

ECE595 / STAT598: Machine Learning I

Lecture 02: Regularized Linear Regression

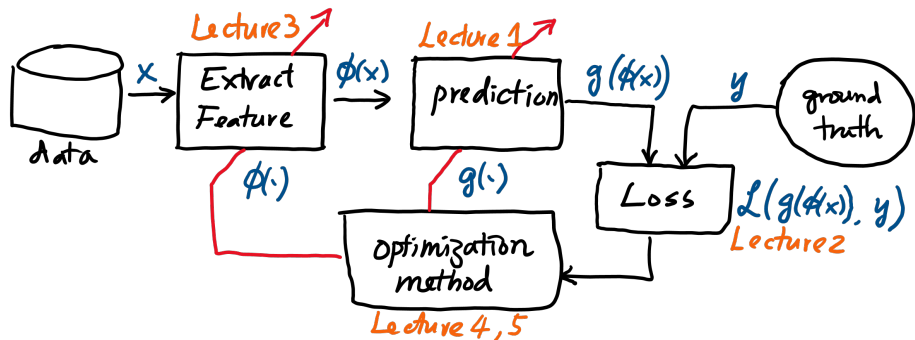
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Outline



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

Ridge Regression

- Applies to both over and under determined systems.
- The loss function of the ridge regression is defined as

$$J(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2 + \lambda\|\boldsymbol{\theta}\|^2$$

- $\|\boldsymbol{\theta}\|^2$ Regularization function
- λ : Regularization parameter
- The solution of the ridge regression is

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}} \left\{ \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2 + \lambda\|\boldsymbol{\theta}\|^2 \right\} \\ &= 2\mathbf{A}^T(\mathbf{A}\boldsymbol{\theta} - \mathbf{y}) + 2\lambda\boldsymbol{\theta} = \mathbf{0},\end{aligned}$$

which gives us $\hat{\boldsymbol{\theta}} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$.

- Probabilistic interpretation: See Appendix.

Change in Eigen-values

Ridge regression improves the eigen-values:

- Eigen-decomposition of $\mathbf{A}^T \mathbf{A}$:

$$\mathbf{A}^T \mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{U}^T \succeq 0,$$

where \mathbf{U} = eigen-vector matrix, \mathbf{S} = eigen-value matrix.

- \mathbf{S} is a diagonal matrix with non-negative entries:

$$\mathbf{S} = \begin{bmatrix} \clubsuit & & & \\ & \clubsuit & & \\ & & \clubsuit & \\ & & & 0 \end{bmatrix}$$

See Tutorial on “Linear Algebra”.

- Therefore, $\mathbf{S} + \lambda \mathbf{I}$ is always positive for any $\lambda > 0$, implying that

$$\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I} = \mathbf{U} (\mathbf{S} + \lambda \mathbf{I}) \mathbf{U}^T \succ 0.$$

Regularization Parameter λ

- The solution of the ridge regression is

$$\hat{\theta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$$

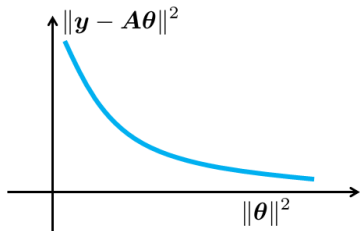
- If $\lambda \rightarrow 0$, then $\hat{\theta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$:

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

- If $\lambda \rightarrow \infty$, then $\hat{\theta} = \mathbf{0}$:

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

- There is a trade-off curve between the two terms by varying λ .



Comparing Vanilla and Ridge

Suppose $\mathbf{y} = \mathbf{A}\boldsymbol{\theta}^* + \mathbf{e}$ for some ground truth $\boldsymbol{\theta}^*$ and noise vector \mathbf{e} . Then, the **vanilla linear regression** will give us

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} \\ &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{A}\boldsymbol{\theta}^* + \mathbf{e}) \\ &= \boldsymbol{\theta}^* + (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{e}\end{aligned}$$

If \mathbf{e} has zero mean and variance σ^2 , we can show that

$$\begin{aligned}\mathbb{E}[\hat{\boldsymbol{\theta}}] &= \boldsymbol{\theta}^*, \\ \text{Cov}[\hat{\boldsymbol{\theta}}] &= \sigma^2 (\mathbf{A}^T \mathbf{A})^{-1}.\end{aligned}$$

