

# UVA CS 4501 - 001 / 6501 – 007

## Introduction to Machine Learning and Data Mining

### Lecture 6: Regression Models with Regularization

Yanjun Qi / Jane

University of Virginia  
Department of  
Computer Science

## Last Lecture Recap

- Linear model is an approximation
- Three ways to moving beyond linearity
  - LR with non-linear basis functions
  - Locally weighted linear regression
  - Regression trees and Multilinear Interpolation (later)

## (1) LR with non-linear basis functions

- LR does not mean we can only deal with linear relationships

$$y = \theta_0 + \sum_{j=1}^m \theta_j \phi_j(x) = \theta^T \phi(x)$$

- We are free to design (non-linear) features (e.g., basis function derived) under LR

where the  $\phi_j(x)$  are fixed basis functions (also define  $\phi_0(x) = 1$ ).

- E.g.: polynomial regression:

$$\phi(x) := [\mathbf{1}, x, x^2, x^3]$$

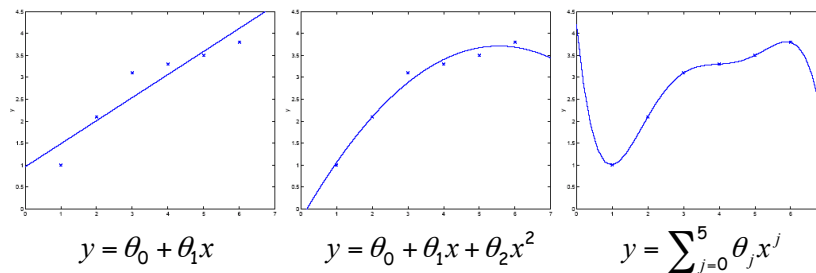
9/9/14

3

## (1) LR With basis functions

### e.g. polynomial regression

Issue: Overfitting OR underfitting



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

**K-fold Cross Validation !!!!**

9/9/14

4

## (2) Locally weighted linear regression

- The algorithm:  
Instead of minimizing

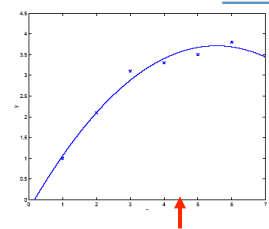
$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

now we fit  $\vartheta$  to minimize

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n w_i (\mathbf{x}_i^T \theta - y_i)^2$$

Where do  $w_i$ 's come from?

$$w_i = \exp\left(-\frac{(\mathbf{x}_i - \mathbf{x})^2}{2\tau^2}\right)$$



- where  $\mathbf{x}$  is the query point for which we'd like to know its corresponding  $y$

→ Essentially we put higher weights on (errors on) training examples that are close to the query point (than those that are further away from the query)

9/9/14

5

## Parametric vs. non-parametric

- Locally weighted linear regression is a **non-parametric** algorithm.
- The (unweighted) linear regression algorithm that we saw earlier is known as a **parametric** learning algorithm
  - because it has a fixed, finite number of parameters (the  $\theta$ ), which are fit to the data;
  - Once we've fit the  $\theta$  and stored them away, we no longer need to keep the training data around to make future predictions.
  - In contrast, to make predictions using locally weighted linear regression, we need to keep the entire training set around.
- The term "**non-parametric**" (roughly) refers to the fact that the amount of stuff we need to keep in order to represent the hypothesis grows with linearly the size of the training set.

9/9/14

6

## Today

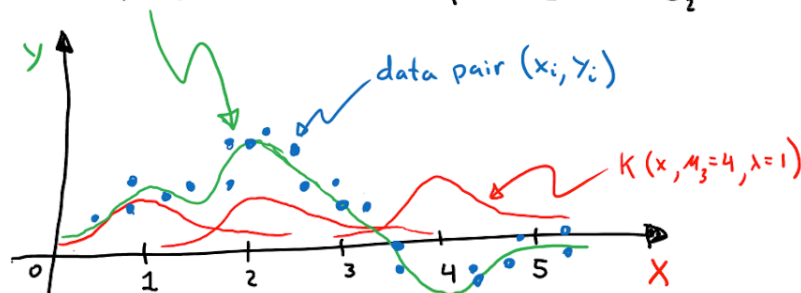
- A bit more about Linear Regression Extension
  - Linear regression with predefined RBF basis
  - Locally weighted regression
- An Exemplar Application of Regression
- Linear Regression Models with Regularization

