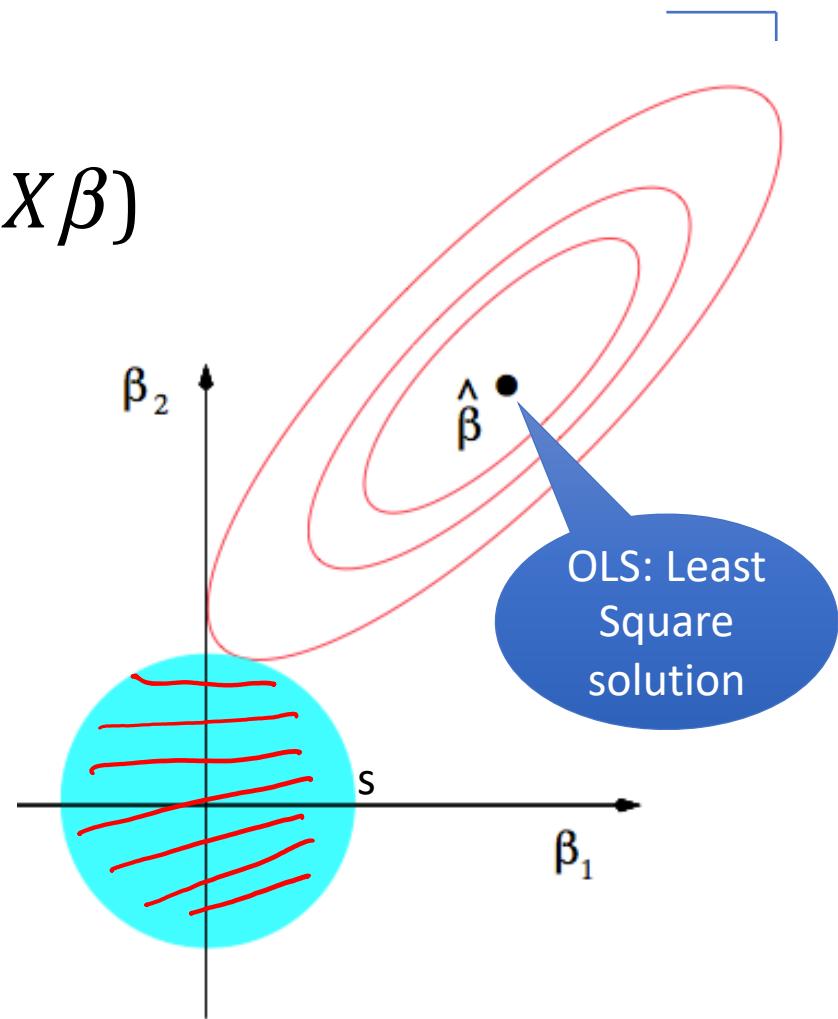


# Objective Function's Contour lines from Ridge Regression

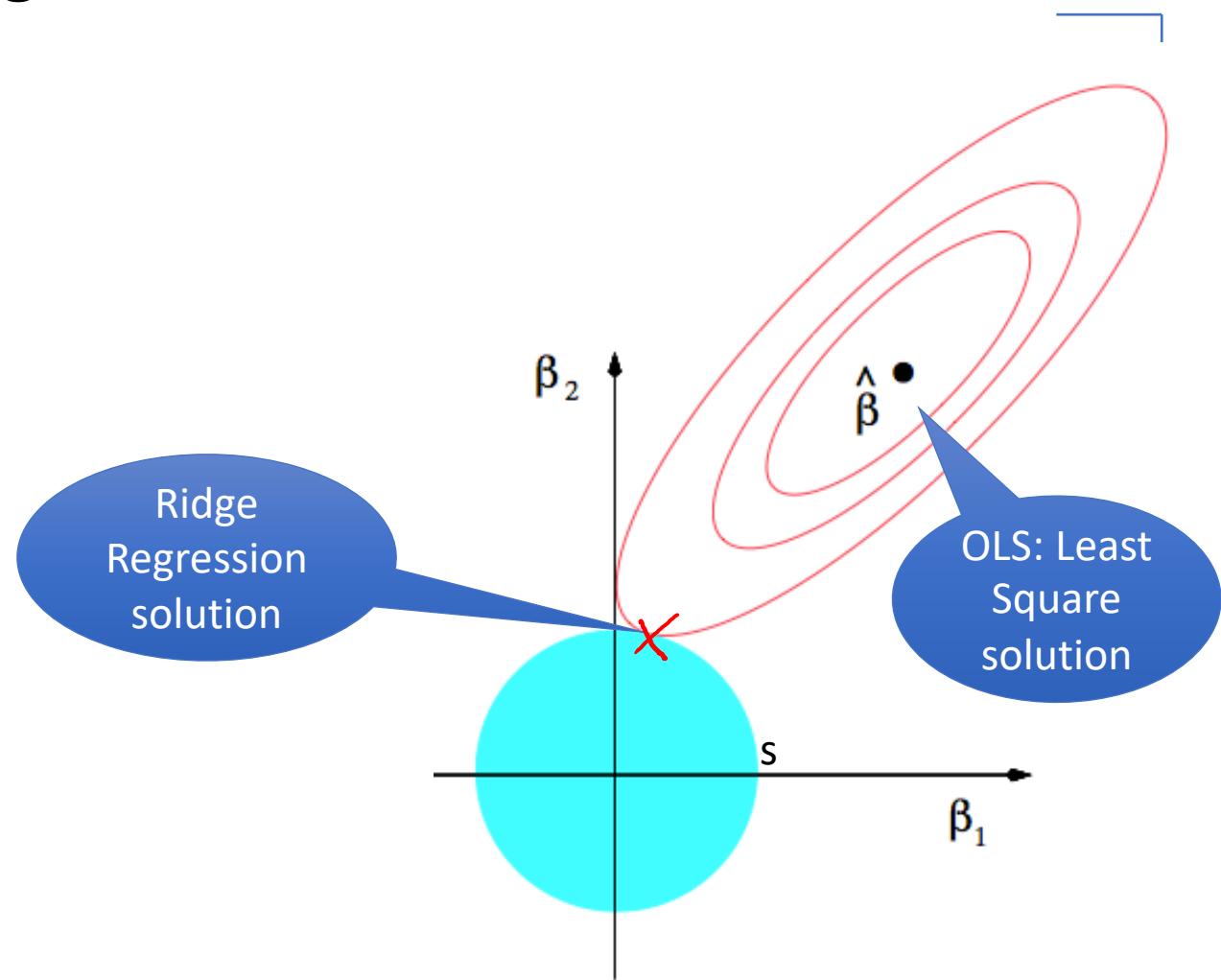
$$\hat{\beta}^{ridge} = \operatorname{argmin}(y - X\beta)^T(y - X\beta)$$

subject to  $\sum_{j=\{1..p\}} \beta_j^2 \leq s^2$

circle  
with radius  
 $s$



# Objective Function's Contour lines from Ridge Regression



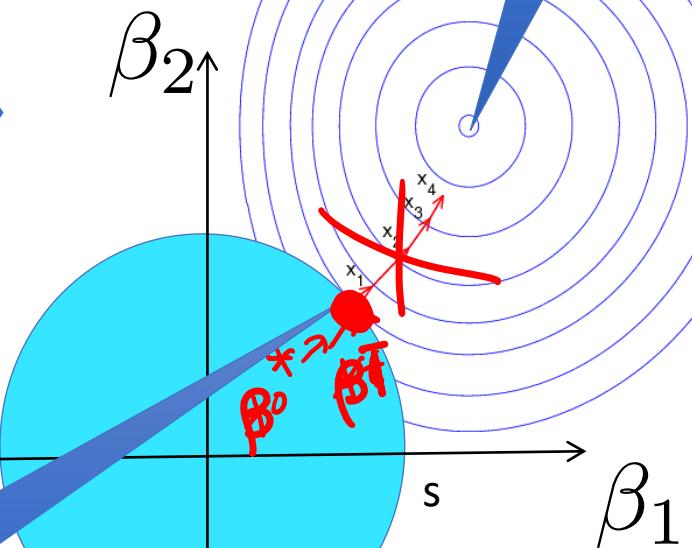
Least Square+L2:  
Ridge solution

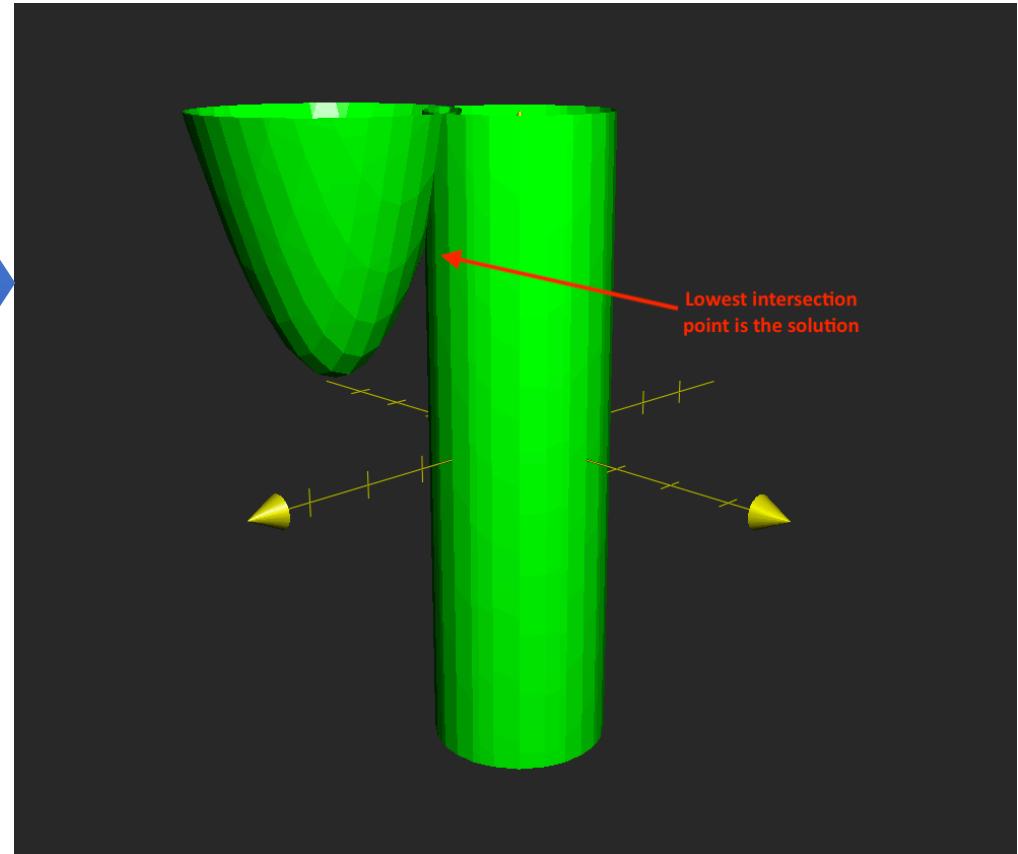


Least  
Square  
solution

Ridge  
Regression  
solution

must within the circle





# Parameter Shrinkage

$$\beta_{OLS} = (X^T X)^{-1} X^T \bar{y}$$

$$\beta_{Rg} = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

When  $X^T X = I \Rightarrow \beta_{Rg} = \frac{1}{1+\lambda} \beta_{OLS}$  [Shrinkage]

When  $X^T X$  general case, see advanced analysis @

Page65 of ESL book @

[http://statweb.stanford.edu/~tibs/ElemStatLearn/printings/ESLII\\_print10.pdf](http://statweb.stanford.edu/~tibs/ElemStatLearn/printings/ESLII_print10.pdf)

$\lambda > 0$

when  $X^T X = I$   
 $\Rightarrow$

$$\beta_{OLS} = X^T y$$

when  $X^T X = I$   
 $\Rightarrow$

$$\boxed{\beta_{Rg}} = \frac{1}{1+\lambda} X^T y = \boxed{\frac{1}{1+\lambda}} \boxed{\beta_{OLS}}$$

$\lambda > 0$

## Extra: two forms of Ridge Regression

- Totally equivalent

$$\left\{ \begin{array}{l} \text{(1) } \underset{\beta}{\operatorname{argmin}} J(\beta) + \lambda \beta^T \beta \\ \text{(2) } \underset{\beta}{\operatorname{argmin}} J(\beta), \text{ s.t. } \beta^T \beta \leq S^2 \end{array} \right.$$

