# UVA CS 6316
# – Fall 2015 Graduate:
# Machine Learning

# Lecture 6: Linear Regression Model
# with Regularizations

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

9/30/15

1

---

# Where we are ? ➔
# Five major sections of this course

❑ Regression (supervised)

❑ Classification (supervised)

❑ Unsupervised models

❑ Learning theory

❑ Graphical models

9/30/15

2

# Today ➔
## Regression (supervised)

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

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

---

# Today

❑ Linear Regression Model with Regularizations
- ❑ Ridge Regression
- ❑ Lasso Regression
- ❑ Elastic net

# Review: Vector norms

A norm of a vector ||x|| is informally a measure of the "length" of the vector.

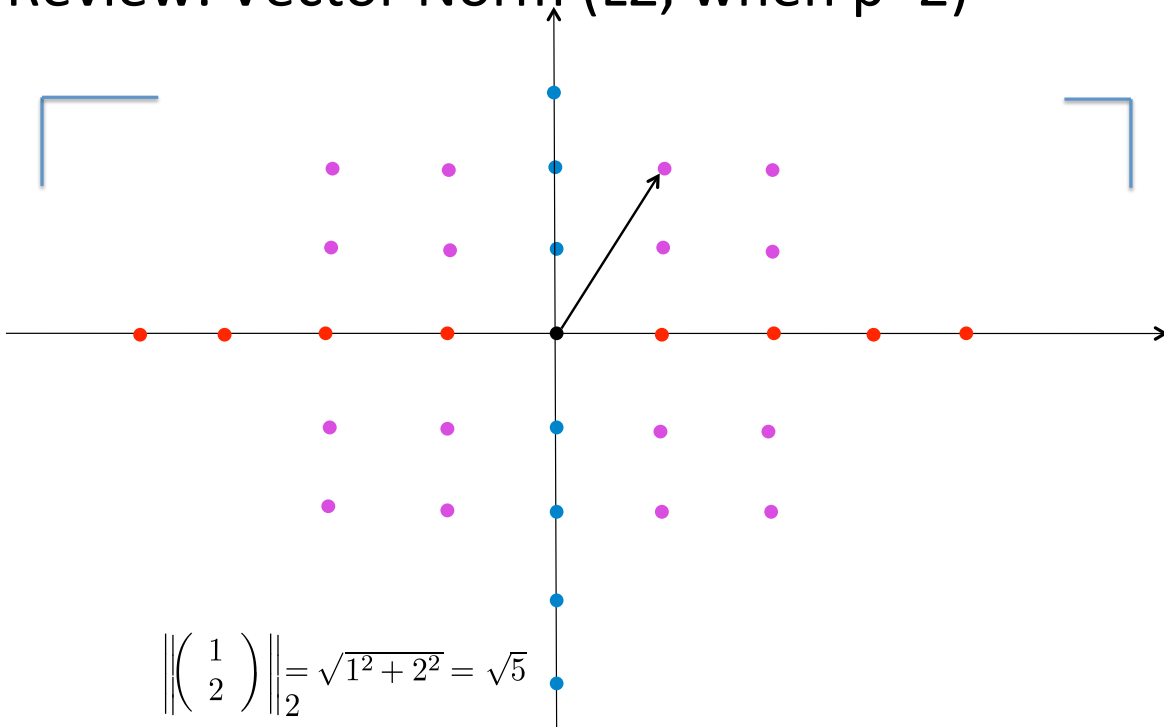$$\|x\|_p = \left( \sum_{i=1}^{n} |x_i|^p \right)^{1/p}$$

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

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

– $L_{infinity}$

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

---

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

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

# Review: Normal equation for LR

- Write the cost function in matrix form:

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

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

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

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

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

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

**The normal equations**

$$\Downarrow$$

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

Assume that $X^T X$ is invertible

9/30/15

7

---

# Comments on the normal equation

$X_{n \times p}$

- In most situations of practical interest, the number of data points $N$ is larger than the dimensionality $p$ of the input space and the matrix $\mathbf{X}$ is of full column rank. If $\nearrow p \times p$ this condition holds, then it is easy to verify that $X^T X$ is necessarily invertible.

