

Machine learning HW3

BN901174 傅啟恩

$y_n \backslash \hat{y}$	+1	-1
+1	0	1
-1	1000	0

avg error

$$J = \begin{cases} +1 & p(-1|x) \times 1 \\ -1 & p(+1|x) \times 1000 \end{cases}$$

$$p(-1|x) + p(+1|x) = 1$$

$$p(-1|x) = 1 - p(+1|x)$$

$$\text{If } (1 - p(+1|x)) \times 1 < p(+1|x) \times 1000$$

$$\Rightarrow 1001 p(+1|x) > 1$$

ideal mini-target

$$p(+1|x) > \alpha, f_{CLA}(x) = 1$$

$$\Rightarrow f_{CLA}(x) = \text{sign}(p(+1|x) - \alpha)$$

$$f_{CLA}(x) = +1$$

$$\alpha = \frac{1}{1001} \quad \#$$

$$\text{otherwise } f_{CLA}(x) = -1$$

$$2. E_{(x,y) \sim p(x,y)} (g(x) \neq y) = E_{x \sim p(x)} (g(x) \neq f(x)) \cdot p(y = -f(x)|x)$$

$$\leftarrow E_{out}(g) + E_{x \sim p(x)} (g(x) \neq f(x)) \cdot p(y = f(x)|x)$$

$$= \frac{(1 - E_{out}(g)) \cdot \epsilon + E_{out}(g) (1 - \epsilon)}{\quad} \quad \#$$

$$3. E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2 = \frac{1}{N} \sum_{n=1}^N (w x_n - y_n)^2$$

$$= \frac{1}{N} \sum_{n=1}^N (w^2 x_n^2 - 2w x_n y_n + y_n^2)$$

$$= \frac{1}{N} \left(w^2 \sum_{n=1}^N x_n^2 - 2w \sum_{n=1}^N x_n y_n + \sum_{n=1}^N y_n^2 \right)$$

$$\nabla E_{in}(w) = \frac{1}{N} \left(2w \sum_{n=1}^N x_n^2 - 2 \sum_{n=1}^N x_n y_n \right) \stackrel{\text{set}}{=} 0$$

$$\text{optimal } w \quad w_{LW} = \frac{\sum_{n=1}^N x_n y_n}{\sum_{n=1}^N x_n^2} \quad \#$$

$$4. E_{in}(w) = \frac{1}{N} \int_0^1 ((w_0 + w_1 x) - (ax^2 + b))^2 dx$$

$$(N=1) = \frac{1}{N} \int_0^1 (-ax^2 + w_1 x + (w_0 - b))^2 dx$$

$$= \frac{1}{N} \int_0^1 (a^2 x^4 - 2aw_1 x^3 - 2a(w_0 - b)x^2 + w_1^2 x^2 + 2w_1(w_0 - b)x + (w_0 - b)^2) dx$$

$$= a^2 \left(\frac{1}{N} \int_0^1 x^4 dx \right) - 2aW_1 \left(\frac{1}{N} \int_0^1 x^3 dx \right) + (-2a(W_0 - b) + W_1^2) \left(\frac{1}{N} \int_0^1 x^2 dx \right) + 2W_1(W_0 - b) \left(\frac{1}{N} \int_0^1 x dx \right) + (W_0 - b)^2 \quad x \sim \text{UNIF}(0,1)$$

$$= \frac{1}{5} a^2 - \frac{1}{2} a W_1 + (-2a(W_0 - b) + W_1^2) \cdot \frac{1}{3} + W_1(W_0 - b) + (W_0 - b)^2$$

$$\frac{\partial \mathcal{E}_M}{\partial W_0} = W_1 - \frac{2}{3} a + 2(W_0 - b) = 0$$

$$2W_1 + 4W_0 = \frac{4}{3} a + 4b$$

$$2W_1 + 3W_0 = \frac{3}{2} a$$

$$\frac{\partial \mathcal{E}_M}{\partial W_1} = \frac{2}{3} W_1 - \frac{1}{2} a + W_0 = 0$$