Therefore, the regression coefficients are unbiased but have large variance. We can further show that the mean-squared error (MSE) is

$$\text{MSE}(\hat{\boldsymbol{\theta}}) = \sigma^2 \text{Tr}\{(\mathbf{A}^T \mathbf{A})^{-1}\}.$$

Comparing Vanilla and Ridge

On the other hand, if we use ridge regression, then

$$\begin{aligned}\hat{\theta}(\lambda) &= (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T (\mathbf{A} \theta^* + \mathbf{e}) \\ &= (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{A} \theta^* + (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{e}.\end{aligned}$$

Again, if \mathbf{e} is zero mean and has a variance σ^2 , then (See Reading List)

$$\begin{aligned}\mathbb{E}[\hat{\theta}(\lambda)] &= (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{A} \theta^* \\ \text{Cov}[\hat{\theta}(\lambda)] &= \sigma^2 (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{A} (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \\ \text{MSE}[\hat{\theta}(\lambda)] &= \sigma^2 \text{Tr}\{\mathbf{W}_\lambda (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{W}_\lambda^T\} + \theta^{*T} (\mathbf{W}_\lambda - \mathbf{I})^T (\mathbf{W}_\lambda - \mathbf{I}) \theta^*,\end{aligned}$$

where $\mathbf{W}_\lambda \stackrel{\text{def}}{=} (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{A}$. In particular, we can show that

[Theorem \(Theobald 1974\)](#)

For $\lambda < 2\sigma^2 \|\theta^\|^{-2}$, it holds that $\text{MSE}(\hat{\theta}(\lambda)) < \text{MSE}(\hat{\theta})$.*

Geometric Interpretation

The following three problems are equivalent

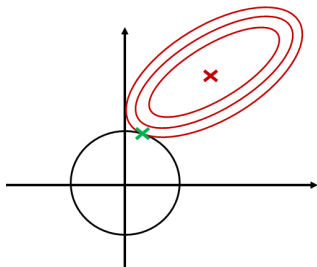
$$\theta_{\lambda}^* = \underset{\theta}{\operatorname{argmin}} \quad \|\mathbf{A}\theta - \mathbf{y}\|^2 + \lambda\|\theta\|^2$$

$$\theta_{\alpha}^* = \underset{\theta}{\operatorname{argmin}} \quad \|\mathbf{A}\theta - \mathbf{y}\|^2 \quad \text{subject to } \|\theta\|^2 \leq \alpha$$

$$\theta_{\epsilon}^* = \underset{\theta}{\operatorname{argmin}} \quad \|\theta\|^2 \quad \text{subject to } \|\mathbf{A}\theta - \mathbf{y}\|^2 \leq \epsilon$$

under an appropriately chosen tuple $(\lambda, \alpha, \epsilon)$.

- Larger λ = Smaller α
- θ^* 's magnitude is tighter bounded



Choosing λ

Because the following three problems are equivalent

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

$$\theta_{\alpha}^* = \underset{\theta}{\operatorname{argmin}} \quad \|\mathbf{A}\theta - \mathbf{y}\|^2 \quad \text{subject to} \quad \|\theta\|^2 \leq \alpha$$

$$\theta_{\epsilon}^* = \underset{\theta}{\operatorname{argmin}} \quad \|\theta\|^2 \quad \text{subject to} \quad \|\mathbf{A}\theta - \mathbf{y}\|^2 \leq \epsilon$$

- We can seek λ that satisfies $\|\theta\|^2 \leq \alpha$:
 - You know how much $\|\theta\|^2$ would be appropriate.
- We can seek λ that satisfies $\|\mathbf{A}\theta - \mathbf{y}\|^2 \leq \epsilon$
 - You know how much $\|\mathbf{A}\theta - \mathbf{y}\|^2$ would be tolerable.
- Other approaches:
 - Akaike's information criterion: Balance model fit with complexity
 - Cross validation: Leave one out
 - Generalized cross-validation: Cross-validation + weight

