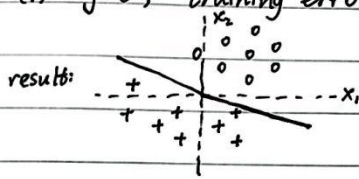
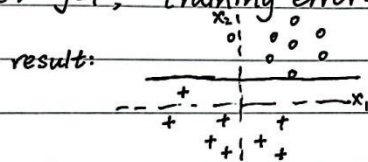
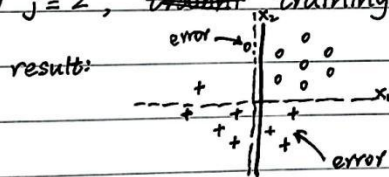


Homework 3

2. (1) $j=0$, training error = 0(2) $j=1$, training error = 0(3) $j=2$, ~~training~~ training error > 0(2) a. A. Training error can be 0 when $w_1=0$.

$$b. P(y=1|x, w) = P(y=0|x, w) = \frac{1}{2}$$

$$\Rightarrow -w_0 - w_1 x_1 - w_2 x_2 = 0$$

So w_0 must be 0.

$$c. P(y=1|x, w) > P(y=0|x, w)$$

$$\Rightarrow \exp(-w_0 - w_1 x_1 - w_2 x_2) > 1$$

$$\Rightarrow w_0 + w_1 x_1 + w_2 x_2 > 0$$

So $w_0 > 0$.3. (1) We want to minimize $\|y - Xw\|_2^2$.

$$\text{So } \frac{\partial \|y - Xw\|_2^2}{\partial w} = 0, \text{ that is, } -2X^T(y - Xw) = 0$$

$$\text{So } \hat{w} = (X^T X)^{-1} X^T y$$

$$\hat{y} = X\hat{w} = X(X^T X)^{-1} X^T y$$

No.

Date

$$(2) P^2 = X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = P$$

$$\lambda x = Px = P^2 x = \lambda^2 x$$

$$\Rightarrow \lambda = \lambda^2 \Rightarrow \lambda = 0 \text{ or } \lambda = 1$$

$$(3) E(\hat{w}) = E((X^T X)^{-1} X^T y)$$

$$= E((X^T X)^{-1} X^T (Xw + \epsilon))$$

$$= E((X^T X)^{-1} X^T Xw)$$

$$= E(w)$$

$$= w$$

$$\text{Var}(\hat{w}) = E[(\hat{w} - E(\hat{w}))(\hat{w} - E(\hat{w}))^T]$$

$$= E[(\hat{w} - w)(\hat{w} - w)^T]$$

$$= E[(X^T X)^{-1} X^T \epsilon (\epsilon^T X (X^T X)^{-1})^T]$$

$$= E[(X^T X)^{-1} X^T \epsilon \epsilon^T X (X^T X)^{-1}]$$

$$= E[(X^T X)^{-1} X^T X (X^T X)^{-1}] \cdot E(\epsilon \epsilon^T)$$

$$= E[(X^T X)^{-1}] \sigma^2$$

$$= (X^T X)^{-1} \sigma^2$$

$$(4) SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})$$

$$= SS_{res} + SS_{reg} + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})$$

$$\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = (y - \hat{y})^T (X\hat{w} - \bar{y} \mathbf{1})$$

$$= (y - \hat{y})^T X\hat{w} - \bar{y} (y - \hat{y})^T \mathbf{1}$$

$$\text{From (1), } \hat{y} = X(X^T X)^{-1} X^T y,$$

$$X^T (y - \hat{y}) = X^T y - X^T X(X^T X)^{-1} X^T y = X^T y - X^T y = 0$$

$$\text{First column of } X \text{ is } \mathbf{1}, \text{ so } \mathbf{1}^T (y - \hat{y}) = 0$$

$$\text{So } \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) = 0$$

$$\Rightarrow SS_{tot} = SS_{res} + SS_{reg}$$

$$4. (1) \sum_{i=1, i \neq k}^n (y_i - x_i^T w)^2 + \lambda \|w\|_2^2 = \|y - Xw\|_2^2 - (y_k - x_k^T w)^2 + \lambda \|w\|_2^2$$

$$\frac{\partial (\|y - Xw\|_2^2 - (y_k - x_k^T w)^2 + \lambda \|w\|_2^2)}{\partial w} = -2X^T(y - Xw) + 2x_k(x_k^T w - y_k) + 2\lambda w = 0$$

$$\Rightarrow (X^T X + \lambda I - x_k x_k^T) w = X^T y - x_k y_k$$

$$\Rightarrow \hat{w}^{[k]} = (X^T X + \lambda I - x_k x_k^T)^{-1} (X^T y - x_k y_k)$$

(2) Using Sherman-Morrison, we get

$$(X^T X + \lambda I - x_k x_k^T)^{-1} = (X^T X + \lambda I)^{-1} + \frac{(X^T X + \lambda I)^{-1} x_k x_k^T (X^T X + \lambda I)^{-1}}{1 - x_k^T (X^T X + \lambda I)^{-1} x_k}$$

$$\text{And we have } x_k^T (X^T X + \lambda I)^{-1} x_k = p_{kk}, \hat{y}_k = x_k^T (X^T X + \lambda I)^{-1} X^T y$$

$$\text{Then } x_k^T \hat{w}^{[k]} - y_k = x_k^T \left[(X^T X + \lambda I)^{-1} + \frac{(X^T X + \lambda I)^{-1} x_k x_k^T (X^T X + \lambda I)^{-1}}{1 - x_k^T (X^T X + \lambda I)^{-1} x_k} \right] (X^T y - x_k y_k) - y_k$$