$$\begin{cases} W_0^* = \frac{1}{6} a + 4b \\ W_1^* = \frac{1}{3} a - 6b \end{cases} \quad \#$$

$$5. \quad y_n' = a y_n + b \quad X^T y = W_{L2N}$$

$$W_{L2N}' = a W_{L2N} + \underbrace{(X^T X)^{-1} X^T b}_{\hat{b}}$$

$$\text{we know that } \hat{b} = \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix}$$

$$\text{and } X = \begin{bmatrix} 1 & -x_1^T & \dots \\ 1 & -x_2^T & \dots \\ \vdots & \vdots & \ddots \\ 1 & -x_N^T & \dots \end{bmatrix}$$

$$\Leftrightarrow X^T X W_b - X^T b = 0,$$

$$\Leftrightarrow X^T (X W_b - \hat{b}) = 0, \quad W_b = (X^T X)^{-1} X^T b$$

$$W_{L2N}' = a W_{L2N} + \left(W_b = \begin{bmatrix} b \\ 0 \\ \vdots \\ 0 \end{bmatrix} \right) = \begin{bmatrix} b \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow \begin{matrix} W_b[0] = b & i=0 \\ W_b[i] = 0 & i \neq 0 \end{matrix}$$

$$6. \quad \nabla \mathcal{E}_M(\hat{w}_t) = \frac{1}{N} \sum_n \frac{1}{1 + \exp(y_n \hat{w}_t^T x_n)} (-y_n \hat{x}_n)$$

$$A_{i,i} = \frac{\partial (\nabla \mathcal{E}_M(\hat{w}_t)_i)}{\partial w_i} = \frac{1}{N} \sum_n \left(\frac{+1}{1 + \exp(y_n \hat{w}_t^T x_n)} \right)^2 y_n^2 (x_{n,i})^2 \exp(y_n \hat{w}_t^T x_n) \\ = \frac{1}{N} \sum_n (h_t(y_n x_n))^2 \exp(y_n \hat{w}_t^T x_n) \cdot y_n^2 (x_{n,i}) (x_{n,i})$$

$$A_{i,j} = \frac{1}{N} \sum_{n=1}^N [h_t(y_n x_n)]^2 \exp(y_n \hat{w}_t^T x_n) \cdot y_n^2 (x_{n,i}) (x_{n,j})$$

$$= \sum_{n=1}^N (x_{n,i})^T (x_{n,j}) D_{nn}, \quad D_{nn} = \begin{bmatrix} D_{11} & \dots & 0 \\ \vdots & D_{22} & \vdots \\ 0 & \dots & D_{NN} \end{bmatrix}$$

$$A = X^T D X$$

$$D_{i,j} = \begin{cases} \frac{1}{N} [\text{ht}(y_n x_n)]^2 \exp(y_n \hat{w}_i^T \hat{x}_n) y_n^2 & i=j \\ 0 & i \neq j \end{cases}$$

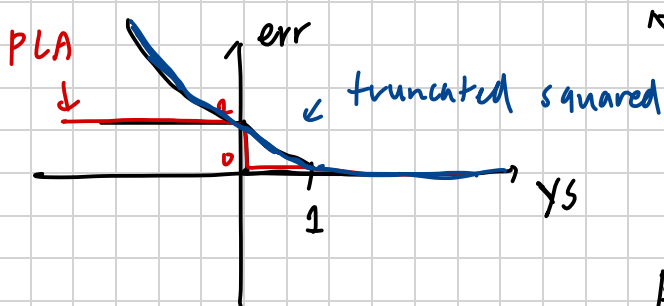
7. $\text{err}(s, y) = \begin{cases} (1-sy)^2 & \text{if } sy \leq 1 \\ 0 & \text{otherwise} \end{cases}$

