

ECE595 / STAT598: Machine Learning I

Lecture 2.1: Regularization - Ridge Regression

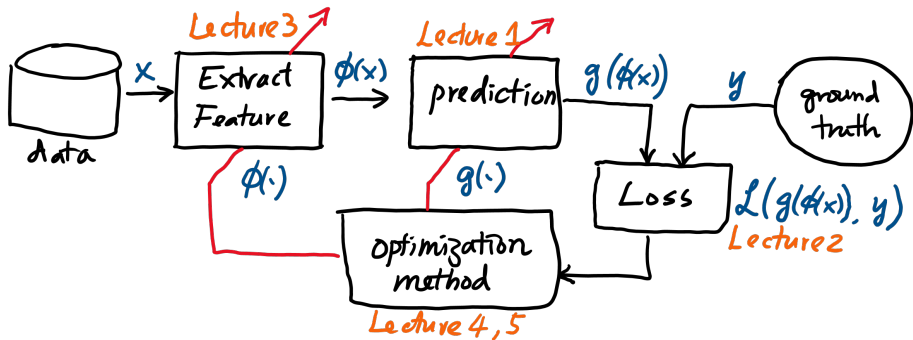
Spring 2020

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

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

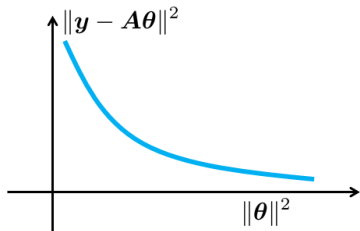
- If $\lambda \rightarrow 0$, then $\hat{\theta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$:

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

- 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

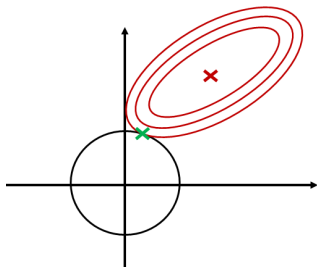
$$\theta_{\lambda}^* = \operatorname{argmin}_{\theta} \|\mathbf{A}\theta - \mathbf{y}\|^2 + \lambda\|\theta\|^2$$

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

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

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