# 1. Ridge Regression

$$E(w) = \sum_{i=1}^{m} \left(w^T \cdot x^{(i)} - y^{(i)}\right)^2 + \lambda \sum_{i=1}^{m} w_i^2$$

Closed-form solution: $w = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T t$

$$E(w) = MSE(w) + \frac{\lambda}{2} \sum_{i=1}^{m} w_i^2$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left(h_w(x') - y_i\right)^2 + \frac{\lambda}{2} \sum_{i=1}^{m} w_i^2$$

$$E(w) = \frac{1}{m} (xw - y)^T (yw - y) + \frac{\lambda}{2} w^T w$$

$$= \left((xw)^T - y^T\right)(yw - y) + \lambda w^T w$$

$$= (xw)^T (yw) - (yw)^T(y)$$
$$\quad - (uw)^T(y) + y^T y + \lambda w + w$$

$$= w^T x^T \, xw - 2(xw)^T y + y^T y + \lambda w^T w$$

For minimizing of $w$

$$\frac{d\hat{w}}{dw} = 0$$

$$-2x^T y + 2(X^T X \lambda I)w = 0$$
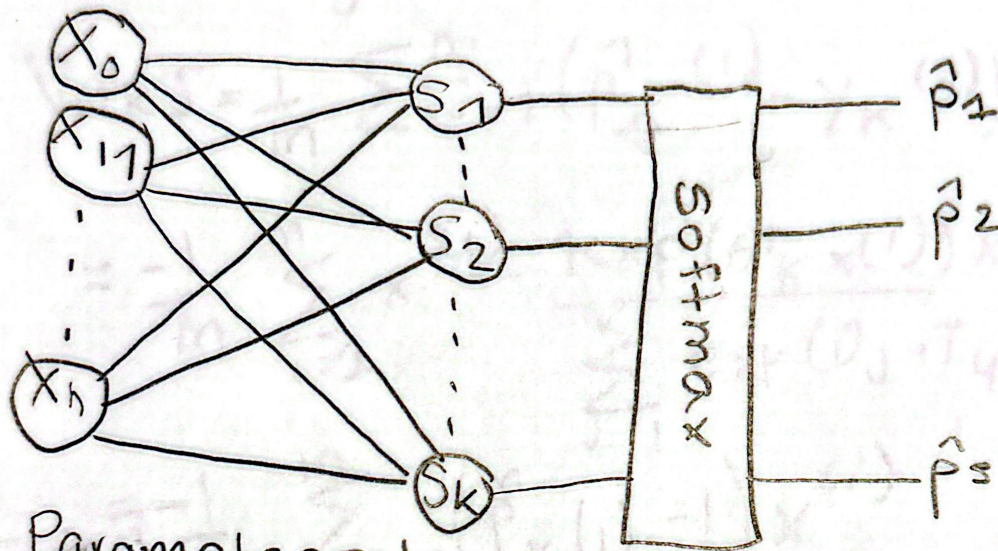
$$-2x^T y = -2(X^T X \lambda I)w$$

$$X^T y = (X^T X \lambda I) w$$

$$(X^T X + \lambda I) w = X^T y$$

1. Now Multiplying Both sides with $(x^T x + \lambda I)^{-1}$

2.

1. $s_k(x) = \theta_k^T \cdot x$    class $1 \leq k \leq K$

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

where $\theta_k$ is the vector param. of input feat for $s_k$.



Parameter estimation is $n+1$ and the parameters are $\theta_1, \theta_2 \ldots \theta_n$.

2. m training samples $\{(x_i, y_i)\}$ $i = 1, 2, \ldots m$
Derive the gradient of $J(\theta)$ regarding $\theta_k$.

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

where $y_k^{(i)} = 1$ if the $i^{th}$ instance belongs to k,

$$= \frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{K} y_k^{(i)} \log\left(\frac{\exp(s_k(x))}{\sum_{j=1}^{K}(s_j(x))}\right)$$

Scanned with CamScanner

$$= \frac{1}{m} \sum_{i=1}^{m} \left( \sum_{k=1}^{K} y_k^{(i)} \log \left( \exp \left( s_k (x^{(i)}) \right) - \right. \right.$$

$$\sum_{k=1}^{K} y_k^{(i)} \log \left( \sum_{j=1}^{K} \exp (s_j (x^{(i)})) \right)$$

$$y_k^{(i)} = 1$$

## Softmax Regression Entropy

$$\nabla_{\theta_k} J = \frac{1}{m} \sum_{i=1}^{m} \left( \hat{p}_k^{(i)} - y_k^{(i)} \right) x^{(i)}$$

$$= \frac{-1}{m} \sum_{i=1}^{m} x^{(i)} - \frac{1 \exp (\theta_k^T x^{(i)}) x^{(i)}}{\sum_{j=1}^{K} \exp (\theta_j \cdot T_y(i))}$$

$$= \frac{-1}{m} \sum_{i=1}^{m} \left( \hat{p}_k (i) - 1 \right) x^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left( \hat{p}_k (i) - y_k^{(i)} \right) x^{(i)} = \boxed{\nabla_{\theta_k} J(\theta)}$$