Applying SGD

↙ truncated squared loss

$$W_{t+1} \leftarrow W_t + \eta (1(y_s \leq 1) 2(y_n - w_t^T x_n) x_n)$$

$$W_{t+1} \leftarrow W_t + \eta (1(y_s \leq 0) y_n x_n)$$



↖ PLA

Compared to PLA, The error is the same \Rightarrow if $y_s = y w^T x \geq 1$

But if $y_s = y w^T x \leq 1$, the truncated squared loss would be the upper bound of PLA method #

8. Find $\nabla E_{in}(w)$

$$W = \begin{pmatrix} w_1 & w_2 & \dots & w_k \\ w_{11} & w_{12} & \dots & w_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{nk} \end{pmatrix}_{(n+1) \times k}$$

$$h_y(x) = \frac{\exp(w_y^T x)}{\sum_{i=1}^K \exp(w_i^T x)}$$

$$\text{and } E_{in}(w) = \frac{1}{N} \sum_n \text{err}(w, x_n, y_n) = \frac{1}{N} \sum_n -\ln h_{y_n}(x_n)$$

$$\text{matrix} = \frac{-1}{N} (\ln h_{y_1}(x_1) + \ln h_{y_2}(x_2) + \dots + \ln h_{y_N}(x_N))$$

$$\ln h_{y_n}(x_n) = (w_{y_n}^T x_n) - \ln \left(\sum_{i=1}^K \exp(w_i^T x_n) \right)$$

First solve $\frac{\partial E_{in}(w)}{\partial w_{i,j}} \leftarrow \text{element / entry}$

$$\frac{\partial (\ln h_{y_n}(x_n))}{\partial w_{i,j}} = \begin{cases} x_{n,i} - \frac{1}{\sum_{i=1}^K \exp(w_i^T x_n)} \exp(w_i^T x_n) \cdot x_{n,i} & \text{if } j = y_n \end{cases}$$

$$\frac{\partial \mathcal{E}_m(w)}{\partial w_{ij}} = \left[- \frac{1}{\sum_{i=1}^K \exp(w_i^T x_n)} \exp(w_{ij} x_{n,j}) x_{n,j} \text{ if } j \neq y_n \right]$$

w_i each element
 is as above
 $i = (0, d+1)$
 $j = (0, K)$
 y_n

$$\frac{\partial \mathcal{E}_m(w)}{\partial w} = \begin{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} & \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \end{bmatrix} \begin{matrix} w_k \\ \vdots \end{matrix}$$

(d+1) x K

13. By problem 6

$$A = X^T D X = \tilde{X}^T \tilde{X} = \sum_{n=1}^N (\hat{x}_{i,n})^T (\hat{x}_{n,j})$$

$$A_{i,j} = \sum_{n=1}^N (\hat{x}_{i,n})^T (\hat{x}_{n,j}) D_{nn}$$

\Rightarrow compare $\tilde{X}^T \tilde{X}$ and $X^T D X$, we can know that

$$(\hat{x}_{i,n}) = (x_{i,n}) \sqrt{D_{nn}}, \quad (\hat{x}_{n,j}) = (x_{n,j}) \sqrt{D_{nn}}$$

$$\text{where } \sqrt{D_{nn}} = \begin{bmatrix} \sqrt{D_{11}} & \cdots & 0 \\ \vdots & \sqrt{D_{22}} & \vdots \\ 0 & \cdots & \sqrt{D_{NN}} \end{bmatrix}$$

$$\tilde{X} = X \cdot \sqrt{D}$$

$$\sqrt{\frac{1}{N} \cdot \exp \cdot h \cdot y_n}$$

$$\sqrt{D_{i,j}} = \begin{cases} \sqrt{\frac{1}{N} [h(y_n x_n)]^2 \exp(y_n \hat{w}_i^T \hat{x}_n) y_n^2} & i=j \\ 0 & i \neq j \end{cases}$$