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

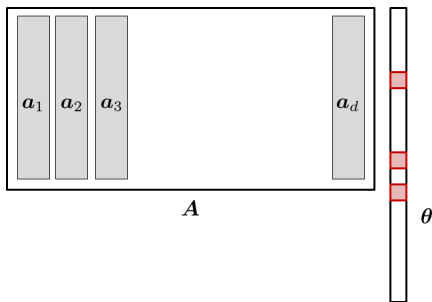
- Ridge Regression
 - Regularization
 - Parameter
- **LASSO Regression**
 - **Sparsity**
 - **Algorithm**
 - **Application**

LASSO Regression

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

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

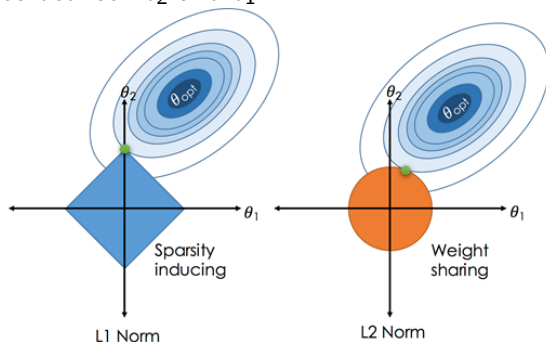
- Intuition behind LASSO: Many features are not active.



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

<https://ieeexplore.ieee.org/document/4392496>

Shrinkage Operator

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

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

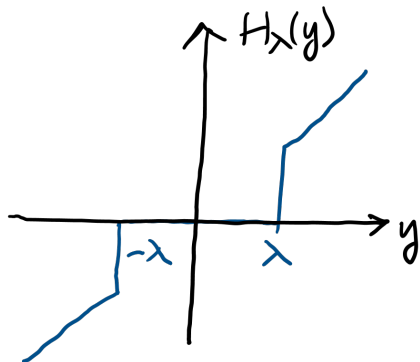
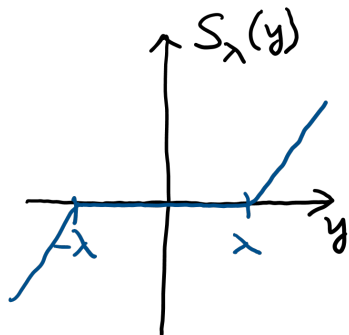
Since the loss is **separable**, the optimization is solved when each individual term is minimized. The individual problem

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

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>

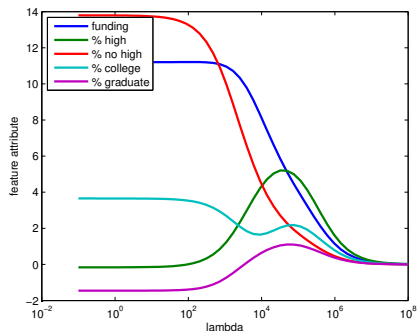
Consider the following two optimizations

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

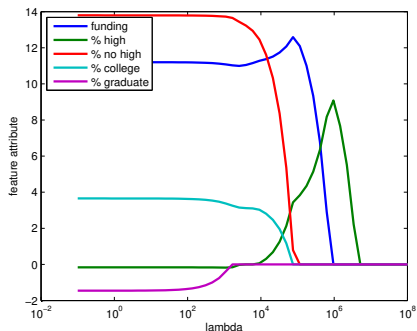
$$\hat{\theta}_2(\lambda) = \underset{\theta}{\operatorname{argmin}} J_2(\theta) \stackrel{\text{def}}{=} \|\mathbf{A}\theta - \mathbf{y}\|^2 + \lambda\|\theta\|^2.$$

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/>

Appendix

Treating Linear Regression as Maximum-Likelihood

- Minimizing $J(\theta)$ is the same as solving a **maximum-likelihood**:

$$\begin{aligned}\theta^* &= \underset{\theta}{\operatorname{argmin}} \quad \|\mathbf{A}\theta - \mathbf{y}\|^2 \\ &= \underset{\theta}{\operatorname{argmin}} \quad \sum_{n=1}^N (\theta^T \mathbf{x}^n - y^n)^2 \\ &= \underset{\theta}{\operatorname{argmax}} \quad \exp \left\{ - \sum_{n=1}^N (\theta^T \mathbf{x}^n - y^n)^2 \right\} \\ &= \underset{\theta}{\operatorname{argmax}} \quad \prod_{n=1}^N \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ - \frac{(\theta^T \mathbf{x}^n - y^n)^2}{2\sigma^2} \right\} \right\}\end{aligned}$$