$n \gg p$   $\operatorname{rank}(X) \leqslant \min(n, p)$

- The assumption that $X^T X$ is invertible implies that it is positive definite (➔ SSE convex), thus the critical point we have found is a minimum.

- What if $\mathbf{X}$ has less than full column rank? ➔ regularization (later).

9/30/15

8

# Ridge Regression / L2

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

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

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

- The ridge estimator is solution from

HW2

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

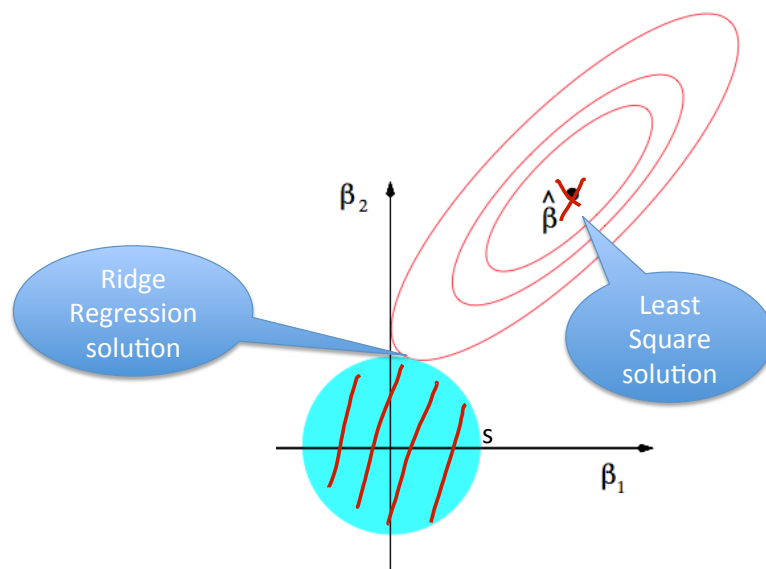to minimize, take derivative and set to zero

- Equivalently

$$\hat{\beta}^{ridge} = \arg\min(y - X\beta)^T(y - X\beta)$$
$$\text{subject to } \sum \beta_j^2 \le s$$

9/30/15

9

---

# Objective Function's Contour lines from Ridge Regression



9/30/15

10

# (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 $\lambda \to \infty$, then $\hat{\beta} \to 0$

---

# Today

❑Linear Regression Model with Regularizations
- ❑ Ridge Regression
- ❑ Lasso Regression
- ❑ Elastic net

# (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

$$\hat{\beta}^{lasso} = \arg\min(y - X\beta)^T(y - X\beta)$$
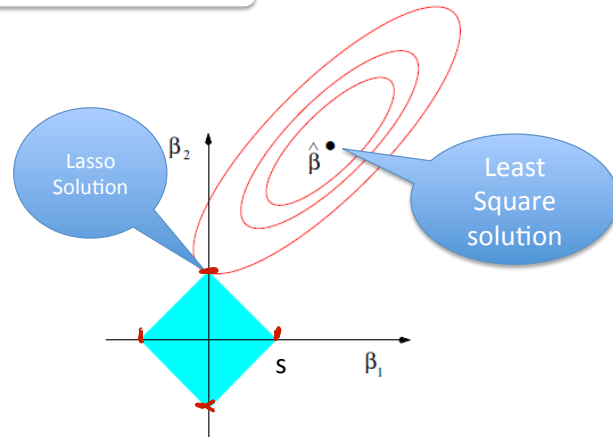
$$\text{subject to } \sum|\beta_j| \le s$$

---

# 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.
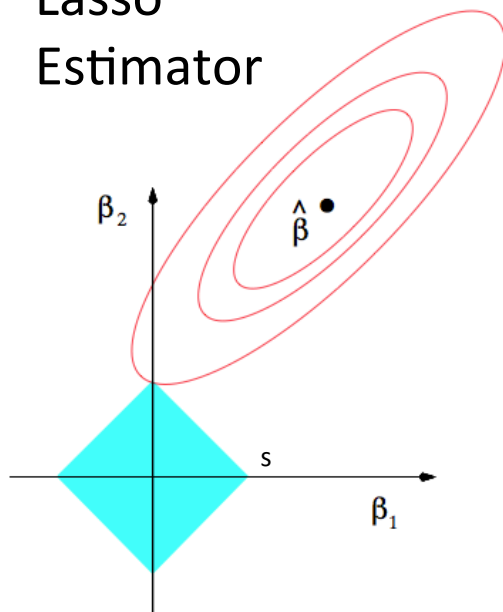
# Lasso (least absolute shrinkage and selection