$$X^T (y_n \cdot y \cdot (X w_t))$$

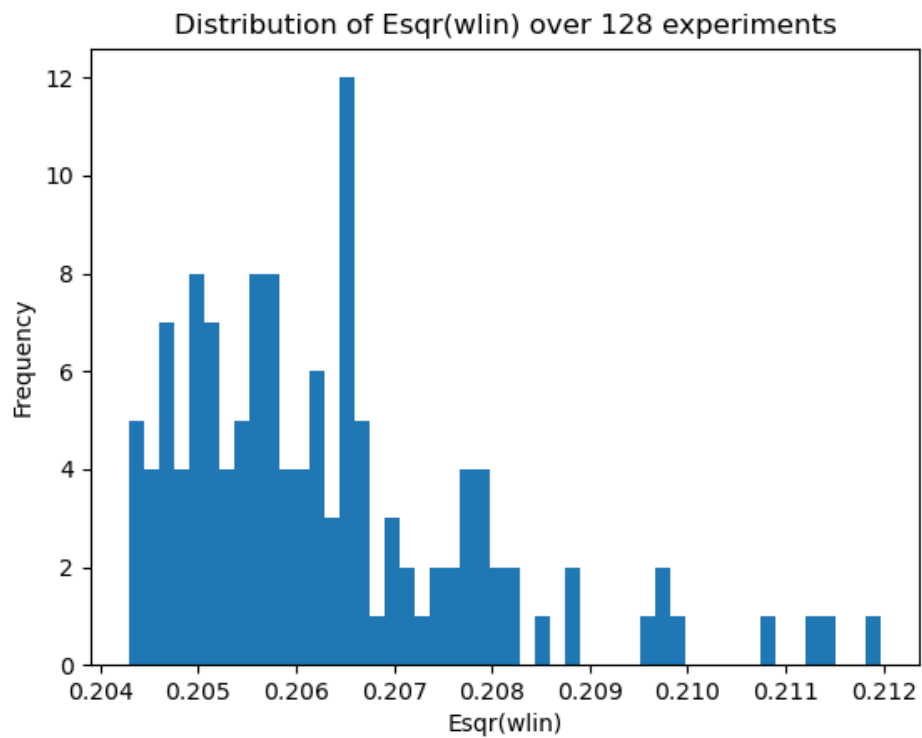
$$\text{and } \frac{1}{N} \sum_n \frac{1}{1 + \exp(y_n w_t^T x_n)} (+ y_n \hat{x}_n) = -\nabla \mathcal{E}_m(w_t) = \tilde{X}^T \hat{y}$$

$$-\nabla \mathcal{E}_m(w_t)_i = \sum_n \hat{x}_{i,n} \hat{y}_n = \sum_n (x_{i,n}) \sqrt{D_{nn}} \hat{y}_n$$