Optimal Solution  $\beta_{Rg}^*$  needs (necessary condition)

$$[\lambda \left( \sum_j (\beta_{Rg})_j^2 - S^2 \right) = 0] \Rightarrow S^2 = \sum_j (\beta_{Rg})_j^2 \quad \lambda > 0$$

$$\text{When } X^T X = I, \quad S^2 = \sum_j (\beta_{Rg})_j^2 = \frac{1}{(1+\lambda)^2} \sum_j (\beta_{OLS})_j^2$$

$$\lambda = \sqrt{\frac{\sum_j (\beta_{OLS})_j^2}{S^2}} - 1 \quad \Rightarrow S^2 \propto \frac{1}{(1+\lambda)^2}$$

<http://stats.stackexchange.com/questions/190993/how-to-find-regression-coefficients-beta-in-ridge-regression>

# Ridge Regression: Squared Loss+L2

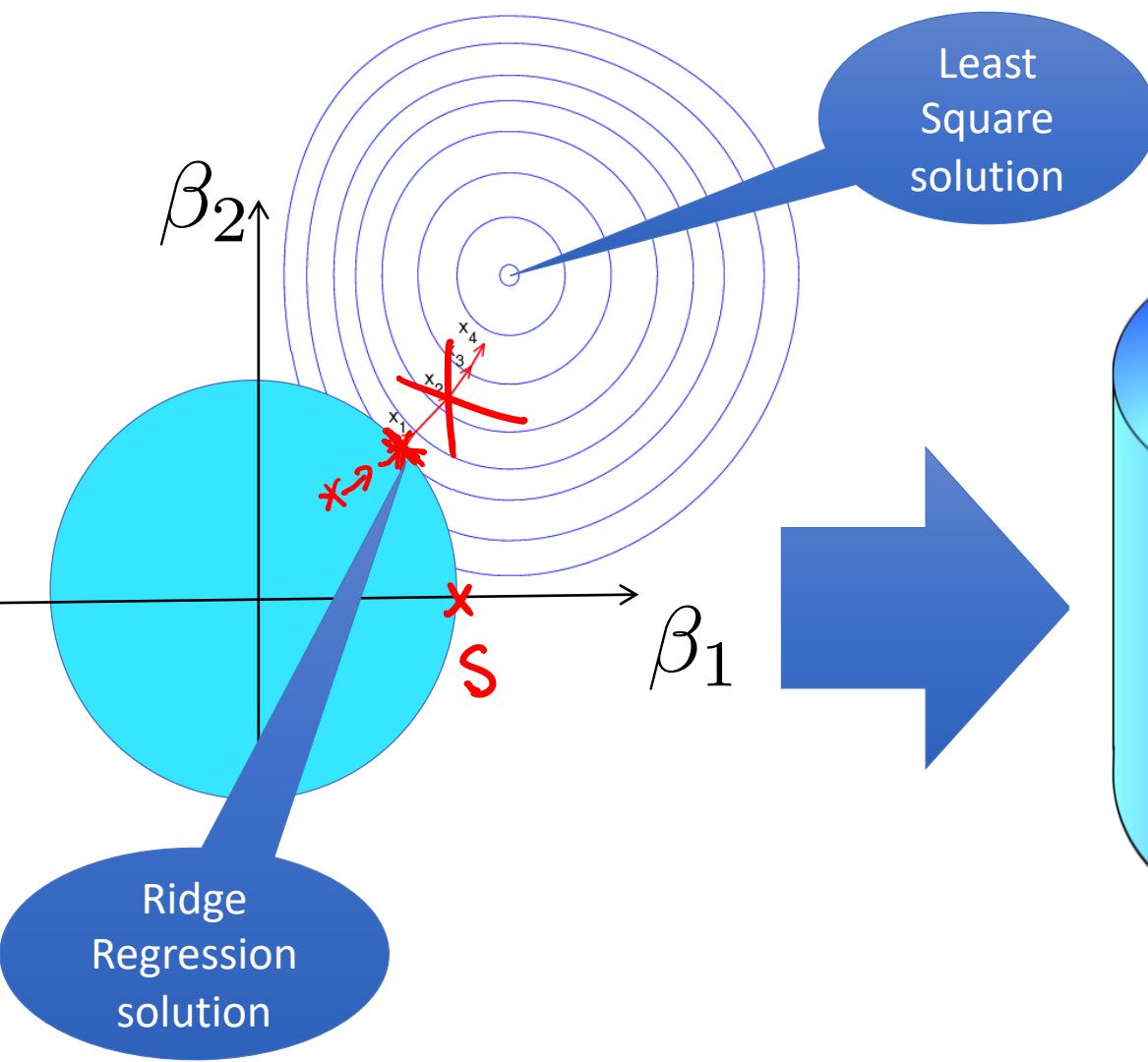
- $\lambda > 0$  penalizes each  $\beta_j$

$$\frac{1}{1 + \lambda} \hat{\beta}_{OLS}$$

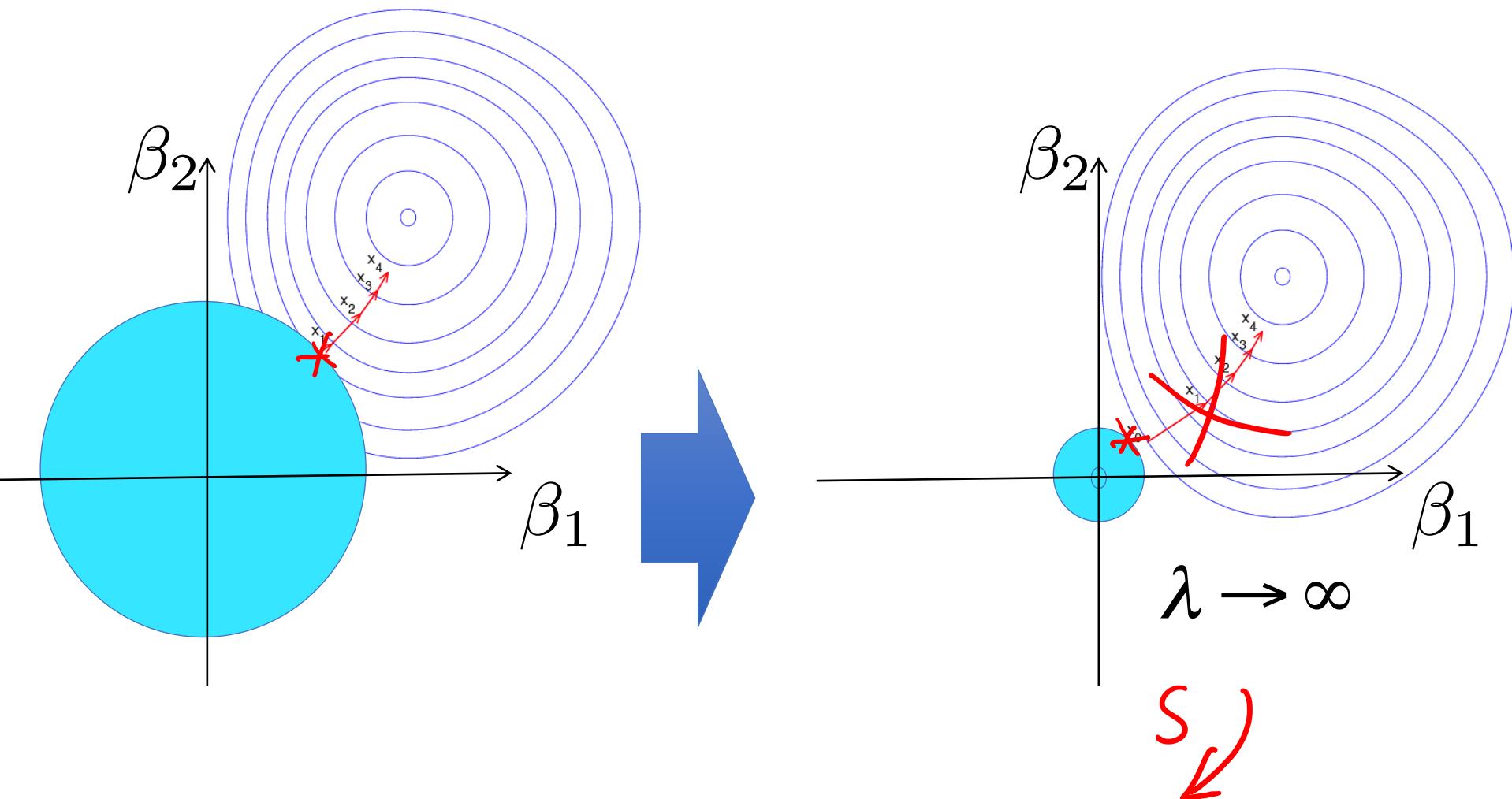
when  $X^T X = I$

- if  $\lambda = 0$  we get the least squares estimator;
- if  $\lambda \rightarrow \infty$ , then  $\beta_j$  to zero

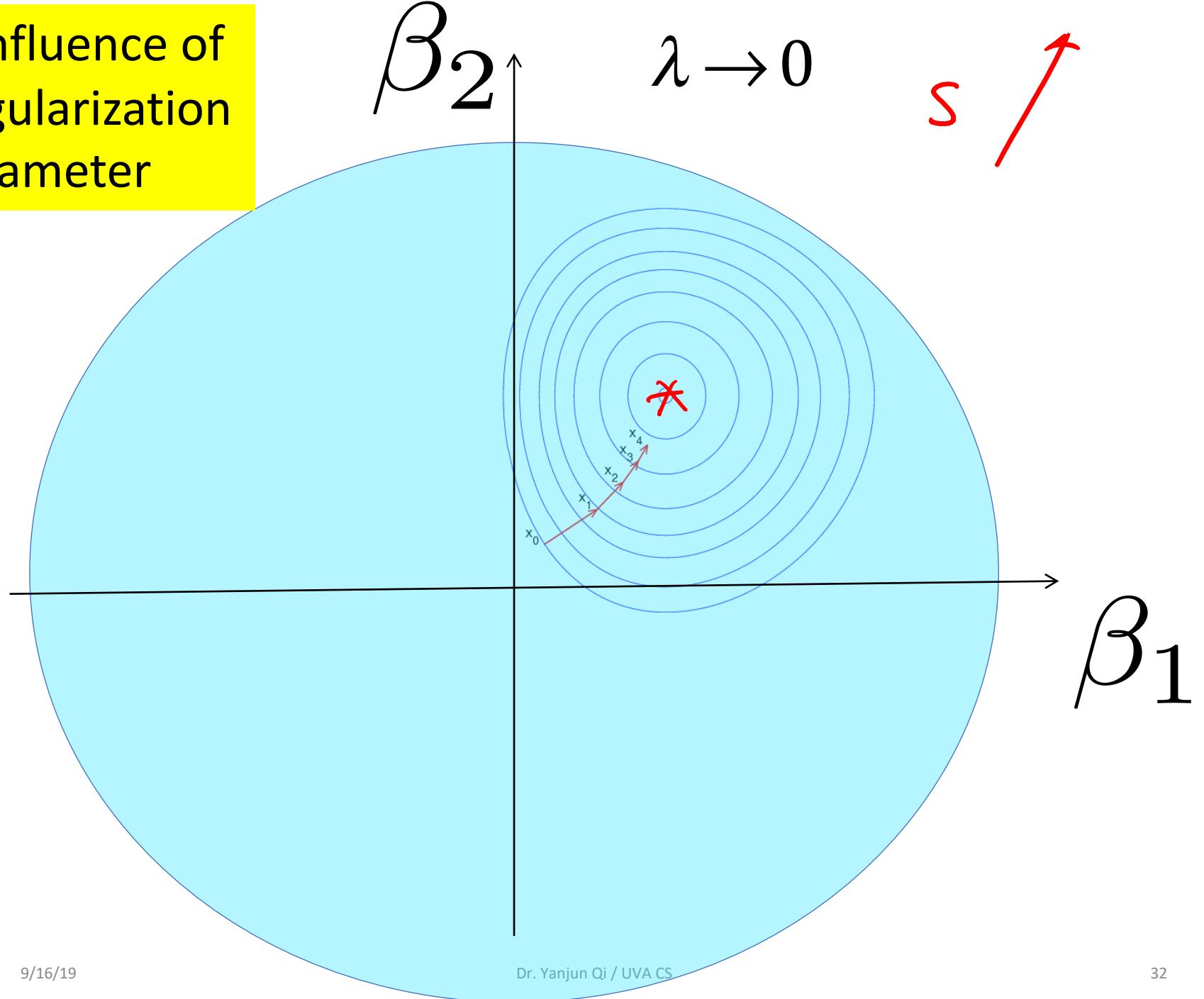
# ✓ Influence of Regularization Parameter



## ✓ Influence of Regularization Parameter



✓ Influence of  
Regularization  
Parameter



# Today

- ❑ Linear Regression Model with Regularizations
- ✓ Review: (Ordinary) Least squares: squared loss (Normal Equation)
- ✓ Ridge regression: squared loss with L2 regularization
- ✓  Lasso regression: squared loss with L1 regularization
- ✓ Elastic regression: squared loss with L1 AND L2 regularization
- ✓ How to Pick Regularization Parameter

## (2) Lasso (least absolute shrinkage and selection operator) / Squared Loss+L1

- The lasso is a shrinkage method like ridge, but acts in a nonlinear manner on the outcome  $y$ .
- The lasso is defined by

$$\hat{\beta}^{lasso} = \operatorname{argmin}_{\beta} (y - X \beta)^T (y - X \beta)$$

subject to  $\sum |\beta_j| \leq s$



L1 norm

By convention, the bias/intercept term is typically not regularized.