$$= \hat{y}_k - p_{kk} y_k + \frac{p_{kk} \hat{y}_k}{1 - p_{kk}} - \frac{p_{kk} y_k}{1 - p_{kk}} - y_k$$

$$= \frac{\hat{y}_k - y_k}{1 - p_{kk}}$$

$$V_o(\lambda) = \frac{1}{n} \sum_{k=1}^n (x_k^T \hat{w}^{[k]} - y_k)^2 = \frac{1}{n} \sum_{k=1}^n \left(\frac{\hat{y}_k - y_k}{1 - p_{kk}} \right)^2$$

$$(3) V(\lambda) = \frac{1}{n} \sum_{k=1}^n \left(\frac{1 - p_{kk}}{\frac{1}{n} \text{tr}(I - P)} \right)^2 \left(\frac{\hat{y}_k - y_k}{1 - p_{kk}} \right)^2$$

$$= \frac{1}{n} \sum_{k=1}^n \left(\frac{\hat{y}_k - y_k}{\frac{1}{n} \text{tr}(I - P)} \right)^2$$

$$= \frac{1}{n} \left(\frac{1}{\frac{1}{n} \text{tr}(I - P)} \right)^2 \sum_{k=1}^n (\hat{y}_k - y_k)^2$$

$$= \frac{n}{(\text{tr}(I) - \text{tr}(P))^2} \| \hat{y} - y \|^2$$

$$= \frac{n}{(n - \text{tr}(P))^2} \| Py - y \|^2$$

$$= \frac{\frac{1}{n} \| (I - P)y \|^2}{(1 - \text{tr}(P)/n)^2}$$

5. LASSO: $\min_{w, z} [\|y - Xw\|_2^2 + \lambda \|z\|_1]$, subject to $w - z = 0$

$$L(w, z, u) = \|y - Xw\|_2^2 + \lambda \|z\|_1 + u^T (w - z) + \frac{1}{2} \|w - z\|_2^2$$

$$\text{Step 1: } w^{(k+1)} = \arg \min_w L(w, z^{(k)}, u^{(k)})$$

$$\frac{\partial L}{\partial w} = -2X^T(y - Xw) + u^{(k)} + (w - z^{(k)}) = 0$$

$$\Rightarrow w^{(k+1)} = (I + 2X^T X)^{-1} (2X^T y - u^{(k)} + z^{(k)})$$

↓ Next Page

No.

Date

$$\text{Step 2: } \mathbf{z}^{(k+1)} = \arg \min_{\mathbf{z}} L(\mathbf{w}^{(k+1)}, \mathbf{z}, \mathbf{u}^{(k)})$$

$$\frac{\partial L}{\partial \mathbf{z}} = \lambda \|\mathbf{z}\|_1 - \mathbf{u}^{(k)} - (\mathbf{w}^{(k+1)} - \mathbf{z})$$

$$\frac{\partial}{\partial z_j} \|\mathbf{z}\|_1 = \begin{cases} 1, & z_j > 0 \\ [-1, 1], & z_j = 0 \\ -1, & z_j < 0 \end{cases}$$

$$\frac{\partial L}{\partial z_j} = 0 \Rightarrow \tilde{z}_j^{(k+1)} = \begin{cases} u_j^{(k)} + w_j^{(k+1)} - \lambda, & u_j^{(k)} + w_j^{(k+1)} > \lambda, \\ 0, & u_j^{(k)} + w_j^{(k+1)} \in [-\lambda, \lambda], \\ u_j^{(k)} + w_j^{(k+1)} + \lambda, & u_j^{(k)} + w_j^{(k+1)} < -\lambda \end{cases}$$

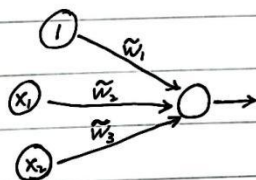
$$\text{Step 3: Iterate } \mathbf{u}^{(k+1)} = \mathbf{u}^{(k)} + \mathbf{w}^{(k+1)} - \mathbf{z}^{(k+1)}$$

$$6. (1) \quad \mathbf{z}_1 = w_1 + w_3 x_1 + w_5 x_2, \quad \mathbf{z}_2 = w_2 + w_4 x_1 + w_6 x_2$$

$$\mathbf{z}' = w_7 + w_8 h_1 + w_9 h_2 = w_7 + w_8 \mathbf{z}_1 + w_9 \mathbf{z}_2$$

$$P(y=1|x, w) = g(\mathbf{z}') = \frac{1}{1 + e^{-(w_7 + w_8(w_1 + w_3 x_1 + w_5 x_2) + w_9(w_2 + w_4 x_1 + w_6 x_2))}}$$

(2)



$$w_7 + w_8 C(w_1 + w_3 x_1 + w_5 x_2) + w_9 C(w_2 + w_4 x_1 + w_6 x_2)$$

$$= (w_7 + w_8 C w_1 + w_9 C w_2) + (w_8 C w_3 + w_9 C w_4) x_1 + (w_8 C w_5 + w_9 C w_6) x_2$$

$$\Rightarrow \tilde{w}_1 = w_7 + w_8 C w_1 + w_9 C w_2$$

$$\tilde{w}_2 = w_3 w_8 C + w_4 w_9 C$$

$$\tilde{w}_3 = w_5 w_8 C + w_6 w_9 C$$

(3) No. Not all activation function is linear. But if all activation functions are linear, it can be without hidden layer.