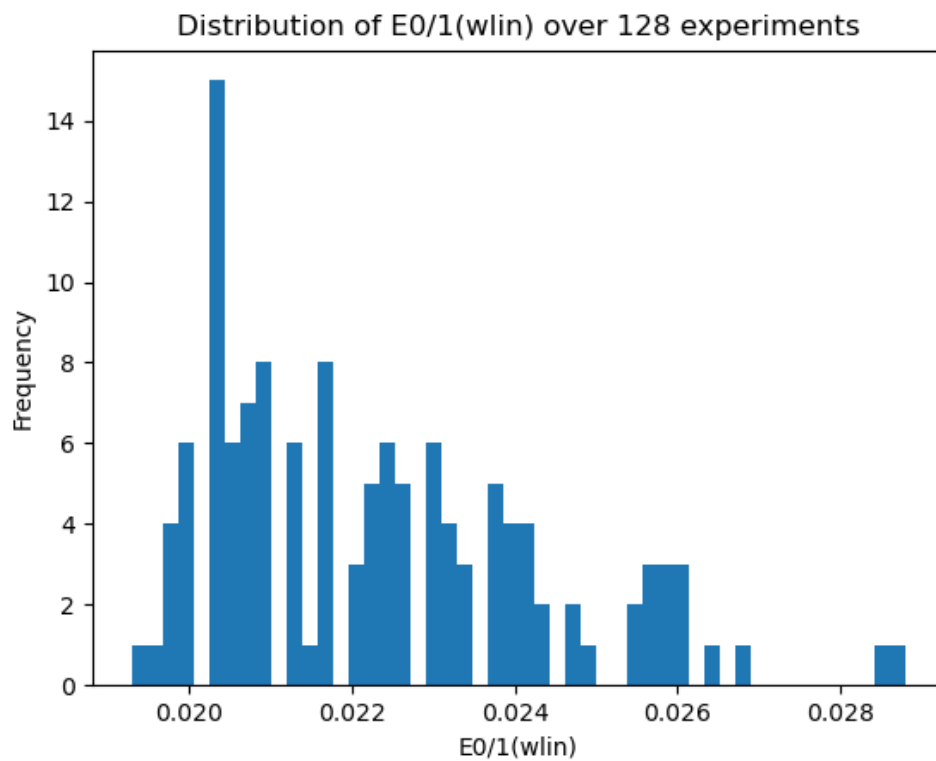
$$\hat{y}_n = \sqrt{\frac{1}{N}} \frac{1}{\sqrt{\exp(y_n w_t^T x_n)}}, \quad \hat{y} = \sqrt{\frac{1}{N}} \cdot \frac{1}{\sqrt{\exp(y \cdot X w_t)}} \quad \#$$

array matrix

9. Median Esqr over 128 experiments: 0.20595399486341354



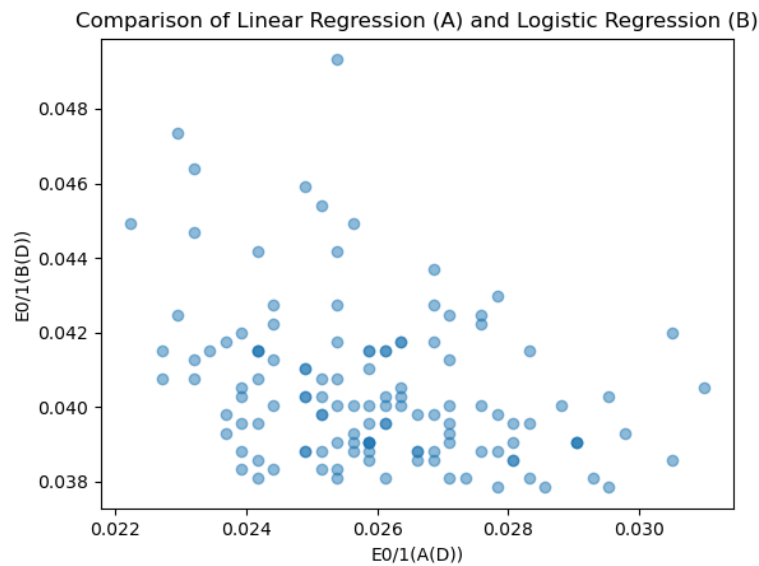
10. Median E0/1 over 128 experiments: 0.02197265625



11.

Median $E0/1(A(D))$ over 128 experiments: 0.022216796875

Median $E0/1(B(D))$ over 128 experiments: 0.0357666015625



12.

Median $E0/1(A(D))$ over 128 experiments: 0.0350341796875

Median $E0/1(B(D))$ over 128 experiments: 0.038818359375

We see that as we add outlier examples, the error rate is getting higher for linear regression and for logistic regression, which is due to new data with different mean and covariance being added. Compared to two types of regressions, we see that linear regression changes more obviously than logistic one.