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

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



Lasso Solution

Least Square solution

9/30/15
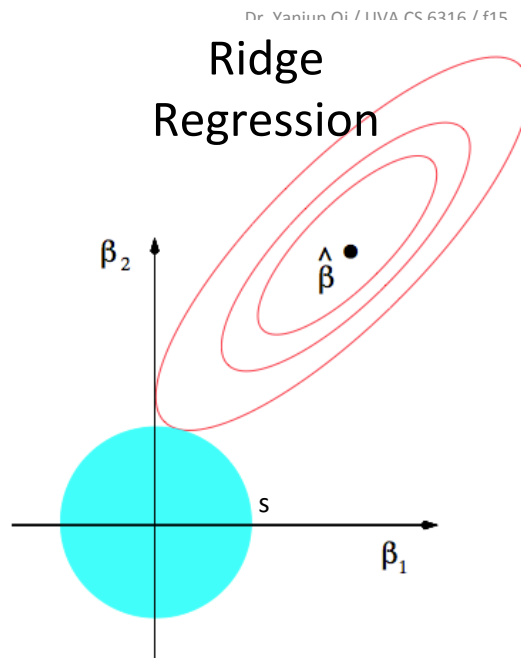
15

---

# 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* $|\beta_1| + |\beta_2| \le t$ *and* $\beta_1^2 + \beta_2^2 \le t^2$, *respectively, while the red ellipses are the contours of the least squares error function.*

# Today

❑Linear Regression Model with Regularizations

   ❑ Ridge Regression

   ❑ Lasso Regression

   ❑ Elastic net

---

# (3) Hybrid of Ridge and Lasso

**Elastic Net regularization**

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

many $\beta_j = 0$

- 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.

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)

$$\hat{\theta} = \underset{\theta=(\beta_0,\beta)}{\text{argmin}} \frac{1}{2n} \sum_{i=1}^{n} \left( y_i - (\beta_0 + x_i^\top \beta) \right)^2 + \lambda P(\beta)$$

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

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.

19

---

# More: A family of shrinkage estimators

$$\beta = \arg\min_\beta \sum_{i=1}^{N} (y_i - x_i^T \beta)^2$$
$$\text{subject to } \sum \left| \beta_j \right|^q \leq s$$

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

Convex

| $q = 4$ | $q = 2$ | $q = 1$ | $q = 0.5$ | $q = 0.1$ |



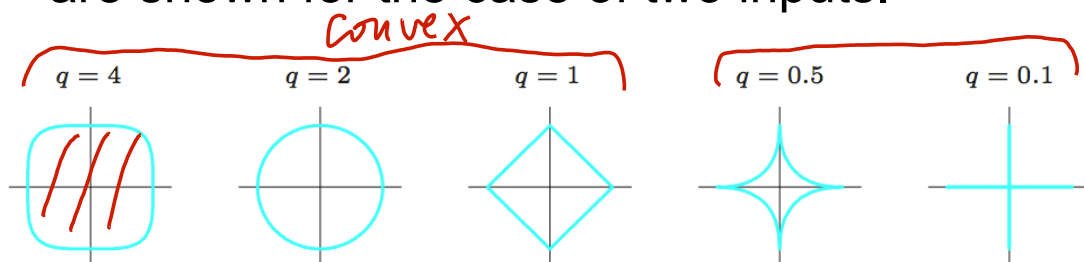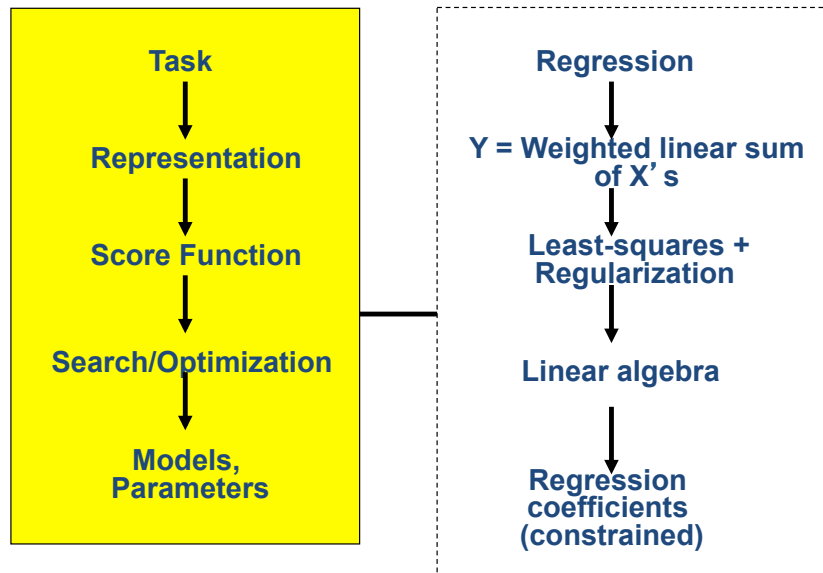**FIGURE 3.12.** *Contours of constant value of $\sum_j |\beta_j|^q$ for given values of q.*