9/9/14

7

### (1) Linear regression with predefined RBF basis functions

$$\hat{y}(x) = e^{-\|x-1\|^2} \theta_1 + e^{-\|x-2\|^2} \theta_2 + e^{-\|x-4\|^2} \theta_3$$



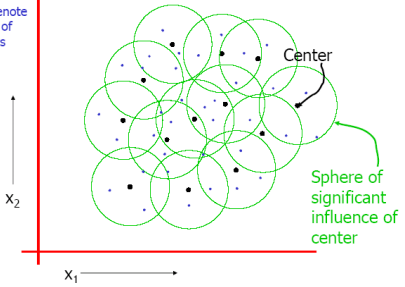
$$\varphi(x) := [1, k(x, 1, 1), k(x, 2, 1), k(x, 4, 1)]$$

9/9/14

YanJun Qi / UVA CS 4501-01-6501-07

## Issue: Choices of Basis Functions: → Good and Bad RBFs

- A good 2D RBF

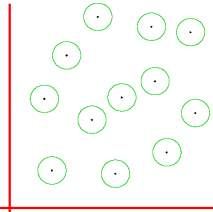


Blue dots denote coordinates of input vectors

Center

Sphere of significant influence of center

- A bad 2D RBFs

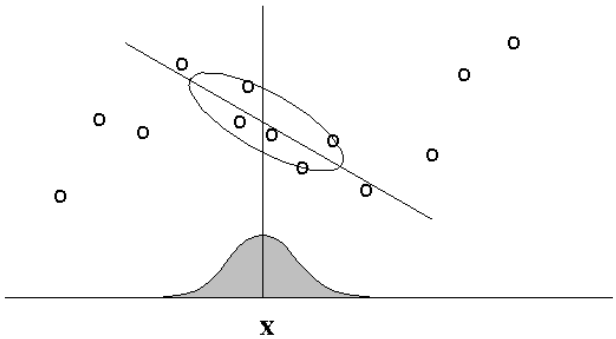


9/9/14

10

## (2) Locally weighted regression

- *aka* locally weighted regression, locally linear regression, LOESS, ...

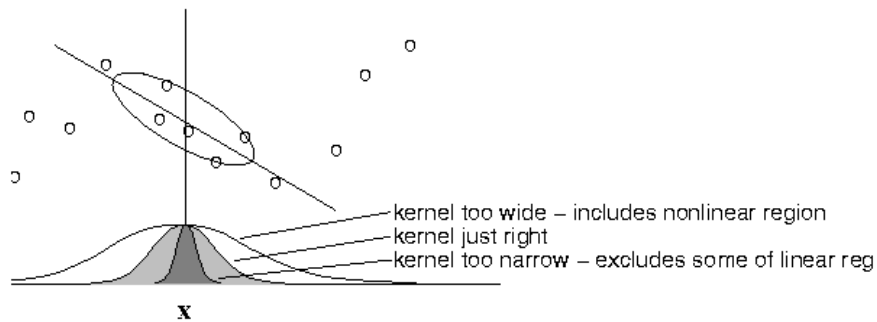


**x**

**Figure 2:** In locally weighted regression, points are weighted by proximity to the current  $x$  in question using a kernel. A regression is then computed using the weighted points.

## (2) Locally weighted linear regression

11



**Figure 3:** The estimator variance is minimized when the kernel includes as many training points as can be accommodated by the model. Here the linear LOESS model is shown. Too large a kernel includes points that degrade the fit; too small a kernel neglects points that increase confidence in the fit.

