

ECE595 / STAT598: Machine Learning I

Lecture 2.2: Regularization - LASSO Regression

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline

Mathematical Background

- Lecture 1: Linear regression: A basic data analytic tool
- **Lecture 2: Regularization: Constraining the solution**
- Lecture 3: Kernel Method: Enabling nonlinearity

Lecture 2: Regularization

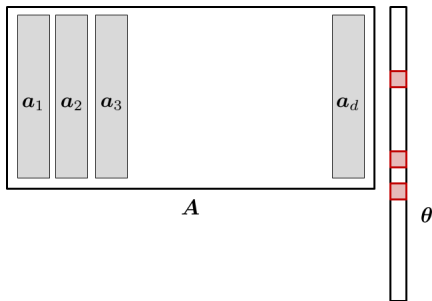
- Ridge Regression
 - Regularization
 - Parameter
- **LASSO Regression**
 - **Sparsity**
 - **Algorithm**
 - **Application**

LASSO Regression

- An alternative to the Ridge Regression is **Least Absolute Shrinkage and Selection Operator (LASSO)**
- The loss function is

$$J(\theta) = \|\mathbf{A}\theta - \mathbf{y}\|^2 + \lambda\|\theta\|_1$$

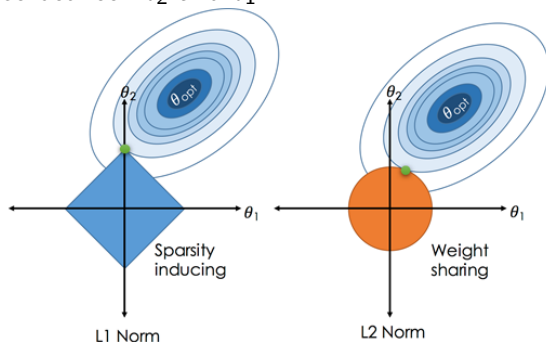
- Intuition behind LASSO: Many features are not active.



Interpreting the LASSO Solution

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \quad \|\mathbf{A}\theta - \mathbf{y}\|^2 + \lambda\|\theta\|_1$$

- $\|\theta\|_1$ promotes sparsity of θ . It is the nearest convex approximation to $\|\theta\|_0$, which is the number of non-zeros.
- The difference between ℓ_2 and ℓ_1 ¹:



¹Figure source: <http://www.ds100.org/>

Why are Sparse Models Useful?



non-zeros = 33.51%



13.58%



1.21%

- Images are sparse in transform domains, e.g., Fourier and wavelet.
- Intuition: There are more low frequency components and less high frequency components.
- Examples above: \mathbf{A} is the wavelet basis matrix. θ are the wavelet coefficients.
- We can truncate the wavelet coefficients and retain a good image.
- Many image compression schemes are based on this, e.g., JPEG, JPEG2000.

LASSO for Image Reconstruction

Image inpainting via KSVD dictionary-learning ²



- \mathbf{y} = image with missing pixels. \mathbf{A} = a matrix storing a set of trained feature vectors (called dictionary atoms). $\boldsymbol{\theta}$ = coefficients.
- minimize $\|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|^2 + \lambda\|\boldsymbol{\theta}\|_1$.
- KSVD = k-means + Singular Value Decomposition (SVD): A method to train the feature vectors that demonstrate sparse representations.