## Regularized multivariate linear regression

Task → Representation → Score Function → Search/Optimization → Models, Parameters

Regression → Y = Weighted linear sum of X's → Least-squares + Regularization → Linear algebra → Regression coefficients (constrained)

---

## Summary:
## Regularized multivariate linear regression

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

- LR estimation:

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

- LASSO estimation:

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

- Ridge regression estimation:

$$\min SSE = \sum_{i=1}^{n} \left( Y - \hat{Y} \right)^2 + \sum_{j=1}^{p} \beta_j^2$$
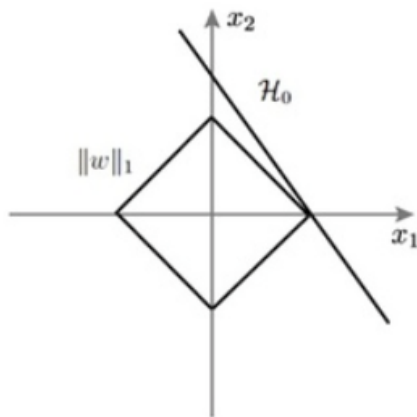
Error on data    +    Regularization
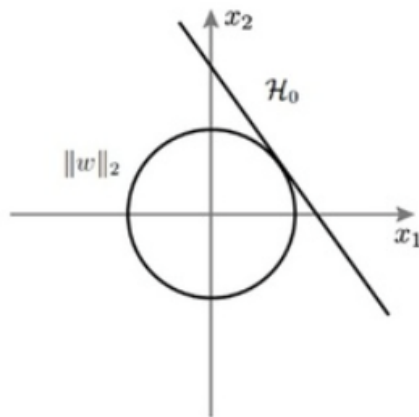
# EXTRA (NOT REQUIRED)

---

# Today

❑ Linear Regression Model with Regularizations
- ❑ Ridge Regression
- ❑ Lasso Regression
    - ❑ Extra: how to perform training
- ❑ Elastic net
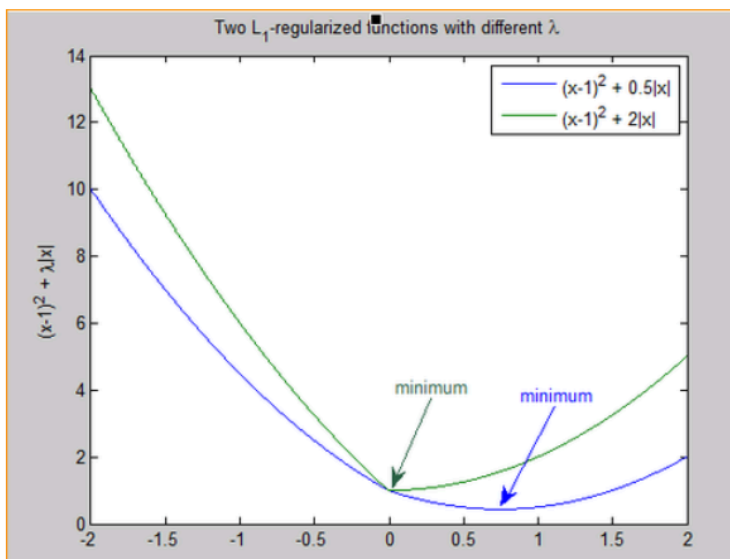
**A**   L1 regularization