## (2) Locally weighted linear regression

e.g. when for only one feature variable

- Separate weighted least squares **at each target point  $x_0$** :

$$\min_{\alpha(x_0), \beta(x_0)} \sum_{i=1}^N K_{\lambda}(x_0, x_i) [y_i - \alpha(x_0) - \beta(x_0)x_i]^2$$

$$\hat{f}(x_0) = \hat{\alpha}(x_0) + \hat{\beta}(x_0)x_0$$

- $b(x)^T = (1, x)$ ;  $B$ :  $N \times 2$  regression matrix with  $i$ -th row  $b(x)^T$ ;  $W_{N \times N}(x_0) = \text{diag}(K_{\lambda}(x_0, x_i)), i = 1, \dots, N$

$$\hat{f}(x_0) = b(x_0)^T (B^T W(x_0) B)^{-1} B^T W(x_0) y$$



$$\text{LR } \hat{f}(x_q) = (x_q)^T \theta^* = (x_q)^T (X^T X)^{-1} X^T \bar{y}$$

9/11/14

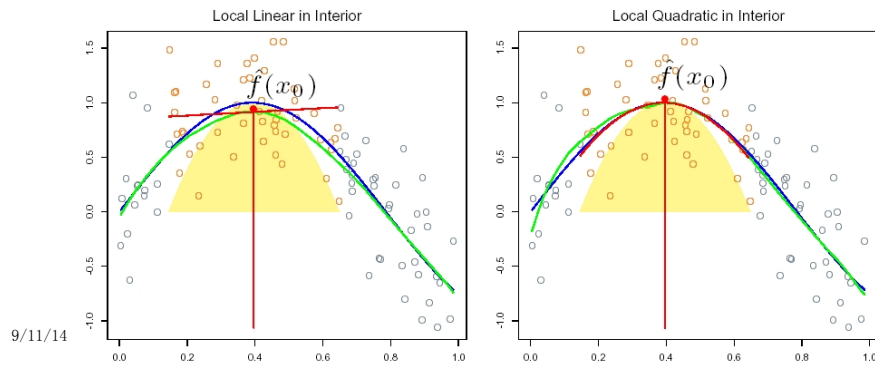
12

## (2). One More → Local Weighted Polynomial Regression

- Local polynomial fits of any degree  $d$

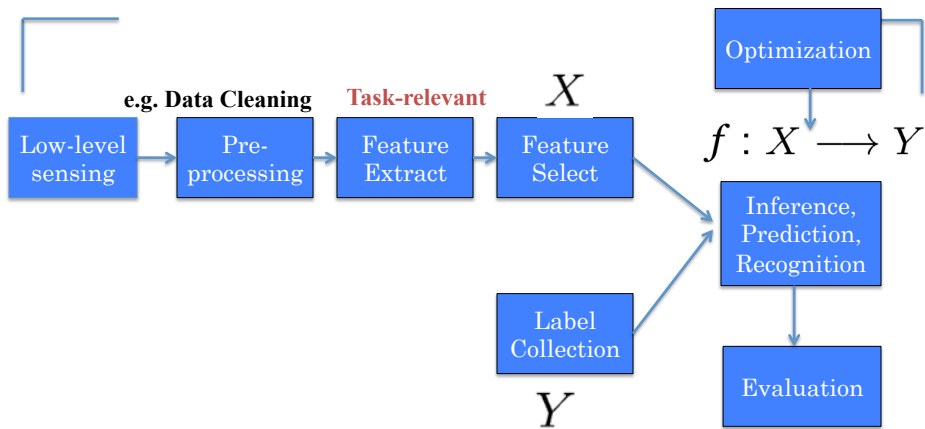
$$\min_{\alpha(x_0), \beta_j(x_0), j=1, \dots, d} \sum_{i=1}^N K_\lambda(x_0, x_i) \left[ y_i - \alpha(x_0) - \sum_{j=1}^d \beta_j(x_0) x_i^j \right]^2$$

$$\hat{f}(x_0) = \hat{\alpha}(x_0) + \sum_{j=1}^d \hat{\beta}_j(x_0) x_0^j$$



9/11/14

## TYPICAL MACHINE LEARNING SYSTEM



8/26/14

14

## Today

- ❑ A bit more about Linear Regression Extension
  - ❑ Linear regression with predefined RBF basis
  - ❑ Locally weighted regression
- ❑ An Exemplar Application of Regression
- ❑ Linear Regression Models with Regularization

## 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.



## I. The Story in Short

/ UVA CS 4501-01-6501-07

- ❖ 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)

9/9/14

Yanjun Qi / UVA CS 4501-01-6501-07

Movie Reviews and Revenues: An Experiment in Text Regression,  
 Proceedings of HLT '10 Human Language Technologies:

