

2) Prove closed-form solution for Ridge Regression:
 $w = (\lambda I + X^T X)^{-1} X^T y$

$$\hookrightarrow \text{def } E(w) = \sum_{i=1}^m (w^T x^{(i)} - y^{(i)})^2 + \lambda \sum_{i=1}^m w_i^2$$

\hookrightarrow Differentiate $\nabla E(w)$

$$\begin{aligned} \hookrightarrow \nabla E(w) &= 2 \sum_i (w^T x^i - y^i) x^i + 2 \lambda w = 0 \\ &= 2 \sum_i (w^T x^i x^i - y^i x^i) + 2 \lambda w \\ &= \frac{2 X^T X w}{2} - \frac{2 X^T y}{2} + \frac{2 \lambda w}{2} = 0 \end{aligned}$$

$$\begin{aligned} \hookrightarrow w X^T X + \lambda w &= X^T y \\ w (X^T X + \lambda I) &= \frac{X^T y}{X^T X + \lambda I} \end{aligned}$$

$$\hookrightarrow w = \frac{X^T y}{X^T X + \lambda I} \quad \blacksquare$$

3) The posterior probability is $\hat{p}_k = \delta(s_k(\pi))_k = \frac{\exp(s_k(\pi))}{\sum_{j=1}^K \exp(s_j(\pi))}$

1. For the softmax regression model, we need to estimate a total of $(n+1) \cdot K$ parameters.

\hookrightarrow Total parameters = $K \cdot (n+1)$

$\hookrightarrow K$ is the total # of classes

$\hookrightarrow n$ is the dimensionality of x

2. Derive $J(\theta)$ with regards to θ_k

$$J(\theta) = \frac{-1}{n} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{p}_k^{(i)})$$

$$\hookrightarrow \hat{p}_k^{(i)} = \frac{\exp(s_k(x^i))}{\sum_{j=1}^K \exp(s_j(x^i))}$$

$$\hookrightarrow dJ(\theta)/d\theta_k = \frac{-1}{n} \sum_{i=1}^m (y_k^{(i)} - \hat{p}_k^{(i)}) \quad \blacksquare$$