**B**   L2 regularization



due to the nature of L_1 norm, the viable solutions are limited to the corners, which are on one axis only - in the above case x1. Value of x2 = 0. This means that the solution has eliminated the role of x2 leading to sparsity
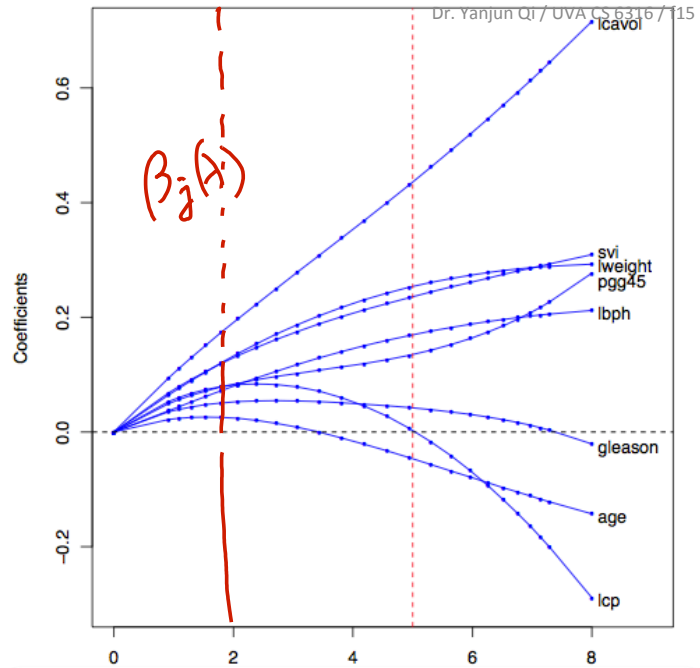
9/30/15

25

---

$L_1$-regularized loss function $F(x) = f(x) + \lambda\|x\|_1$ is non-smooth. It's not differentiable at 0. Optimization theory says that the optimum of a function is either the point with 0-derivative or one of the irregularities (corners, kinks, etc.). So, it's possible that the optimal point of $F$ is 0 even if 0 isn't the stationary point of $f$. In fact, it would be 0 if $\lambda$ is large enough (stronger regularization effect). Below is a graphical illustration.

http://www.quora.com/What-is-the-difference-between-L1-and-L2-regularization



26

# Regularization path of a Ridge Regression



$\lambda \to \infty$

$\lambda = 0$

9/30/15

# Regularization path of a Lasso Estimator



$\lambda \to \infty$

Shrinkage Factor s

$\lambda = 0$

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.*

9/30/15

# How to Learn Parameter for Lasso

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

$$\text{subject to } \sum \left| \beta_j \right| \leq s$$

- $\ell_1$-**norm is non differentiable!**

  – cannot compute the gradient of the absolute value

  $\Rightarrow$ **Directional derivatives** (or subgradient)

---

$$RSS\_loss\,(\lambda) = (Y - X\beta)^T(Y - X\beta) + \lambda \sum_{\partial=1}^{P} |\beta_\partial|$$

$$= \sum_{i=1}^{n} (Y_i - x_i^T\beta)^2 + \lambda \sum_{\partial=1}^{P} |\beta_\partial|$$

$$= \left[ \sum_{i=1}^{n} (Y_i - x_{ij}\beta_\partial - x_{i\{-\partial\}}^T \beta_{-\partial})^{②} \right] + \underbrace{\lambda \sum_{j=1}^{P} |\beta_j|}$$
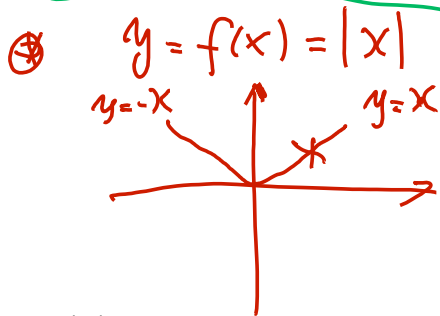
if $\beta = (\beta_1, \beta_2, \beta_3)$

