

Regularization 1

Saturday, 18 February 2023 9:54 AM

Example: L^2 -regularized Linear Regression

$$\textcircled{1} R(w; X, y) = \frac{1}{2} \|Xw - y\|^2$$

Least Squares Formula: $\hat{w} = \underbrace{(X^T X)^{-1}}_H X^T y$

$$\textcircled{2} \tilde{R}(w; X, y) = \frac{1}{2} \|Xw - y\|^2 + \underbrace{\frac{1}{2} \alpha \|w\|^2}_{\text{regularizer}}$$

$$\text{Set } \nabla_w \tilde{R} = 0 \Rightarrow X^T (X\tilde{w} - y) + \alpha \tilde{w} = 0$$

$$(X^T X + \alpha I) \tilde{w} = \underbrace{X^T y}_{\substack{m \times N \\ N \times 1}} \in \mathbb{R}^m$$

$$\Rightarrow \tilde{w} = (H + \alpha I)^{-1} X^T y$$

Observe that H is symmetric positive definite, so.

$$\Rightarrow \begin{array}{ll} \text{eigenvectors:} & \{u_1, \dots, u_m\} \\ \text{eigenvalues:} & \{\lambda_1, \dots, \lambda_m\}, \end{array} \rightarrow \begin{array}{l} \text{orthonormal} \\ \text{basis of } \mathbb{R}^m \end{array} \quad (\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m)$$

$$X^T y = \sum_{i=1}^m \beta_i u_i \quad (\beta_i = u_i^T (X^T y))$$

$$\text{Hence, } \textcircled{1} \hat{w} = H^{-1} X^T y = H^{-1} \sum_{i=1}^m \beta_i u_i$$

$$\begin{aligned} &= \sum_{i=1}^m \beta_i \underbrace{H^{-1} u_i}_{u_i / \lambda_i} \\ &= \sum_{i=1}^m \frac{\beta_i}{\lambda_i} \cdot u_i \end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \quad \tilde{\omega} &= (H + \alpha I)^{-1} X^T y = (H + \alpha I)^{-1} \sum_{i=1}^m \beta_i u_i \\
 &= \sum_{i=1}^m \beta_i \underbrace{(H + \alpha I)^{-1} u_i}_{u_i / (\lambda_i + \alpha)} \\
 &= \sum_{i=1}^m \frac{\beta_i}{\lambda_i + \alpha} \cdot u_i
 \end{aligned}$$

Example: L^1 vs L^2 Regularization

$$R(\theta) = \frac{1}{2} \sum_{i=1}^m \lambda_i (\theta_i - \theta_i^*)^2$$

$\theta^* = (\theta_1^*, \dots, \theta_m^*) \in \mathbb{R}^m$
is given.

L^2 -regularized Problem:

$$\tilde{R}(\theta) = R(\theta) + \frac{1}{2} \alpha \|\theta\|^2 = \frac{1}{2} \sum_{i=1}^m \left[\lambda_i (\theta_i - \theta_i^*)^2 + \alpha \theta_i^2 \right]$$

$$\nabla_{\theta} \tilde{R} = 0 \Rightarrow \lambda_i (\theta_i - \theta_i^*) + \alpha \theta_i = 0$$

$$\therefore \theta_i = \frac{\lambda_i}{\lambda_i + \alpha} \theta_i^*$$

- $\theta_i \approx \theta_i^*$ if $\alpha \ll \lambda_i$
- $\theta_i \approx 0$ if $\alpha \gg \lambda_i$

L^1 -Regularized Problem:

$$\tilde{R}(\theta) = R(\theta) + \alpha \sum_{i=1}^m |\theta_i|$$

$$= \sum_{i=1}^m \left[\underbrace{\frac{1}{2} \lambda_i (\theta_i - \theta_i^*)^2}_{R_i(\theta_i)} + \alpha |\theta_i| \right]$$

Case 1: $\theta_i > 0$ attains minimum.

$$R_i(\theta_i) = \frac{1}{2} \lambda_i (\theta_i - \theta_i^*)^2 + \alpha \theta_i$$

$$\nabla_{\theta_i} R_i = 0 \Rightarrow \lambda_i (\theta_i - \theta_i^*) + \alpha = 0$$

$$\theta_i = \theta_i^* - \alpha / \lambda_i$$

(valid if $\theta_i^* > \alpha / \lambda_i$)

Case 2: $\theta_i < 0$ attains minimum

$$R_i(\theta_i) = \frac{1}{2} \lambda_i (\theta_i - \theta_i^*)^2 - \alpha \theta_i$$

$$\nabla_{\theta_i} R_i = 0 \Rightarrow \lambda_i (\theta_i - \theta_i^*) - \alpha = 0$$

$$\theta_i = \theta_i^* + \alpha / \lambda_i$$

(valid if $\theta_i^* < -\alpha / \lambda_i$)

Case 3: $\theta_i = 0$

valid if $-\alpha / \lambda_i \leq \theta_i^* \leq \alpha / \lambda_i$

Example: Early stopping for a linear regression

$$R(\theta) = \frac{1}{2} \lambda (\theta - \theta^*)^2 \Rightarrow \nabla R(\theta) = \lambda (\theta - \theta^*)$$

Consider GD on this:

$$\begin{aligned} \theta_{k+1} &= \theta_k - \varepsilon \lambda (\theta_k - \theta^*) \\ &= (1 - \varepsilon \lambda) \theta_k + \varepsilon \lambda \theta^* \end{aligned}$$

$$\begin{aligned} \Rightarrow \theta_k &= (1 - \varepsilon \lambda) \theta_{k-1} + \varepsilon \lambda \theta^* \\ &= (1 - \varepsilon \lambda) [(1 - \varepsilon \lambda) \theta_{k-2} + \varepsilon \lambda \theta^*] + \varepsilon \lambda \theta^* \\ &= \dots \\ &= (1 - \varepsilon \lambda)^k \theta_0 + [1 - (1 - \varepsilon \lambda)^k] \theta^* \end{aligned}$$

Stop at iteration τ :

$$\hat{\theta} = \theta_\tau = (1 - \varepsilon \lambda)^\tau \theta_0 + [1 - (1 - \varepsilon \lambda)^\tau] \theta^*$$

Consider a "variant" of L^2 -regularization

$$\tilde{R}(\theta) = \frac{1}{2}\lambda(\theta - \theta^*)^2 + \frac{1}{2}\alpha(\theta - \theta_0)^2$$

$$\nabla \tilde{R}(\theta) = 0 \Rightarrow \lambda(\tilde{\theta} - \theta^*) + \alpha(\tilde{\theta} - \theta_0) = 0$$

$$\tilde{\theta} = \frac{\alpha}{\lambda + \alpha} \cdot \theta_0 + \left(1 - \frac{\alpha}{\lambda + \alpha}\right) \theta^*$$

$$\tilde{\theta} = \hat{\theta} \text{ if } \frac{\alpha}{\alpha + \lambda} = (1 - \epsilon\lambda)^\tau \text{ i.e. } \alpha = \frac{\lambda(1 - \epsilon\lambda)^\tau}{1 - (1 - \epsilon\lambda)^\tau}$$