²Figure is taken from Mairal, Elad, Sapiro, IEEE T-IP 2008

Shrinkage Operator

The LASSO problem can be solved using a shrinkage operator. Consider a simplified problem (with $\mathbf{A} = \mathbf{I}$)

$$\begin{aligned} J(\boldsymbol{\theta}) &= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\theta}\|^2 + \lambda \|\boldsymbol{\theta}\|_1 \\ &= \sum_{j=1}^d \left\{ \frac{1}{2} (y_j - \theta_j)^2 + \lambda |\theta_j| \right\} \end{aligned}$$

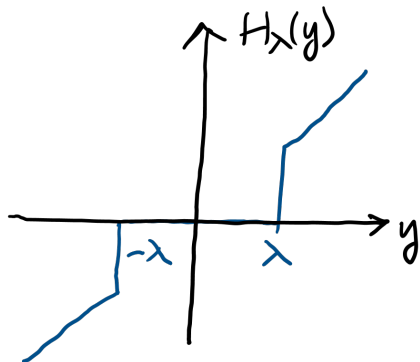
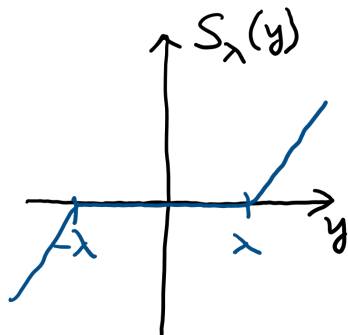
Since the loss is **separable**, the optimization is solved when each individual term is minimized. The individual problem

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \left\{ \frac{1}{2} (y - \theta)^2 + \lambda |\theta| \right\} \\ &= \max(|y| - \lambda, 0) \operatorname{sign}(y) \\ &\stackrel{\text{def}}{=} \mathcal{S}_\lambda(y). \end{aligned}$$

Proof: See Appendix.

Shrinkage VS Hard Threshold

- The shrinkage operator looks as follows.
- Any number between $[-\lambda, \lambda]$ is “shrink” to zero.
- Try compare with the hard threshold operator $\mathcal{H}_\lambda(y) = y \cdot \mathbf{1}\{|y| \geq \lambda\}$



Algorithms to Solve LASSO Regression

In general, the LASSO problem requires iterative algorithms:

- ISTA Algorithm (Daubechies et al. 2004)
 - For $k = 1, 2, \dots$
 - $\mathbf{v}^k = \boldsymbol{\theta}^k - 2\gamma \mathbf{A}^T (\mathbf{A}\boldsymbol{\theta}^k - \mathbf{y})$.
 - $\boldsymbol{\theta}^{k+1} = \max(|\mathbf{v}^k| - \lambda, 0) \text{sign}(\mathbf{v}^k)$.
- FISTA Algorithm (Beck-Teboulle 2008)
 - For $k = 1, 2, \dots$
 - $\mathbf{v}^k = \boldsymbol{\theta}^k - 2\gamma \mathbf{A}^T (\mathbf{A}\boldsymbol{\theta}^k - \mathbf{y})$.
 - $\mathbf{z}^k = \max(|\mathbf{v}^k| - \lambda, 0) \text{sign}(\mathbf{v}^k)$.
 - $\boldsymbol{\theta}^{k+1} = \alpha_k \boldsymbol{\theta}^k + (1 - \alpha_k) \mathbf{z}^k$.
- ADMM Algorithm (Eckstein-Bertsekas 1992, Boyd et al. 2011)
 - For $k = 1, 2, \dots$
 - $\boldsymbol{\theta}^{k+1} = (\mathbf{A}^T \mathbf{A} + \rho \mathbf{I})^{-1} (\mathbf{A}^T \mathbf{y} + \rho \mathbf{z}^k - \mathbf{u}^k)$
 - $\mathbf{z}^{k+1} = \max(|\boldsymbol{\theta}^{k+1} + \mathbf{u}^k / \rho| - \lambda / \rho, 0) \text{sign}(\boldsymbol{\theta}^{k+1} + \mathbf{u}^k / \rho)$
 - $\mathbf{u}^{k+1} = \mathbf{u}^k + \rho(\boldsymbol{\theta}^{k+1} - \mathbf{z}^{k+1})$
- And many others.