$\Rightarrow \beta_{-2} = (\beta_1, \beta_3)$

$\Rightarrow \dfrac{\partial \ell}{\partial \beta_\partial} = \sum_{i=1}^{n} \underbrace{2(Y_i - x_{ij}\beta_j - x_{i\{-\partial\}}^T \beta_{-\partial})} \underbrace{(-x_{ij})}$

$$+ \lambda \dfrac{\partial}{\partial \beta_\partial} |\beta_\partial|$$

$$= 2 \sum_{i=1}^{n} x_{ij}^2 \beta_j - 2 \sum_{i=1}^{n} \left( y_i - x_{i\{-j\}}^T \beta_{\{-j\}} \right) x_{ij}$$

$$\underbrace{\phantom{2 \sum_{i=1}^{n} x_{ij}^2 \beta_j}}_{a_j} \qquad \underbrace{+ \lambda \frac{\partial}{\partial \beta_j} |\beta_j|}_{C_j}$$

$$= a_j \beta_j - C_j + \lambda \frac{\partial}{\partial \beta_j} |\beta_j| \underset{\text{set to}}{=\!=} 0$$

convext $\Rightarrow$ unique

④ $y = f(x) = |x|$

$y = -x$ ↗ $y = x$

$x$

$$\partial f(x) = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \end{cases}$$

---

$$\frac{\partial \ell}{\partial \beta_j} = \begin{cases} a_j \beta_j - C_j - \lambda = 0, & \text{if } \beta_j < 0 \\ a_j \beta_j - C_j + \lambda, & \text{if } \beta_j > 0 \\ [a_j \beta_j - C_j - \lambda, \ a_j \beta_j - C_j + \lambda], & \text{if } \beta_j = 0 \end{cases}$$

Set to 0

$$\hat{\beta}_j = \begin{cases} \dfrac{C_j + \lambda}{a_j}, & \text{if } C_j + \lambda < 0 \Rightarrow C_j < -\lambda \\[2mm] \dfrac{C_j - \lambda}{a_j}, & \text{if } C_j > \lambda \\[2mm] 0, & \text{if } -\lambda \leqslant C_j \leqslant \lambda \end{cases}$$

Soft thresholding

# Coordinate descent based Learning of Lasso

1. Initialize $\beta$

2. Repeat until converged

3. For $j = 1, 2, \ldots, P$ do

$$a_j = 2 \sum_{i=1}^{n} x_{ij}^2$$

$$c_j = 2 \sum_{i=1}^{n} x_{ij} \left( y_i - x_i^T \beta + x_{ij} \beta_j \right)$$

if $c_j < -\lambda$

$$\beta_j = (c_j + \lambda)/a_j$$

else if, $c_j > \lambda$

$$\beta_j = (c_j - \lambda)/a_j$$
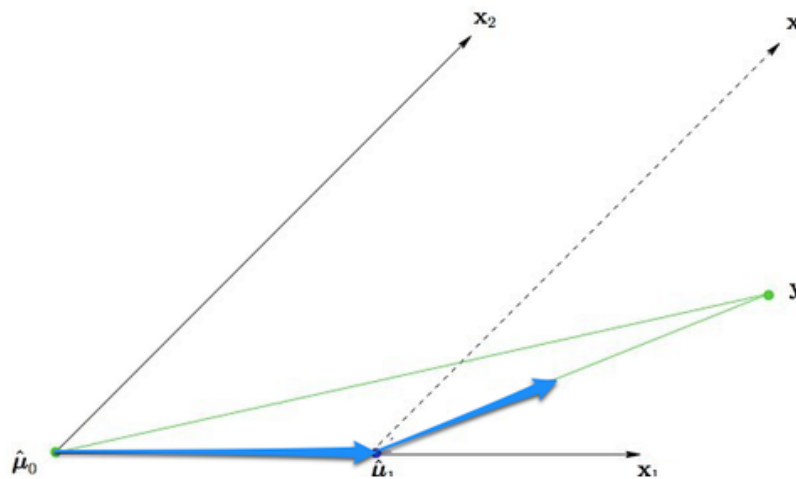
else,     soft-thresholding

$$\beta_j = 0$$

33

> Coordinate descent (WIKI)➔ one does line search along one coordinate direction at the current point in each iteration. One uses different coordinate directions cyclically throughout the procedure.