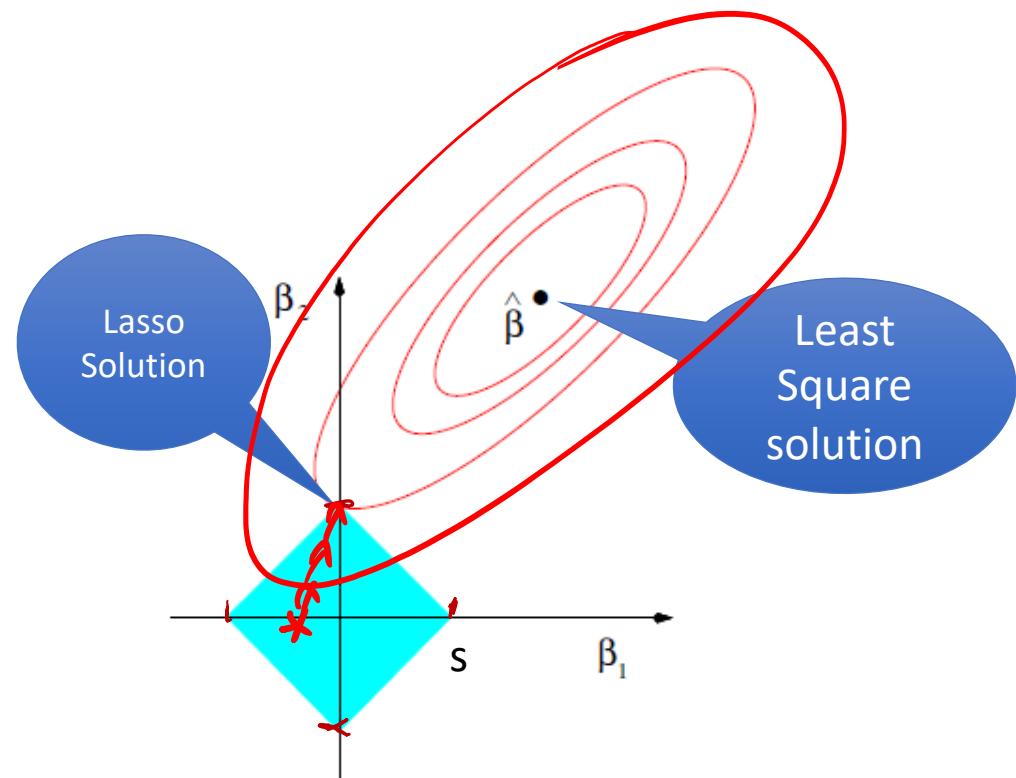
Here we assume data has been centered ... therefore no bias term

# Lasso (least absolute shrinkage and selection operator)

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}.$$

$$\beta^{\text{lasso}} = [0, s, 0]^T$$

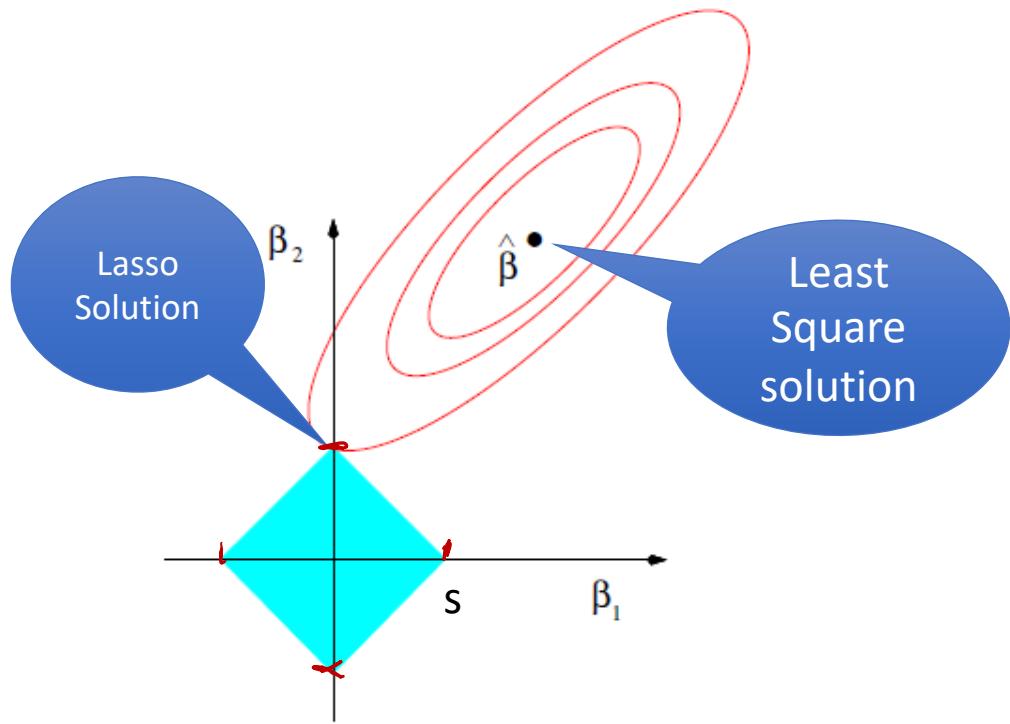
- Suppose in 2 dimension
- $\beta = (\beta_1, \beta_2)$
- $|\beta_1| + |\beta_2| = \text{const}$
- $|\beta_1| + |-\beta_2| = \text{const}$
- $|- \beta_1| + |\beta_2| = \text{const}$
- $|- \beta_1| + |-\beta_2| = \text{const}$



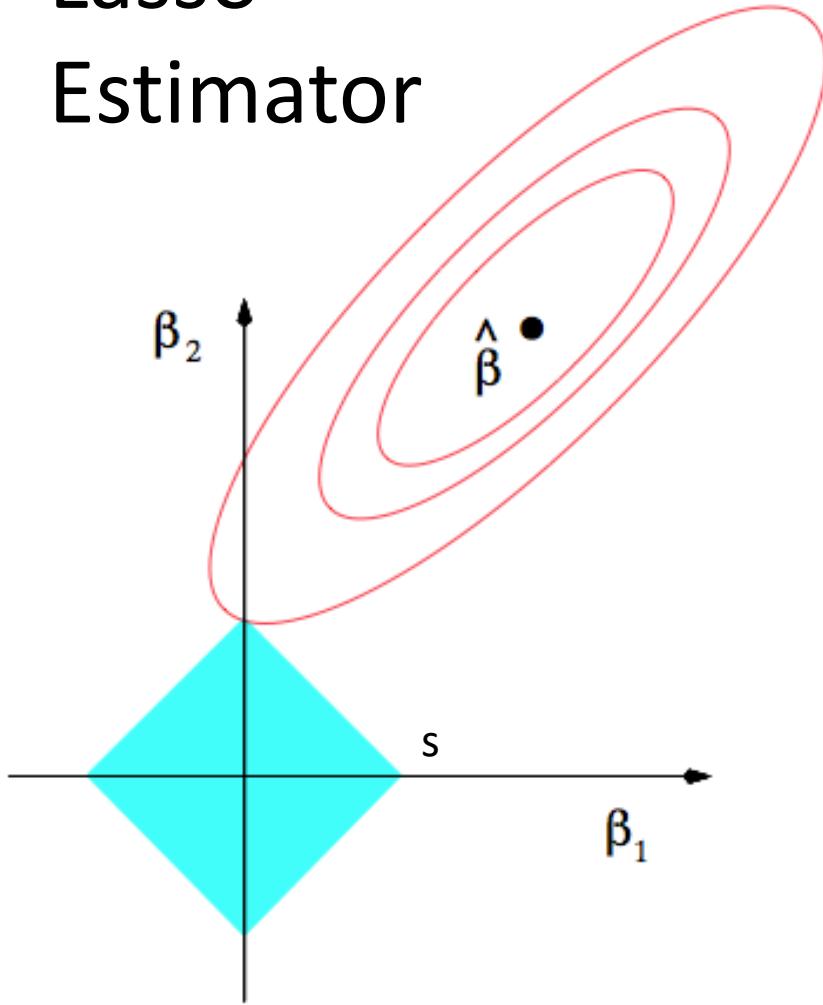
$$\hat{y} = \sum_{j=1}^p \beta_j x_j$$

when many  $\beta_j$  are zero  
 $\Rightarrow$  select feature

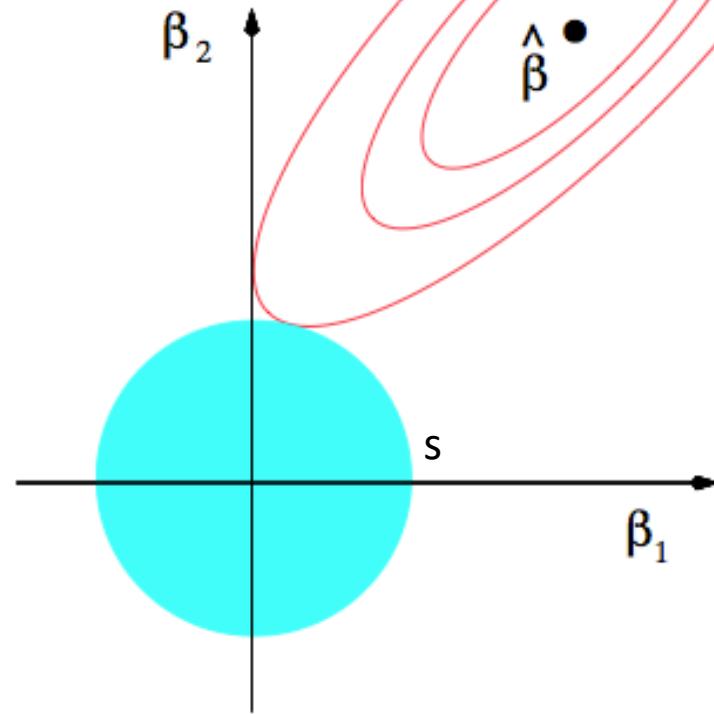
- In the Figure, the solution has eliminated the role of  $x_2$ , leading to sparsity



# Lasso Estimator



# Ridge Regression



**FIGURE 3.11.** Estimation picture for the lasso (left) and ridge regression (right). Shown are contours of the error and constraint functions. The solid blue areas are the constraint regions  $|\beta_1| + |\beta_2| \leq t$  and  $\beta_1^2 + \beta_2^2 \leq t^2$ , respectively, while the red ellipses are the contours of the least squares error function.