## III. Model

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

$$\hat{\theta} = \operatorname{argmin}_{\theta=(\beta_0, \beta)} \frac{1}{2n} \sum_{i=1}^n \left( y_i - (\beta_0 + \mathbf{x}_i^\top \beta) \right)^2 + \lambda P(\beta)$$

$$P(\beta) = \sum_{j=1}^p \left( \frac{1}{2}(1 - \alpha)\beta_j^2 + \alpha|\beta_j| \right)$$

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

18

Yanjun Qi / UVA CS 4501-01-6501-07

Movie Reviews and Revenues: An Experiment in Text Regression,  
 Proceedings of HLT '10 Human Language Technologies:

### IV. Features

|             |  |
|-------------|--|
| <b>I</b>    | Lexical n-grams (1,2,3)  |
| <b>II</b>   | Part-of-speech n-grams (1,2,3)   |
| <b>III</b>  | Dependency relations (nsubj,advmod,...)  |
| <b>Meta</b> | U.S. origin, running time, budget (log), # 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

9/9/14
19

Yanjun Qi / UVA CS 4501-01-6501-07

➔

### VIII. Get the Data!

[www.ark.cs.cmu.edu/movie\\$-data](http://www.ark.cs.cmu.edu/movie$-data)

### V. What May Have Brought You to movies

documentary

running time N

philosophical

bogeyman

straightforward

arthouse

is rated R

The feature weights can be directly interpreted as U.S. dollars contributed to the predicted value  $\hat{y}$  by each occurrence of the feature.

feature weight in dollars

blooper

poop

Will Smith

torso

this series

midlife crisis

anticipation

## Today

- A bit more about Linear Regression Extension
  - Linear regression with predefined RBF basis
  - Locally weighted regression
  
- An Exemplar Application of Regression
  
- Linear Regression Model with Regularizations
  - Ridge Regression
  - Lasso Regression

9/9/14

21

## Review: Vector norms

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

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

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

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

–  $L_{\text{infinity}}$

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

8/28/14

22

Yanjun Qi / UVA CS 4501-01-6501-07

## Review: Vector Norm (L2, when p=2)

$$\left\| \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right\|_2 = \sqrt{1^2 + 2^2} = \sqrt{5}$$

8/28/14 23

Yanjun Qi / UVA CS 4501-01-6501-07

## Review: Normal equation for LR

- Write the cost function in matrix form:

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

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

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

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

The normal equations

$$\theta^* = (X^T X)^{-1} X^T \bar{\mathbf{y}}$$

Assume that  $X^T X$  is invertible

9/2/14 24

## (1) Ridge Regression / L2

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

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_p x_p \quad \text{Basic Model,}$$

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

- The ridge estimator is solution from

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

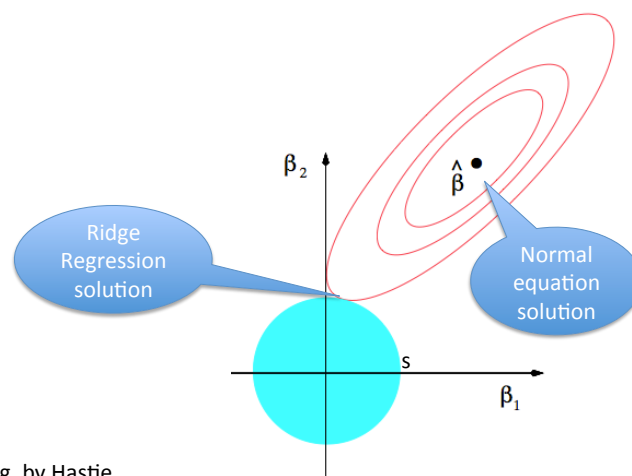
- Equivalently

$$\hat{\beta}^{ridge} = \operatorname{argmin} (y - X\beta)^T (y - X\beta) \\ \text{subject to } \sum \beta_j^2 \leq s$$