---

# LARS: Least Angle Regression (state-of-art LASSO solver algorithm)

# Lasso 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

---

# Today

❑Linear Regression Model with Regularizations
    ❑ Ridge Regression
    ❑ Lasso Regression
       ❑ Extra: how to perform training
    ❑ Elastic net
       ❑ Extra: how to perform training

# (3) Elastic Net:
# 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.

9/30/15

37

# Naïve elastic net

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

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

data error    L2    L1

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

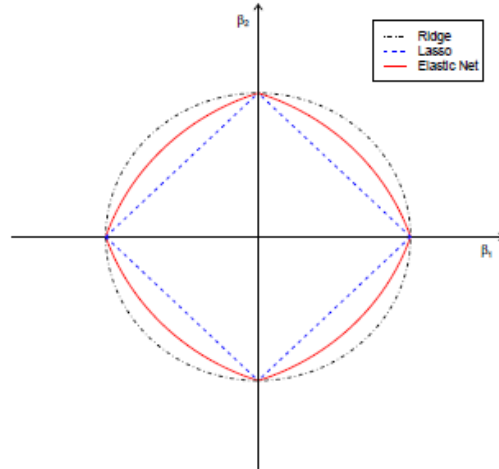- The naive elastic net estimator is the minimizer of equation

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

- Let $\quad \alpha = \lambda_2/(\lambda_1 + \lambda_2)$

$$\hat{\beta} = \arg\min_{\beta} |\mathbf{y} - \mathbf{X}\beta|^2, \qquad \text{subject to } (1-\alpha)|\beta|_1 + \alpha|\beta|^2 \leq t \text{ for some } t.$$

# Geometry of elastic net

2-dimensional illustration $\alpha = 0.5$



# Connecting LASSO and Elastic net

- Lemma: Given $(\lambda_1, \lambda_2)$, define an artificial data set $(y^*, X^*)$

$$\mathbf{X}^*_{(n+p)\times p} = (1+\lambda_2)^{-1/2}\begin{pmatrix}\mathbf{X}\\\sqrt{\lambda_2}\mathbf{I}\end{pmatrix}, \qquad \mathbf{y}^*_{(n+p)} = \begin{pmatrix}\mathbf{y}\\0\end{pmatrix}.$$

*(handwritten annotations: n×p, n×1, (n+p)×p, (n+p)×1)*

Let $\gamma = \lambda_1/\sqrt{(1+\lambda_2)}$ and $\beta^* = \sqrt{(1+\lambda_2)}\beta$. Then the naïve elastic net criterion can be written as

$$L(\gamma,\beta) = L(\gamma,\beta^*) = \left|\mathbf{y}^* - \mathbf{X}^*\beta^*\right|^2 + \gamma\left|\beta^*\right|_1. \;\Rightarrow\; \beta^*$$

- Let,

$$\hat{\beta}^* = \arg\min_{\beta^*} L\{(\gamma,\beta^*)\};$$

- Then

elastic $\quad \hat{\beta} = \dfrac{1}{\sqrt{(1+\lambda_2)}}\hat{\beta}^*$ lasso augmented

*(handwritten: $X\beta = y$, $n\times p$, $p\times 1$, $n\times 1$, $(n+p)\times p$, $p\times 1$, $(n+p)\times 1$)*

# Advantage of Elastic net

$p \gg n$

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

$\Rightarrow X \; n \times p$

- In the augmented formulation, $\Rightarrow X^*$
  $(n+p) \times p$
  - sample size n+p and X* has rank p
  - ➔ can potentially select all the predictors