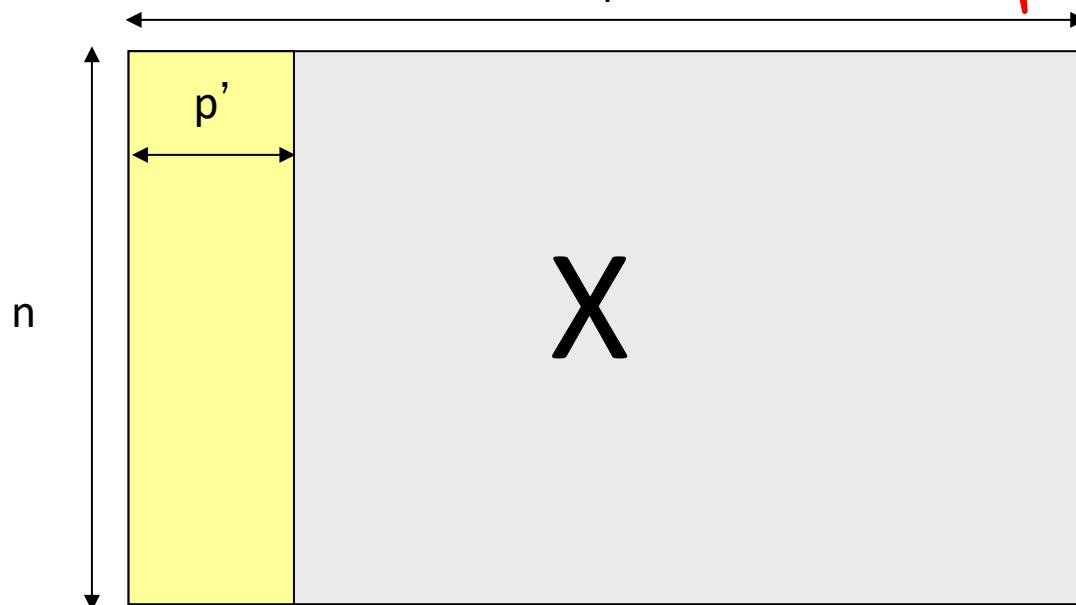
# Lasso (least absolute shrinkage and selection operator)

- Notice that ridge penalty is replaced by  $\sum |\beta_j|$
- Due to the nature of the constraint, if tuning parameter is chosen small enough, then the lasso will set some coefficients exactly to zero.

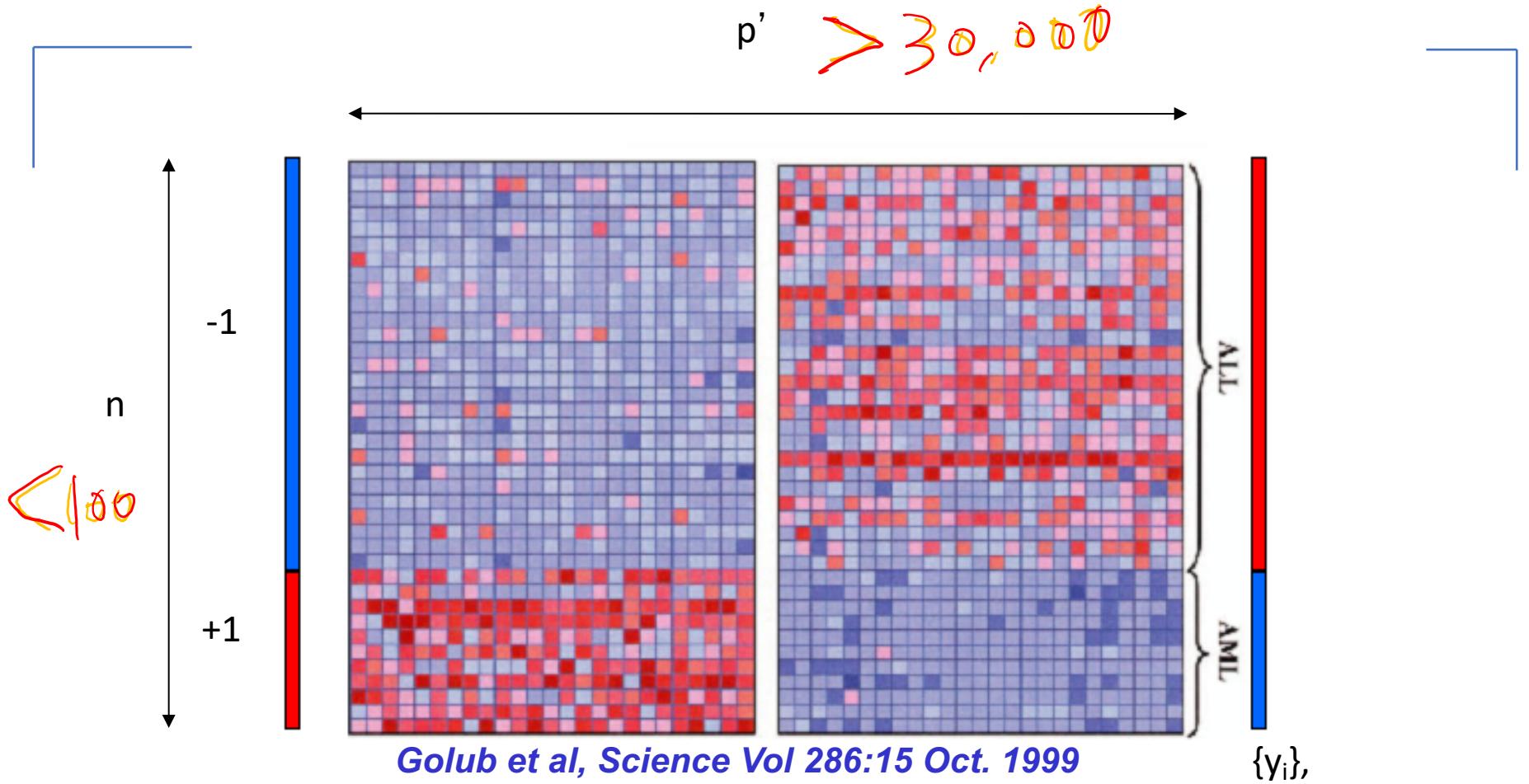
$$\sum \beta_j^2$$

# Lasso: Implicit Feature Selection

$p \rightarrow p' \Rightarrow \begin{cases} \text{easy to understand} \\ \text{Computational efficient} \end{cases}$



## e.g., Leukemia Diagnosis



$$(\Sigma^T \Sigma + \lambda I)^{-1} \Sigma y$$

When  $n < p$ ,  $O(p^3)$

Computationally,

$$\Rightarrow \underbrace{\Sigma^T \Sigma}_{p \times n \quad n \times p} : O(n p^2)$$

choose to

$$\Rightarrow \underbrace{(\Sigma^T \Sigma + \lambda I)^{-1}}_{p \times p} : O(p^3)$$

make  $p \downarrow$   
if we can

$$\Rightarrow \Sigma y : O(n p)$$

operational mode

$$\underbrace{\Sigma y}_{n \times p \quad p \times 1} : O(n' p)$$

★

# Today



Linear Regression Model with Regularizations

- ✓ Review: (Ordinary) Least squares: squared loss (Normal Equation)
- ✓ Ridge regression: squared loss with L2 regularization
- ✓ Lasso regression: squared loss with L1 regularization
- Elastic regression: squared loss with L1 AND L2 regularization
- ✓ How to Pick Regularization Parameter

# Lasso for when $p > n$

- Prediction **accuracy and model interpretation** are two important aspects of regression models.
- LASSO does **shrinkage and variable selection** simultaneously for better prediction and model interpretation.

## Disadvantage:

- In  $p > n$  case, lasso selects at most  $n$  variable before it saturates
- If there is a group of variables among which the pairwise correlations are very high, then lasso select one from the group

## (3) Hybrid of Ridge and Lasso : Elastic Net regularization

- L1 part of the penalty generates a sparse model
- L2 part of the penalty (extra):
  - Remove the limitation of the number of selected variables
  - Encouraging group effect
  - Stabilize the L1 regularization path

# Naïve elastic net

- For any non negative fixed  $\lambda_1$  and  $\lambda_2$ , naive elastic net criterion:

$$L(\lambda_1, \lambda_2, \beta) = |\mathbf{y} - \mathbf{X}\beta|^2 + \lambda_2|\beta|^2 + \lambda_1|\beta|_1,$$

$$|\beta|^2 = \sum_{j=1}^p \beta_j^2, \quad |\beta|_1 = \sum_{j=1}^p |\beta_j|.$$

- The naive elastic net estimator is the minimizer of equation

$$\hat{\beta} = \arg \min_{\beta} \{L(\lambda_1, \lambda_2, \beta)\}.$$

- Let

$$\alpha = \lambda_2 / (\lambda_1 + \lambda_2)$$

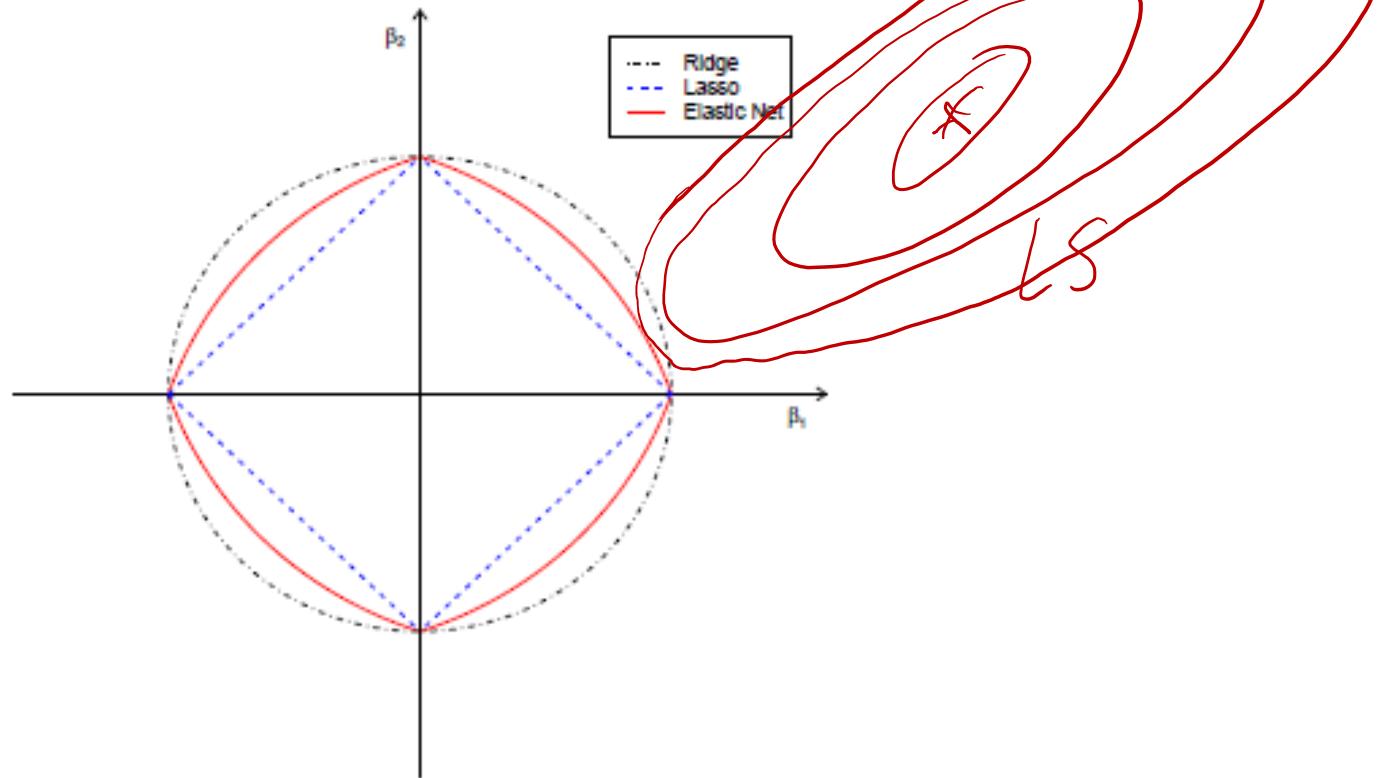
$$\frac{\lambda_2}{\lambda_1 + \lambda_2}$$

$$\hat{\beta} = \arg \min_{\beta} |\mathbf{y} - \mathbf{X}\beta|^2,$$

subject to  $(1 - \alpha)|\beta|_1 + \alpha|\beta|^2 \leq t$  for some  $t$ .

# Geometry of elastic net

2-dimensional illustration  $\alpha = 0.5$



# e.g. A Practical Application of Regression Model

## **Movie Reviews and Revenues: An Experiment in Text Regression\***

**Mahesh Joshi Dipanjan Das Kevin Gimpel Noah A. Smith**

Language Technologies Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213, USA

{maheshj, dipanjan, kgimpel, nasmith}@cs.cmu.edu

### **Abstract**

We consider the problem of predicting a movie's opening weekend revenue. Previous work on this problem has used metadata about a movie—e.g., its genre, MPAA rating, and cast—with very limited work making use of text *about* the movie. In this paper, we use the text of film critics' reviews from several sources to predict opening weekend revenue. We describe a new dataset pairing movie reviews with metadata and revenue data, and show that review text can substitute for metadata, and even improve over it, for prediction.