- Assume noise is i.i.d. Gaussian with variance σ^2 .
- See Tutorial on Probability

Likelihood Function

- **Likelihood:**

$$p_{X|\Theta}(x|\theta) = \text{probability density of } x \text{ given } \theta$$

- **Prior:**

$$p_{\Theta}(\theta) = \text{probability density of } \theta$$

- **Posterior:**

$$p_{\Theta|X}(\theta|x) = \text{probability density of } \theta \text{ given } x$$

- Bayes Theorem

$$\begin{aligned} p_{\Theta|X}(\theta|x) &= \frac{p_{X|\Theta}(x|\theta)p_{\Theta}(\theta)}{p_X(x)} \\ &= \frac{p_{X|\Theta}(x|\theta)p_{\Theta}(\theta)}{\int p_{X|\Theta}(x|\theta)p_{\Theta}(\theta)d\theta} \end{aligned}$$

Treating Linear Regression as Maximum-a-Posteriori

- We can modify the MLE by adding a prior

$$p_{\Theta}(\boldsymbol{\theta}) = \exp \left\{ -\frac{\rho(\boldsymbol{\theta})}{\beta} \right\}.$$

- Then, we have a MAP problem:

$$\begin{aligned}\boldsymbol{\theta}^* &= \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{n=1}^N \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(\boldsymbol{\theta}^T \mathbf{x}^n - y^n)^2}{2\sigma^2} \right\} \right\} \exp \left\{ -\frac{\rho(\boldsymbol{\theta})}{\beta} \right\} \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \frac{1}{2\sigma^2} \sum_{n=1}^N (\boldsymbol{\theta}^T \mathbf{x}^n - y^n)^2 + \frac{1}{\beta} \rho(\boldsymbol{\theta}) \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2 + \lambda \rho(\boldsymbol{\theta}), \quad \text{where } \lambda = 2\sigma^2/\beta.\end{aligned}$$

- $\rho(\cdot)$ is called **regularization function**.

Ridge Regression interpreted via a Gaussian prior

- One option: Choose a Gaussian prior

$$\exp \left\{ -\frac{\rho(\boldsymbol{\theta})}{\beta} \right\} = \exp \left\{ -\frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2} \right\}$$

- Then, the MAP becomes

$$\begin{aligned}\boldsymbol{\theta}^* &= \operatorname{argmax}_{\boldsymbol{\theta}} \prod_{n=1}^N \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(\boldsymbol{\theta}^T \mathbf{x}^n - y^n)^2}{2\sigma^2} \right\} \right\} \exp \left\{ -\frac{\|\boldsymbol{\theta}\|^2}{2\sigma_0^2} \right\} \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \sum_{n=1}^N (\boldsymbol{\theta}^T \mathbf{x}^n - y^n)^2 + \underbrace{\frac{\sigma^2}{\sigma_0^2}}_{=\lambda} \|\boldsymbol{\theta}\|^2 \\ &= \operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{A}\boldsymbol{\theta} - \mathbf{y}\|^2 + \lambda \|\boldsymbol{\theta}\|^2\end{aligned}$$

- This is exactly the **ridge regression**.

Proof of the Shrinkage Operator

$$\text{Let } J(\theta) = \frac{1}{2}(\theta - y)^2 + \lambda|\theta|.$$

$$0 = \frac{d}{d\theta}J(\theta) = (\theta - y) + \lambda\text{sign}(\theta).$$

- If $\theta > 0$, then $\theta = y - \lambda$. But since $\theta > 0$, it holds that $y > \lambda > 0$.
- If $\theta < 0$, then $\theta = y + \lambda$. But since $\theta < 0$, it holds that $y < -\lambda < 0$.
- If $\theta = 0$, then $\theta = y$. But since $\theta = 0$, it holds that $y = 0$.
- So the solution is

$$\hat{\theta} = \begin{cases} y - \lambda, & \text{if } y > 0, \\ 0 & \text{if } y = 0, \\ y + \lambda, & \text{if } y < 0. \end{cases}$$

- This is the same as

$$\hat{\theta} = \max(|y| - \lambda, 0)\text{sign}(y).$$