- Naïve elastic net can perform automatic variable selection like lasso

# Grouping Effect of Naïve Elastic net

*Theorem 1.* Given data $(\mathbf{y}, \mathbf{X})$ and parameters $(\lambda_1, \lambda_2)$, the response $\mathbf{y}$ is centred and the predictors $\mathbf{X}$ are standardized. Let $\hat{\beta}(\lambda_1, \lambda_2)$ be the naïve elastic net estimate. Suppose that $\hat{\beta}_i(\lambda_1, \lambda_2)\, \hat{\beta}_j(\lambda_1, \lambda_2) > 0$. Define

$$D_{\lambda_1,\lambda_2}(i,j) = \frac{1}{|\mathbf{y}|_1}|\hat{\beta}_i(\lambda_1,\lambda_2) - \hat{\beta}_j(\lambda_1,\lambda_2)|;$$

then

$$D_{\lambda_1,\lambda_2}(i,j) \leqslant \frac{1}{\lambda_2}\sqrt{\{2(1-\rho)\}},$$

where $\rho = \mathbf{x}_i^T \mathbf{x}_j$, the sample correlation.

- D is the difference between the coefficient paths of predictors i and j.
- If $x_i$ and $x_j$ are high correlated $\rho=1$, this theorem provides a quantitative description for the grouping effect of Naive Elastic Net.

# Elastic Net:
# Re-scaling of Naive Elastic Net

- **Deficiency of the Naive Elastic Net:** Empirical evidence shows the Naive Elastic Net does not perform satisfactorily. The reason is that there are two shrinkage procedures (Ridge and LASSO) in it. Double shrinkage introduces unnecessary bias.

- Re-scaling of Naive Elastic Net gives better performance, yielding the Elastic Net solution:

$$\hat{\beta}(\text{ENet}) = (1 + \lambda_2) \cdot \hat{\beta}(\text{Naive ENet})$$

- Reason: Undo shrinkage.

---

*Theorem 2.* Given data $(\mathbf{y}, \mathbf{X})$ and $(\lambda_1, \lambda_2)$, then the elastic net estimates $\hat{\beta}$ are given by

$$\hat{\beta} = \arg\min_{\beta} \beta^{\mathrm{T}} \left( \frac{\mathbf{X}^{\mathrm{T}}\mathbf{X} + \lambda_2 \mathbf{I}}{1 + \lambda_2} \right) \beta - 2\mathbf{y}^{\mathrm{T}}\mathbf{X}\beta + \lambda_1 |\beta|_1. \tag{14}$$

It is easy to see that

$$\hat{\beta}(\text{lasso}) = \arg\min_{\beta} \beta^{\mathrm{T}} (\mathbf{X}^{\mathrm{T}}\mathbf{X}) \beta - 2\mathbf{y}^{\mathrm{T}}\mathbf{X}\beta + \lambda_1 |\beta|_1. \tag{15}$$

Hence theorem 2 interprets the elastic net as a stabilized version of the lasso. Note that $\hat{\Sigma} = \mathbf{X}^{\mathrm{T}}\mathbf{X}$ is a sample version of the correlation matrix $\Sigma$ and

$$\frac{\mathbf{X}^{\mathrm{T}}\mathbf{X} + \lambda_2 \mathbf{I}}{1 + \lambda_2} = (1 - \gamma)\hat{\Sigma} + \gamma \mathbf{I}$$

- Rescaling after the elastic net penalization is mathematically equivalent to replacing $\Sigma$ with its shrunken version in the lasso.

# Computation of Elastic Net

- First solve the Naive Elastic Net problem, then rescale it.
- For fixed $\lambda_2$, the Naive Elastic Net problem is equivalent to a LASSO problem, with a huge design matrix if p >> n
- LASSO already has an efficient solver called LARS (Least Angle Regression).

# Today Recap

❑Linear Regression Model with Regularizations
- ❑ Ridge Regression
- ❑ Lasso Regression
  - ❑ Extra: how to perform training
- ❑ Elastic net
  - ❑ Extra: how to perform training

# Extra: Shrinkage Bias Term ?

- If the data is not centered, there exists bias term
    - http://stats.stackexchange.com/questions/86991/reason-for-not-shrinking-the-bias-intercept-term-in-regression

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

- We normally assume we centered x and y. If this is true, no need to have bias term, e.g., for lasso,

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

# 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*