Proceedings of  
HLT '2010  
Human  
Language  
Technologies:

## I. The Story in Short

- ❖ Use metadata and critics' reviews to predict opening weekend revenues of movies
- ❖ Feature analysis shows what aspects of reviews predict box office success

## II. Data

- ❖ 1718 Movies, released 2005-2009
- ❖ Metadata (genre, rating, running time, actors, director, etc.): [www.metacritic.com](http://www.metacritic.com)
- ❖ Critics' reviews (~7K): Austin Chronicle, Boston Globe, Entertainment Weekly, LA Times, NY Times, Variety, Village Voice
- ❖ Opening weekend revenues and number of opening screens: [www.the-numbers.com](http://www.the-numbers.com)

e.g., Movie Reviews and Revenues: An Experiment in Text Regression, Proceedings of HLT '10 (1.7k n / >3k features)

## IV. Features

|      |  |
|------|--|
| I    | Lexical n-grams (1,2,3)  |
| II   | Part-of-speech n-grams (1,2,3)   |
| III  | Dependency relations (nsubj,advmod,...)  |
| Meta | U.S. origin, running time, budget (log),<br># of opening screens, genre, MPAA rating, holiday release (summer, Christmas, Memorial day,... ), star power (Oscar winners, high-grossing actors) |

e.g. counts  
of a ngram in  
the text

### III. Model

- ❖ Linear regression with the elastic net (Zou and Hastie, 2005)

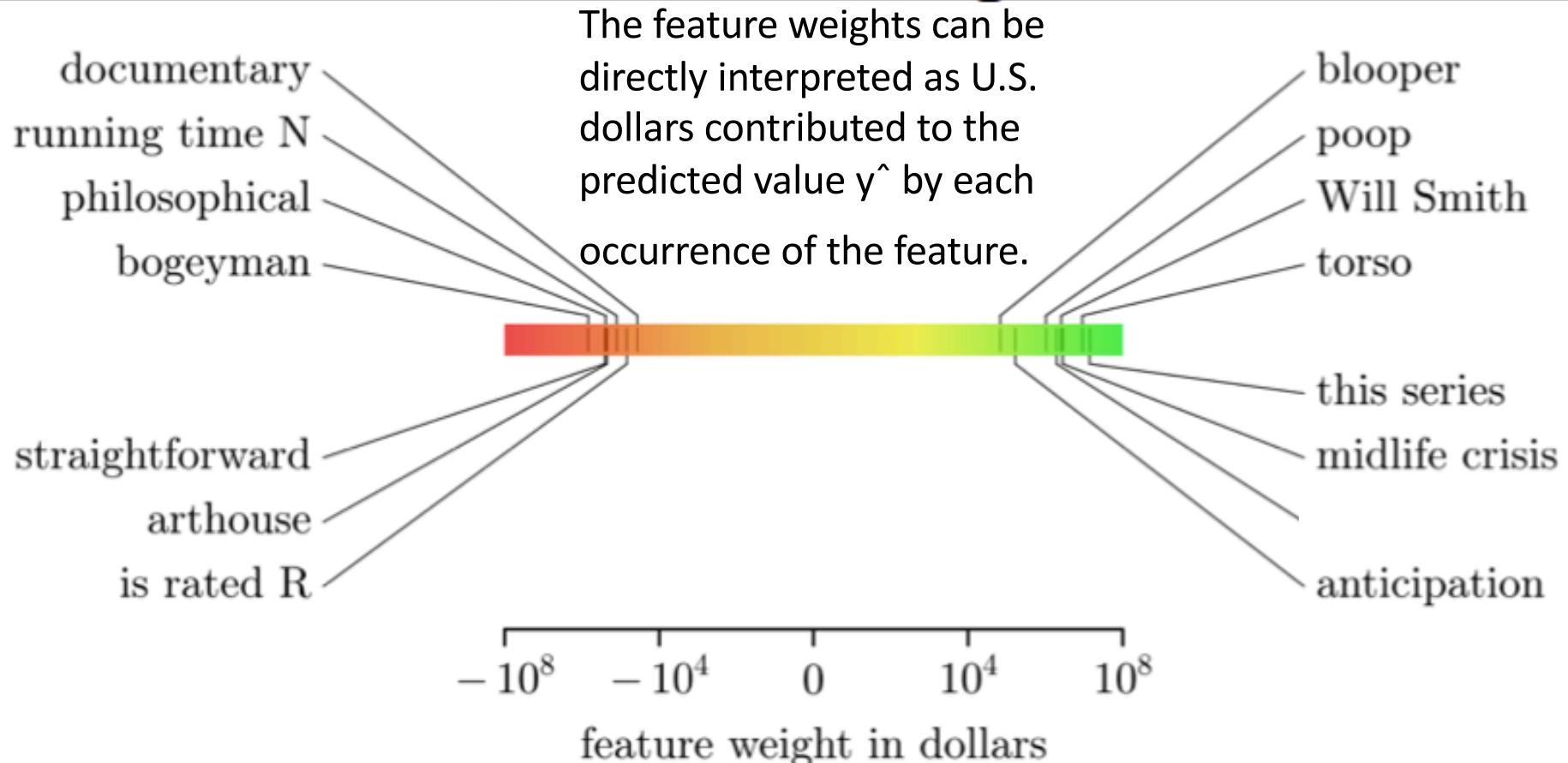
$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}=(\beta_0, \boldsymbol{\beta})}{\operatorname{argmin}} \frac{1}{2n} \left[ \sum_{i=1}^n \left( y_i - (\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}) \right)^2 \right] + \lambda P(\boldsymbol{\beta})$$
$$P(\boldsymbol{\beta}) = \sum_{j=1}^p \left( \frac{1}{2}(1-\alpha)\beta_j^2 + \alpha|\beta_j| \right)$$

Use linear regression to directly predict the opening weekend gross earnings, denoted as  $y$ , based on features  $x$  extracted from the movie metadata and/or the text of the reviews.

## VIII. Get the Data!

[www.ark.cs.cmu.edu/movies-data](http://www.ark.cs.cmu.edu/movies-data)

## V. What May Have Brought You to movies



# Advantage of Elastic net (Extra)

$P \gg n$

- Native Elastic set can be converted to lasso with augmented data form

$\Rightarrow X_{n \times p}$  (when  $n < p$ )

- In the augmented formulation,
  - sample size  $n+p$  and  $X^*$  has rank  $p$
  - $\rightarrow$  can potentially select all the predictors
- Naïve elastic net can perform automatic variable selection like lasso

# Summary: Regularized multivariate linear regression

• Model:  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_p x_p$

- LR estimation:

$$\arg \min \sum \left( Y - \hat{Y} \right)^2$$

- LASSO estimation:

$$\arg \min \sum_{i=1}^n \left( Y - \hat{Y} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- Ridge regression estimation:

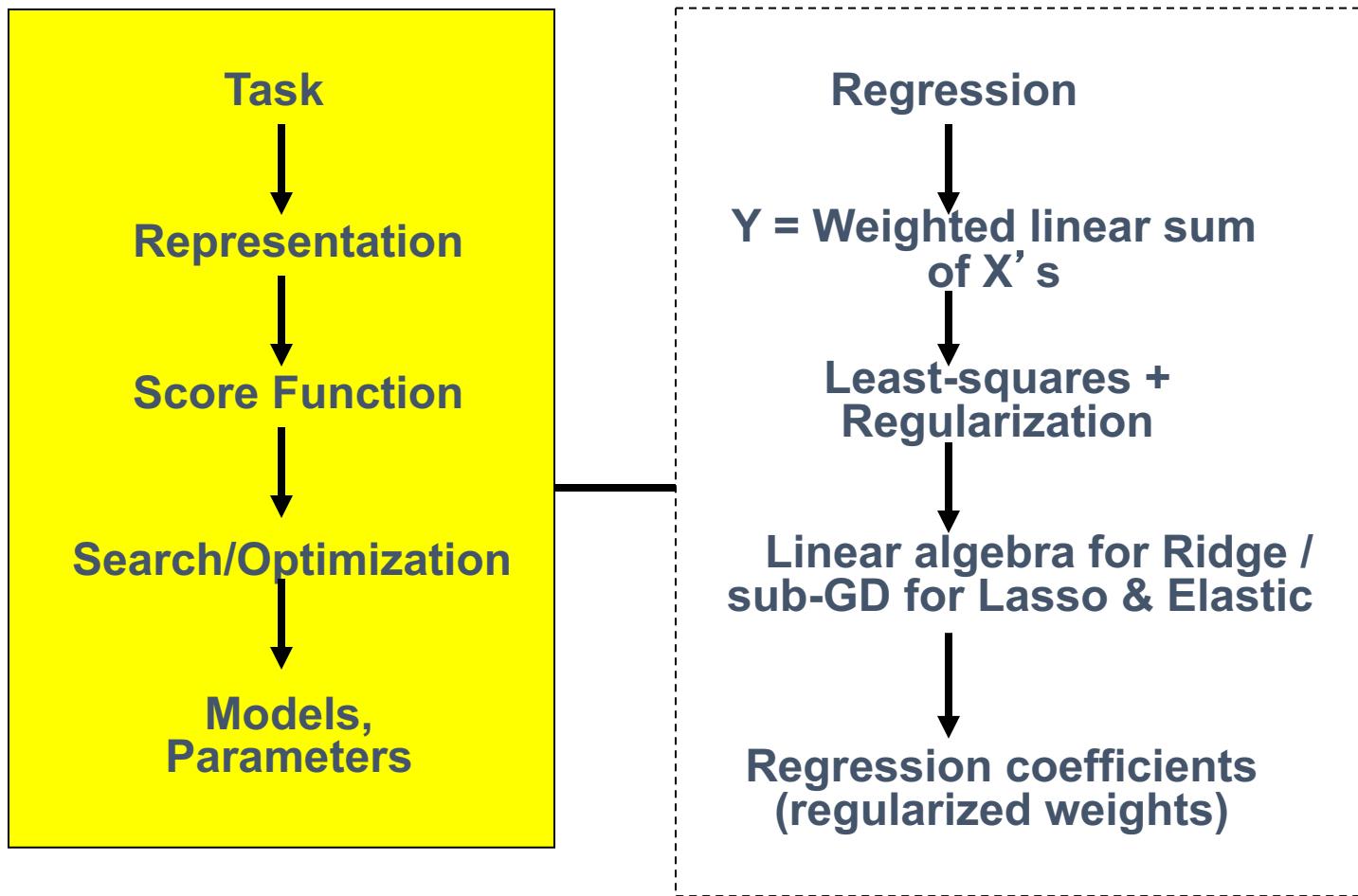
$$\arg \min \sum_{i=1}^n \left( Y - \hat{Y} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

53/54

Error on data

+ Regularization

# Regularized multivariate linear regression



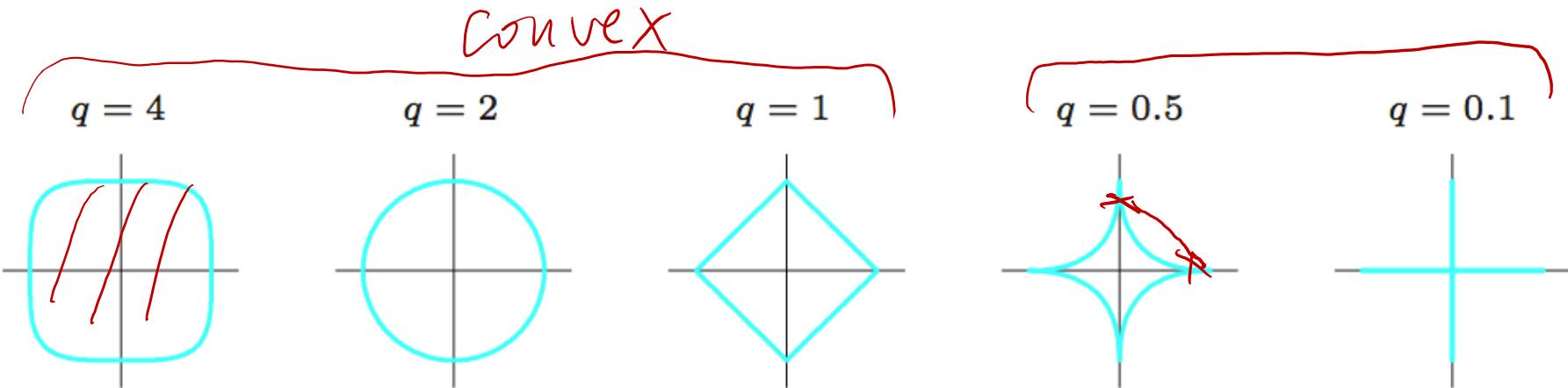
$$\min J(\beta) = \sum_{i=1}^n \left( Y - \hat{Y} \right)^2 + \lambda \left( \sum_{j=1}^p \beta_j^q \right)^{1/q}$$