25

## Objective Function's Contour lines from Ridge Regression

Yanjun Qi / UVA CS 4501-01-6501-07



Elements of Statistical Learning, by Hastie,  
Tibshirani and Friedman

26

## (1) Ridge Regression / L2

- The parameter  $\lambda > 0$  penalizes  $\beta_j$  proportional to its size  $\beta_j^2$
- Solution is  $\hat{\beta}_\lambda = (X^T X + \lambda I)^{-1} X^T y$
- where  $I$  is the identity matrix.
- Note  $\lambda = 0$  gives the least squares estimator;
- if  $\lambda \rightarrow \infty$ , then  $\hat{\beta} \rightarrow 0$

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

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

$$\hat{\beta}^{lasso} = \arg \min (y - X\beta)^T (y - X\beta)$$

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

## Lasso (least absolute shrinkage and selection operator)

- Notice that ridge penalty  $\sum \beta_j^2$  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.

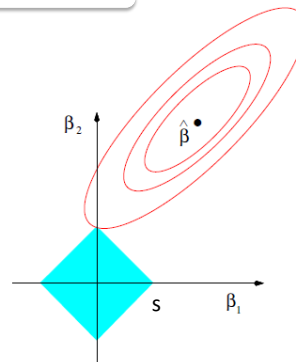
9/11/14

29

## Lasso (least absolute shrinkage and selection)

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

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



Elements of Statistical Learning, by Hastie, Tibshirani and Friedman

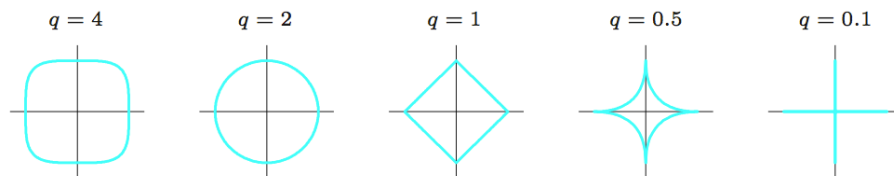
30

## (3) A family of shrinkage estimators

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

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

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



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

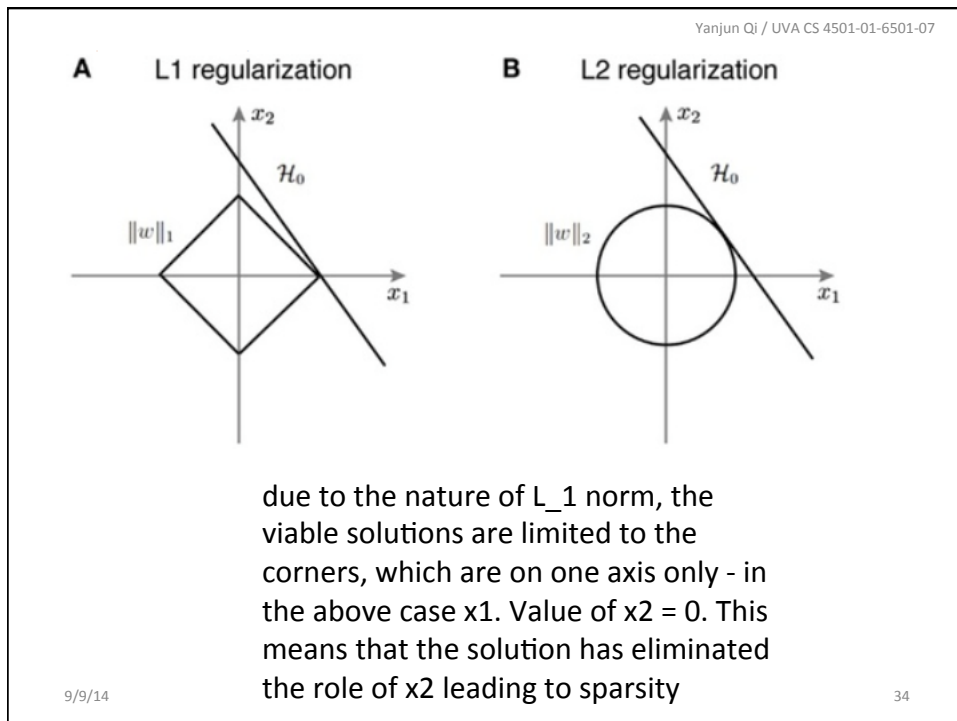
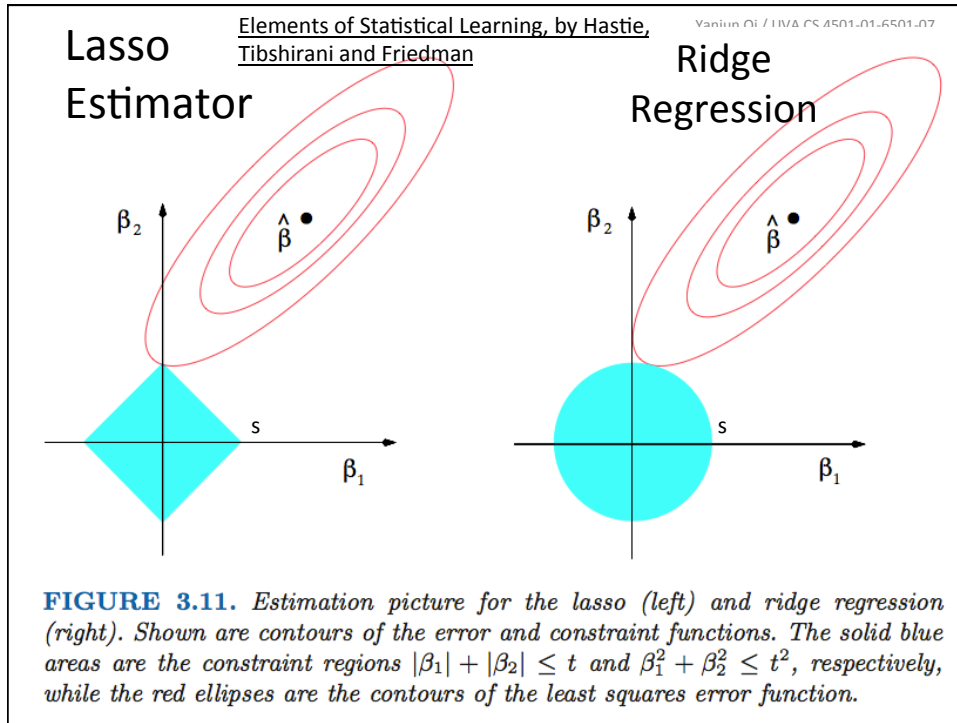
## In the example: Hybrid of Ridge and Lasso