Example: Crime Rate Data

| city | funding | hs | not-hs | college | college4 | crime rate |
|----------|----------|----------|----------|----------|----------|------------|
| 1 | 40 | 74 | 11 | 31 | 20 | 478 |
| 2 | 32 | 72 | 11 | 43 | 18 | 494 |
| 3 | 57 | 70 | 18 | 16 | 16 | 643 |
| 4 | 31 | 71 | 11 | 25 | 19 | 341 |
| 5 | 67 | 72 | 9 | 29 | 24 | 773 |
| \vdots | \vdots | \vdots | \vdots | \vdots | | |
| 50 | 66 | 67 | 26 | 18 | 16 | 940 |

<https://web.stanford.edu/~hastie/StatLearnSparsity/data.html>

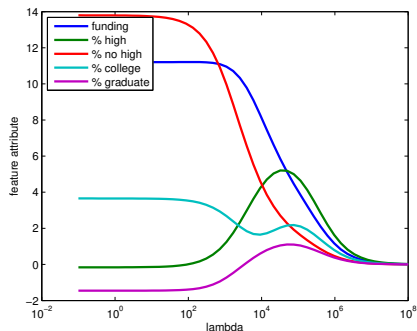
Consider the following two optimizations

$$\hat{\theta}_1(\lambda) = \operatorname{argmin}_{\theta} J_1(\theta) \stackrel{\text{def}}{=} \|\mathbf{A}\theta - \mathbf{y}\|^2 + \lambda\|\theta\|_1,$$

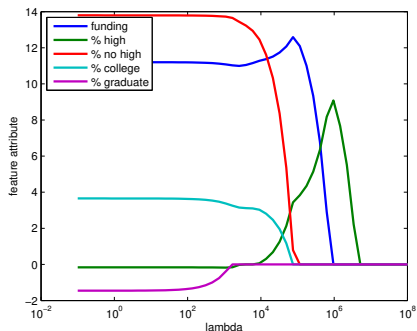
$$\hat{\theta}_2(\lambda) = \operatorname{argmin}_{\theta} J_2(\theta) \stackrel{\text{def}}{=} \|\mathbf{A}\theta - \mathbf{y}\|^2 + \lambda\|\theta\|^2.$$

Comparison between ℓ_1 and ℓ_2 norm

- Plot $\hat{\theta}_1(\lambda)$ and $\hat{\theta}_2(\lambda)$ vs. λ .
- LASSO tells us which factor appears first.
- If we are allowed to use only one feature, then % high is the one.
- Two features, then % high + funding.



Ridge



LASSO

Pros and Cons

Ridge Regression

- (+) Analytic solution, because the loss function is differentiable.
- (+) As such, a lot of well-established theoretical guarantees.
- (+) Algorithm is simple, just one equation.
- (-) Limited interpretability, since the solution is usually a dense vector.
- (-) Does not reflect the nature of certain problems, e.g., sparsity.

LASSO

- (+) Proven applications in many domains, e.g., images and speeches.
- (+) Echoes particularly well in modern deep learning where parameter space is huge.
- (+) Increasing number of theoretical guarantees for special matrices.
- (+) Algorithms are available.
- (-) No closed-form solution. Algorithms are iterative.

Reading List

Ridge Regression

- Stanford CS 229 Note on Linear Algebra
<http://cs229.stanford.edu/section/cs229-linalg.pdf>
- Lecture Note on Ridge Regression
<https://arxiv.org/pdf/1509.09169.pdf>
- Theobald, C. M. (1974). Generalizations of mean square error applied to ridge regression. *Journal of the Royal Statistical Society. Series B (Methodological)*, 36(1), 103-106.

LASSO Regression

- ECE/STAT 695 (Lecture 1)
<https://engineering.purdue.edu/ChanGroup/ECE695.html>
- Statistical Learning with Sparsity (Chapter 2)
<https://web.stanford.edu/~hastie/StatLearnSparsity/>
- Elements of Statistical Learning (Chapter 3.4)
<https://web.stanford.edu/~hastie/ElemStatLearn/>