# More: A family of shrinkage estimators

$$\beta = \arg \min_{\beta} \sum_{i=1}^N (y_i - x_i^T \beta)^2$$

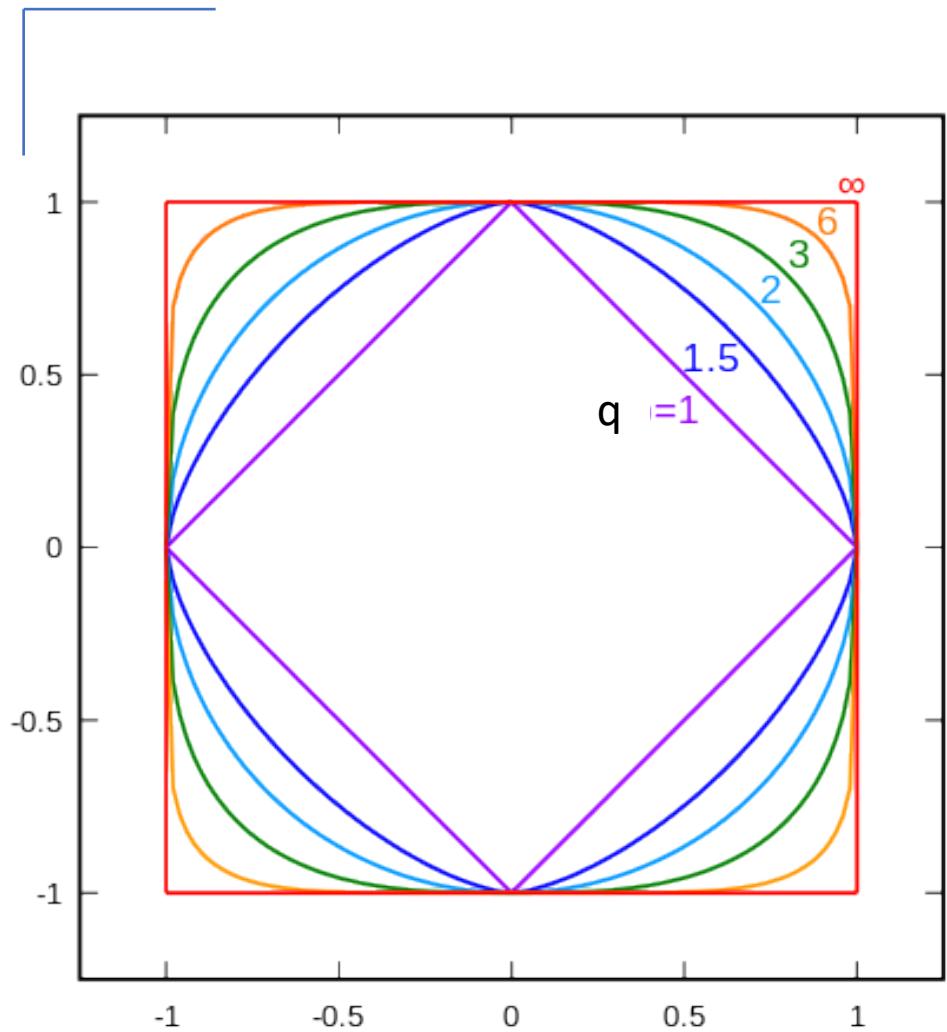
- for  $q \geq 0$ , contours of constant value of  $\sum_j |\beta_j|^q$  are shown for the case of two inputs. subject to  $\sum_j |\beta_j|^q \leq s$

$$\sum_j |\beta_j|^q$$



**FIGURE 3.12.** Contours of constant value of  $\sum_j |\beta_j|^q$  for given values of  $q$ .

# norms visualized



all p-norms penalize larger weights

$q < 2$  tends to create sparse  
(i.e. lots of 0 weights)

$q > 2$  tends to push for similar weights

We aim to make the learned model

- 1. Generalize Well
- 2. Computational Scalable and Efficient
- 3. Robust / Trustworthy / **Interpretable**
  - Especially for some domains, this is about trust!

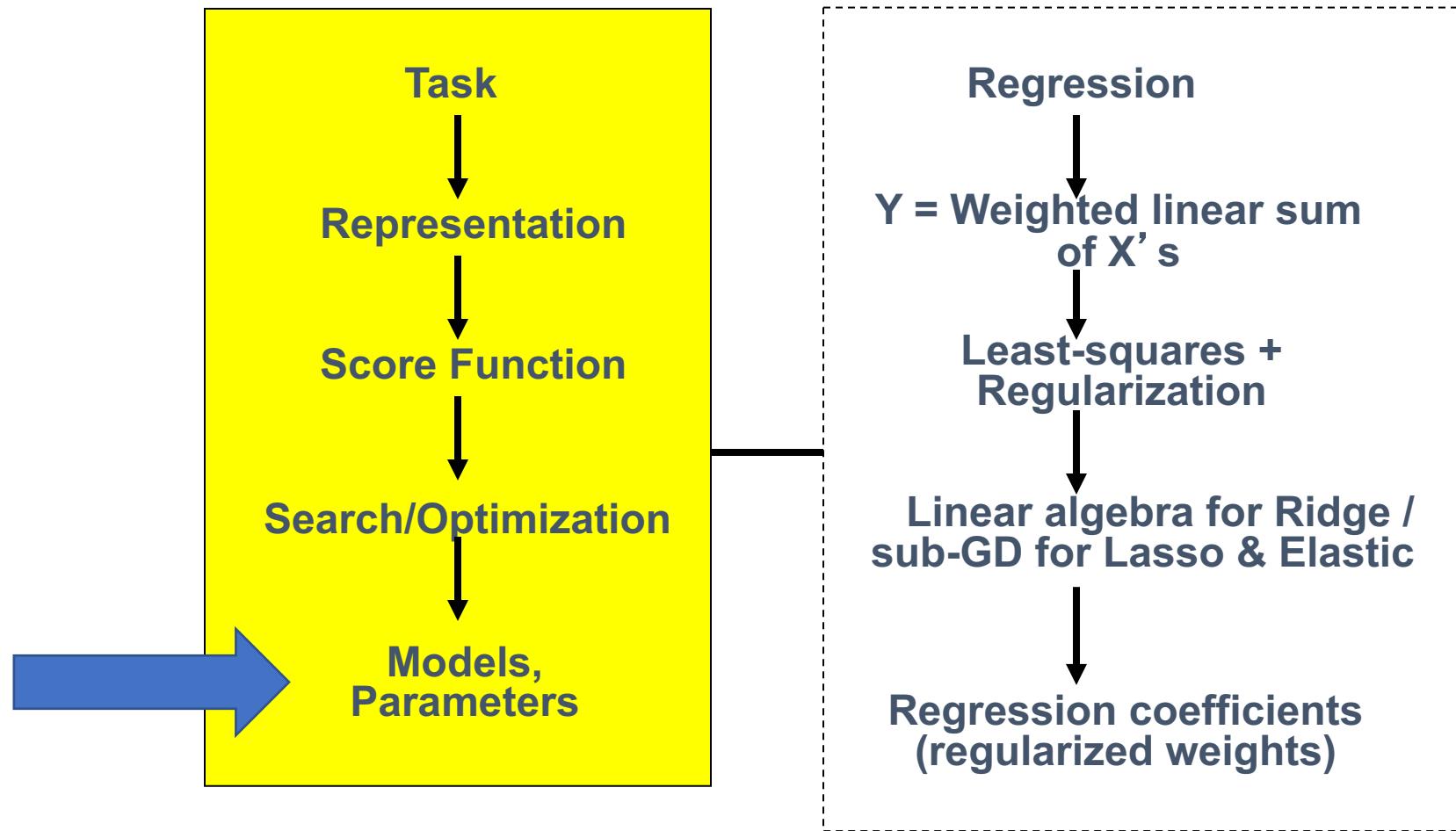
# Today



Linear Regression Model with Regularizations

- ✓ Review: (Ordinary) Least squares: squared loss (Normal Equation)
- ✓ Ridge regression: squared loss with L2 regularization
- ✓ Lasso regression: squared loss with L1 regularization
- ✓ Elastic regression: squared loss with L1 AND L2 regularization
- ✗ How to pick Regularization Parameter

# Regularized multivariate linear regression



$$\min J(\beta) = \sum_{i=1}^n \left( Y - \hat{Y} \right)^2 + \lambda \left( \sum_{j=1}^p \beta_j^q \right)^{1/q}$$

# Common regularizers

L2: Squared weights penalizes large values more

L1: Sum of weights will penalize small values more

$$\sum_j |\beta_j|$$

$$\sum_j \beta_j^2$$

Generally, we don't want huge weights

If weights are large, a small change in a feature can result in a large change in the prediction

Might also prefer weights of 0 for features that aren't useful

# Why to Use simpler models?

- Because:
  - Simpler to use (lower computational complexity)
  - Easier to train (needs less examples)
  - Less sensitive to noise
  - Easier to explain (more interpretable)
  - Generalizes better (lower variance <- Occam's razor) --- **More in future lectures!!!**

# Model Selection & Generalization

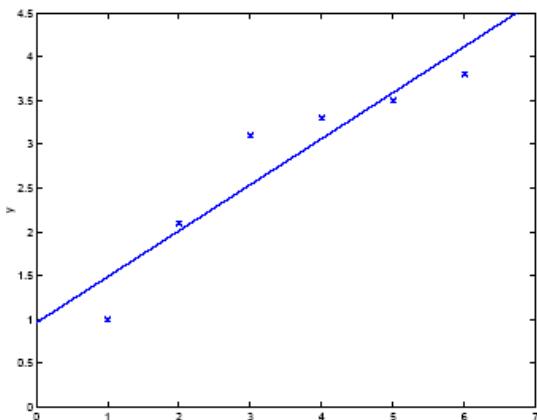
- **Generalisation**: learn function / hypothesis from **past data** in order to “explain”, “predict”, “model” or “control” **new** data examples
- Underfitting: when model is too simple, both training and test errors are large
- Overfitting: when model is too complex and test errors are large although training errors are small.
  - After learning knowledge, model tends to learn “**noise**”

# Issue: Overfitting and underfitting

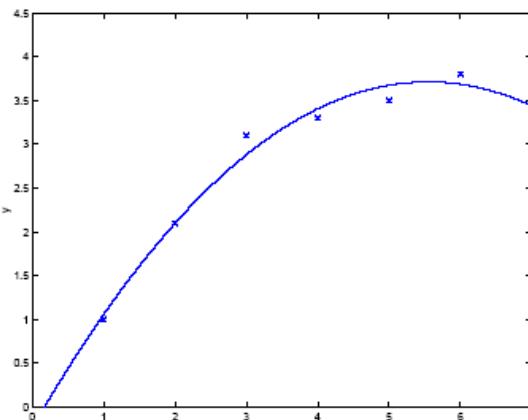
Under fit

Looks good

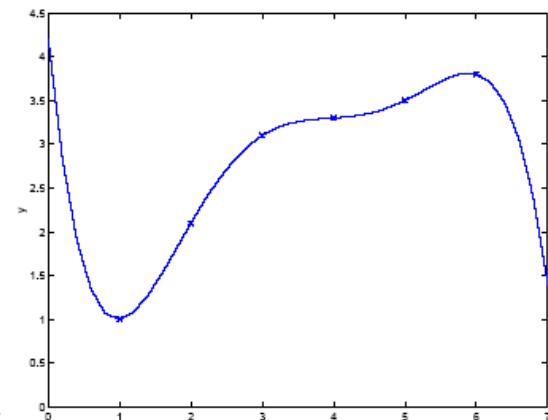
Over fit



$$y = \theta_0 + \theta_1 x$$



$$y = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$y = \sum_{j=0}^5 \theta_j x^j$$

Generalisation: learn function / hypothesis from **past data** in order to “explain”, “predict”, “model” or “control” **new** data examples

9/16/

K-fold Cross Validation !!!!