### Elastic Net regularization

$$\hat{\beta} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda_2 \|\beta\|^2 + \lambda_1 \|\beta\|_1$$

- The  $\ell_1$  part of the penalty generates a sparse model.
- The quadratic part of the penalty
  - Removes the limitation on the number of selected variables;
  - Encourages *grouping effect*;
  - Stabilizes the  $\ell_1$  regularization path.





## Summary: Regularized multivariate linear regression

- **Model:**  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$
- LR estimation:  $\min SSE = \sum (Y - \hat{Y})^2$
  - LASSO estimation:  $\min SSE = \sum_{i=1}^n (Y - \hat{Y})^2 + \sum_{j=1}^p |\beta_j|$
  - Ridge regression estimation:  $\min SSE = \sum_{i=1}^n (Y - \hat{Y})^2 + \sum_{j=1}^p \beta_j^2$
- Error on data      +      Regularization

35/54

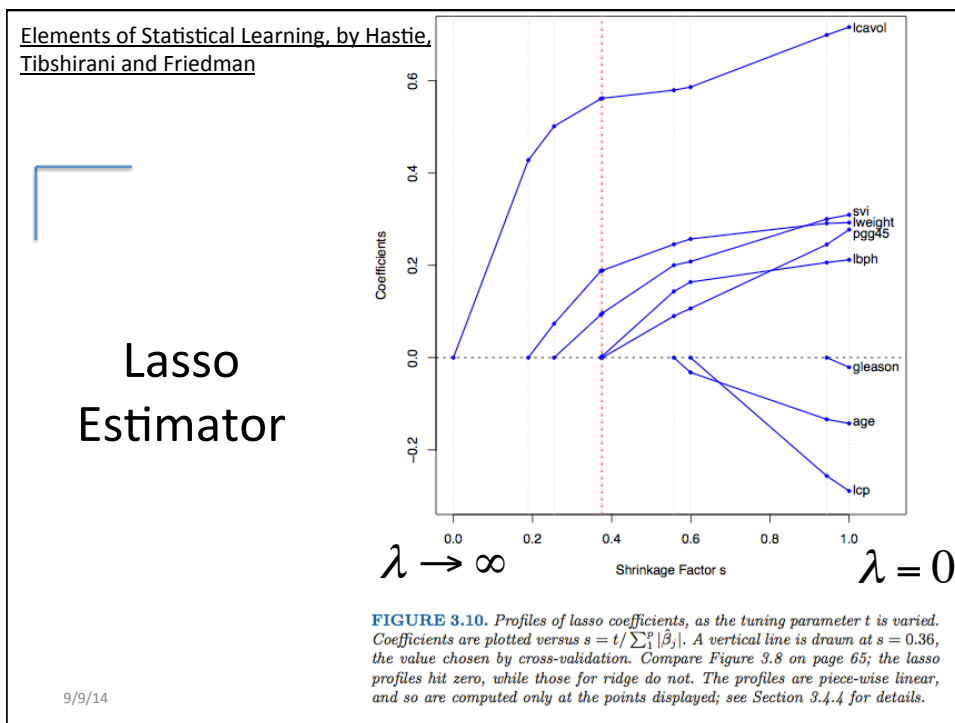
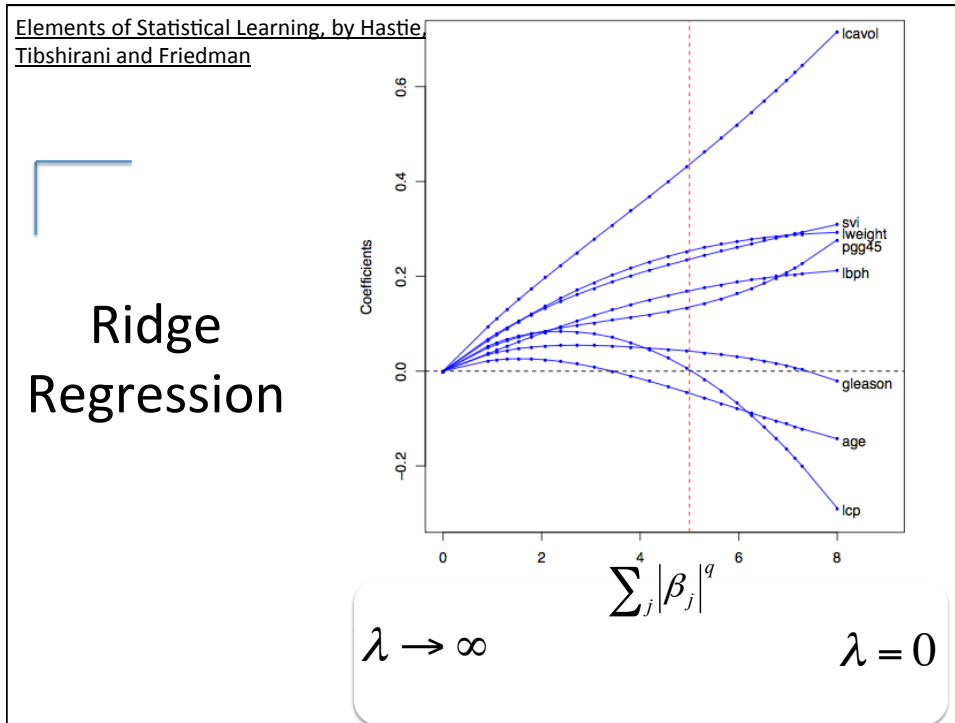
Yanjun Qi / UVA CS 4501-01-6501-07

## Extra

- Not required, though roughly covered during class
  - Subgradient
  - Coordinate descent based learning for Lasso

9/9/14

36



## Today's Recap

- A bit more about Linear Regression Extension
  - Linear regression with predefined RBF basis
  - Locally weighted regression
  
- An Exemplar Application of Regression
  - Text based movie open weekend revenue prediction
  
- Linear Regression Models with Regularization
  - Ridge Regression
  - Lasso

## References

- Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- Elements of Statistical Learning, by Hastie, Tibshirani and Friedman (page 61-69)