# Overfitting: Handled by Regularization

A **regularizer** is an additional criteria to the loss function to make sure that we don't overfit

It's called a regularizer since it tries to keep the parameters more normal/regular

It is a bias on the model forces the learning to prefer certain types of weights over others, e.g.,

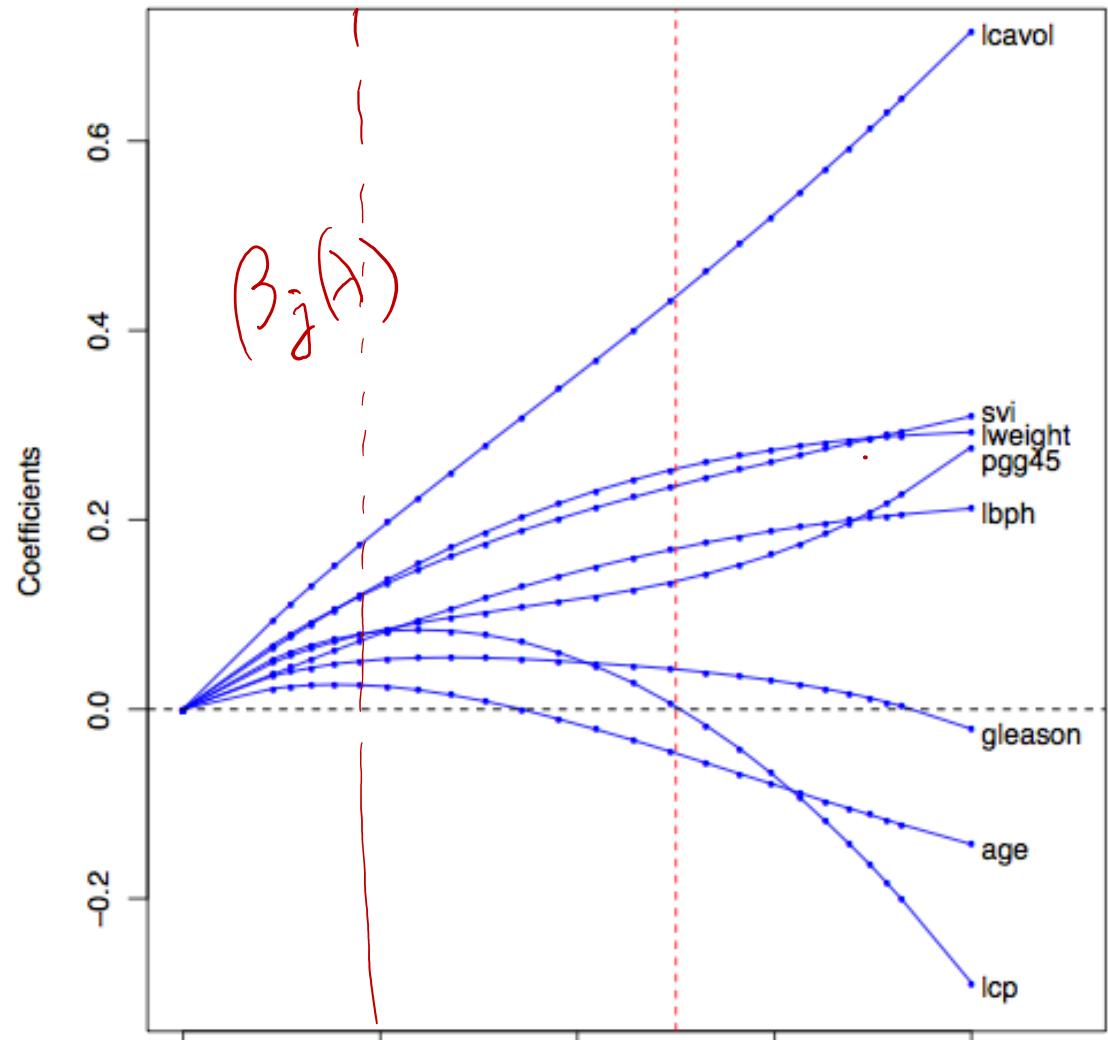
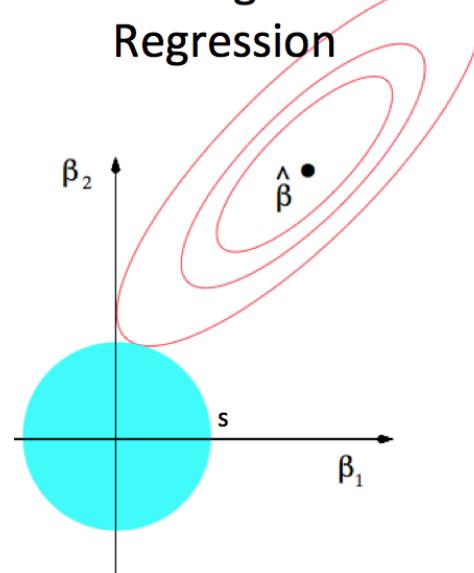
$$\hat{\beta}^{ridge} = \operatorname{argmin}_{\beta} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \beta^T \beta$$

# WHY and How to Select $\lambda$ ?

- 1. Generalization ability  
→ k-folds CV to decide
- 2. Control the bias and Variance of the model (details in future lectures)

# Regularization path of a Ridge Regression

Ridge Regression



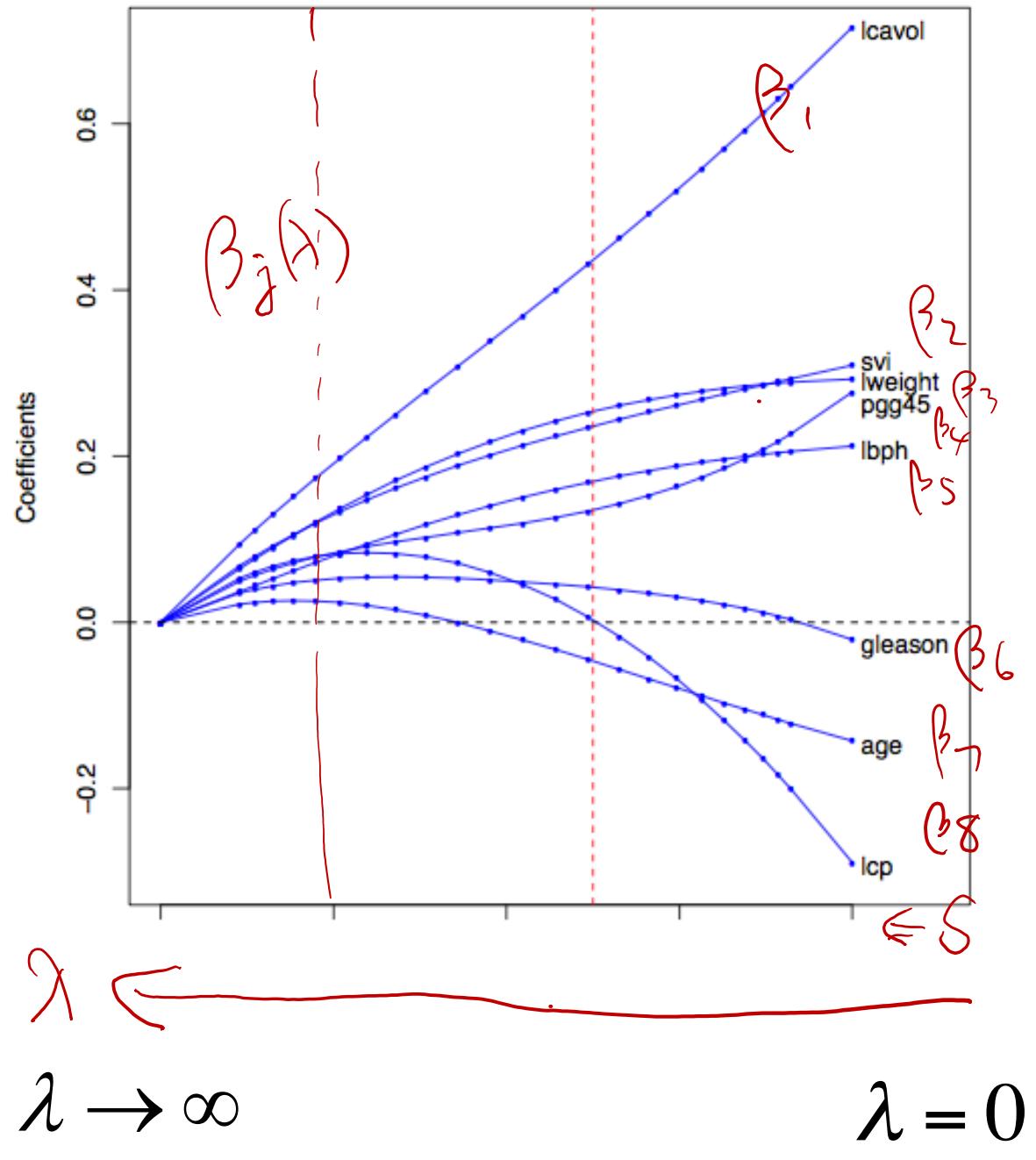
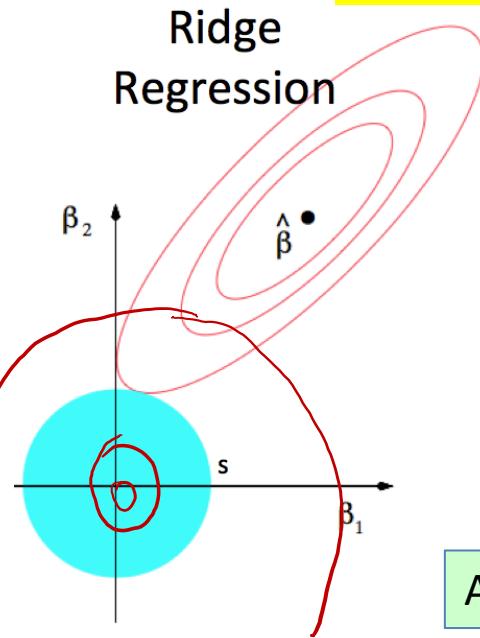
$\lambda \rightarrow \infty$

$\lambda = 0$

# Regularization path of a Ridge Regression

When  
 $\tilde{X}^T \tilde{X} = I \Rightarrow \frac{1}{1+\lambda} \beta_{OLS}$

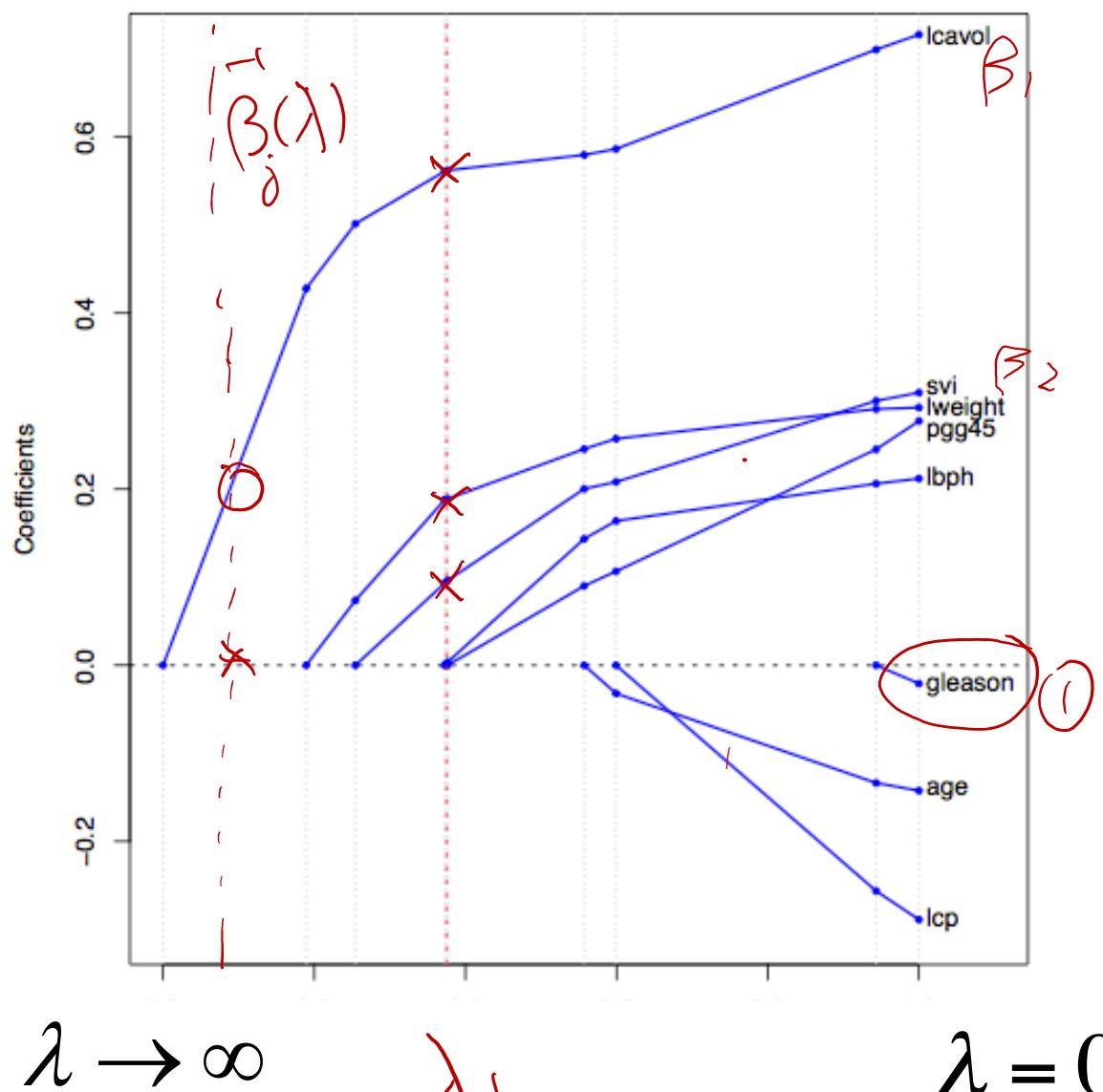
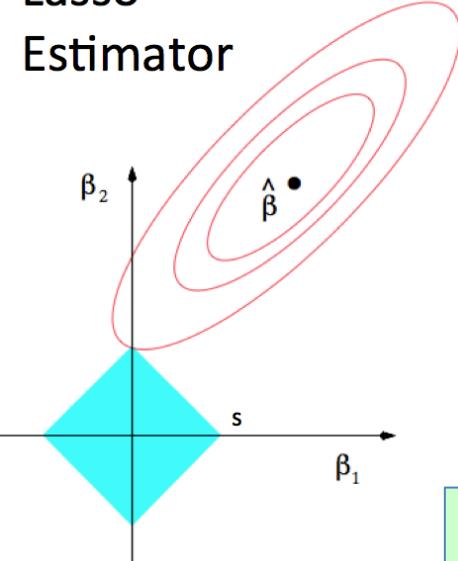
Weight Decay



# Regularization path of a Lasso Regression

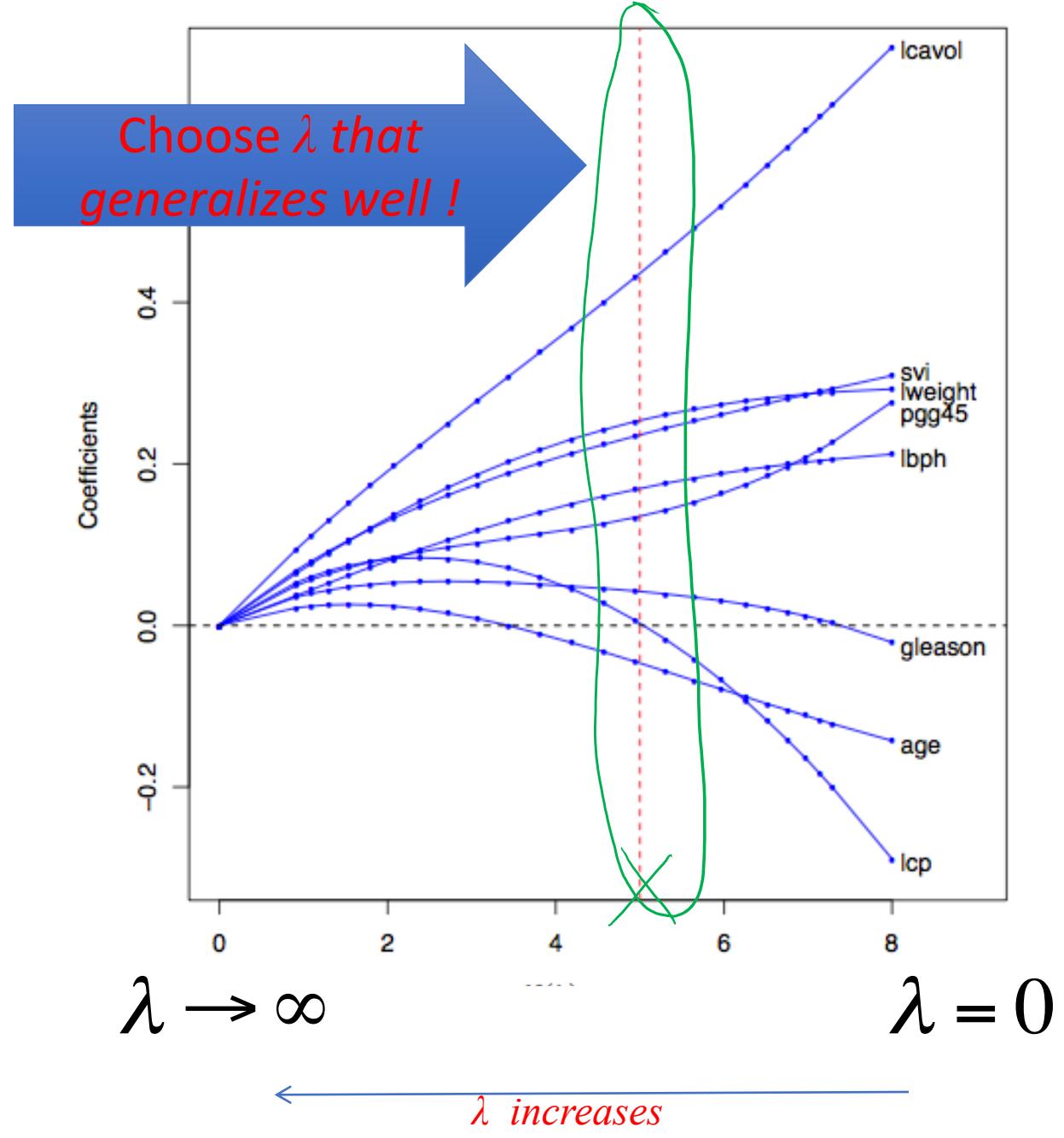
when varying  $\lambda$ ,  
how  $\beta_j$  varies.

Lasso  
Estimator



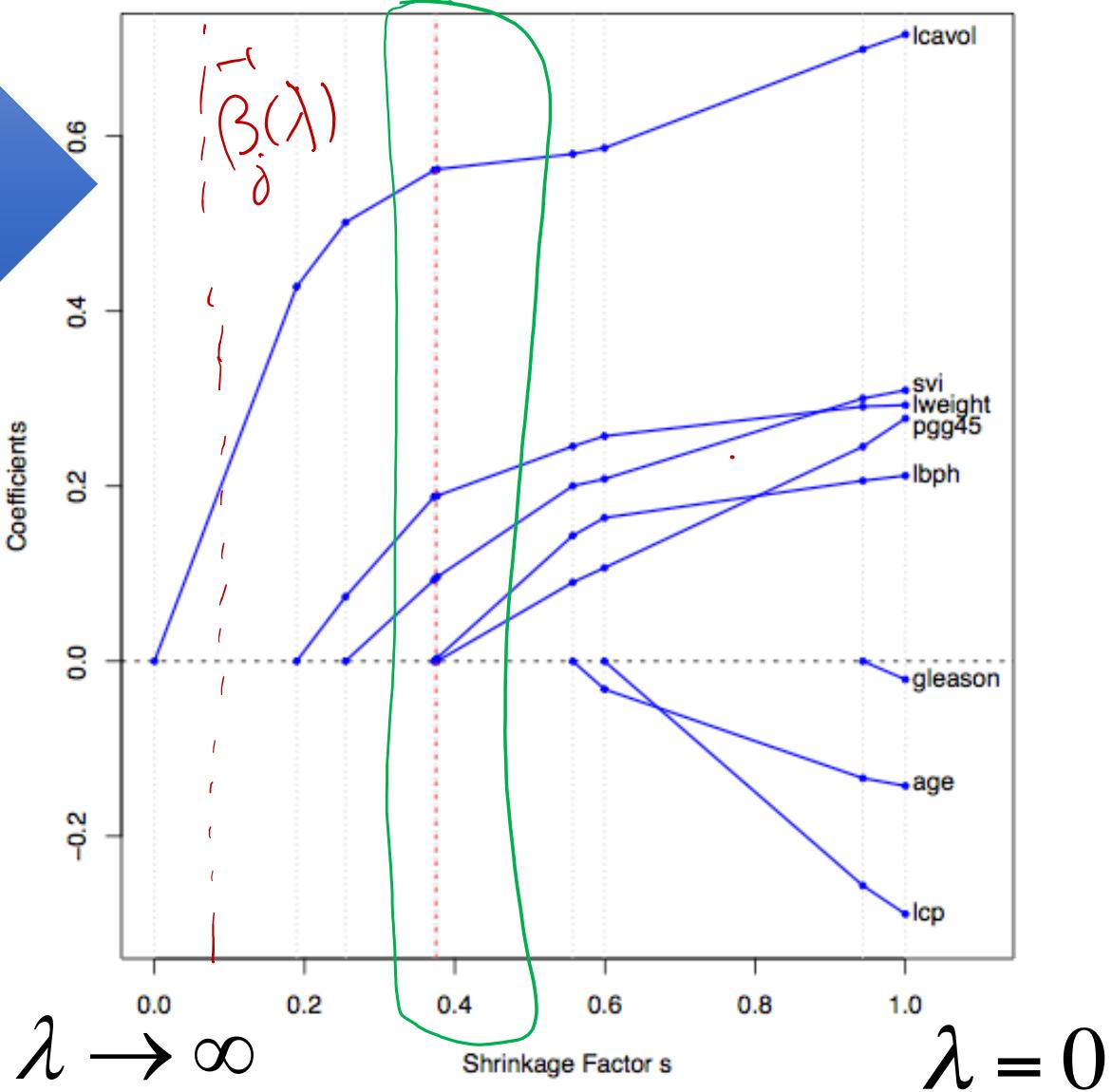
An example with 8 features

An example  
of  
Ridge Regression  
  
when varying  
 $\lambda$ , how  $\beta_j$   
varies.



Choose  $\lambda$  that generalizes well!

when varying  $\lambda$ ,  
how  $\beta_j$  varies.



**FIGURE 3.10.** Profiles of lasso coefficients, as the tuning parameter  $t$  is varied. Coefficients are plotted versus  $s = t / \sum_1^p |\hat{\beta}_j|$ . A vertical line is drawn at  $s = 0.36$ , the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso profiles hit zero, while those for ridge do not. The profiles are piece-wise linear, and so are computed only at the points displayed; see Section 3.4.4 for details.

An example with 8 features

# Today Recap



- Linear Regression Model with Regularizations
- ✓ Review: (Ordinary) Least squares: squared loss (Normal Equation)
- ✓ Ridge regression: squared loss with L2 regularization
- ✓ Lasso regression: squared loss with L1 regularization
- ✓ Elastic regression: squared loss with L1 AND L2 regularization
- ✓ Influence of Regularization Parameter

# Regression (supervised)

- ❑ Four ways to train / perform optimization for linear regression models
  - ❑ Normal Equation
  - ❑ Gradient Descent (GD)
  - ❑ Stochastic GD
  - ❑ Newton's method

} Variations of  $\underset{\theta}{\operatorname{arg\min}} L(\theta)$

- ❑ Supervised regression models
  - ❑ Linear regression (LR)
  - ❑ LR with non-linear basis functions
  - ❑ Locally weighted LR
  - ❑ LR with Regularizations

} Variations of  $f(x)$   
→ Variations of  $L(\theta)$

# References

- ❑ Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- ❑ Prof. Nando de Freitas's tutorial slide
- ❑ **Regularization and variable selection via the elastic net**, Hui Zou and Trevor Hastie, *Stanford University, USA*
- ❑ *ESL book: Elements of Statistical Learning*

# Extra More

- Optimization of regularized regressions:
  - See L6-extra slide
- Relation between  $\lambda$  and  $s$ 
  - See L6-extra slide
- Why Elastic Net has a few nice properties
  - See L6-extra slide