

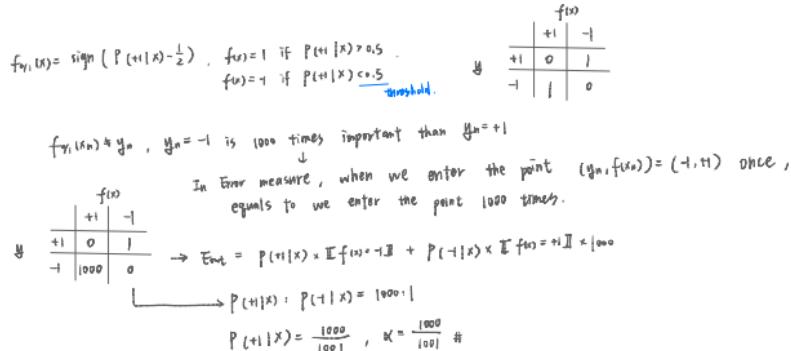
1. (20 points) Consider a binary classification problem, where $\mathcal{Y} = \{-1, +1\}$. Assume a noisy scenario where the data is generated i.i.d. from some $P(\mathbf{x}, y)$. In class, we discussed that when the 0/1 error function (i.e. classification error) is considered, calculating the "ideal mini-target" on each \mathbf{x} reveals the hidden target function of

$$f_{0/1}(\mathbf{x}) = \operatorname{argmax}_{y \in \{-1, +1\}} P(y|\mathbf{x}) = \operatorname{sign}\left(P(+1|\mathbf{x}) - \frac{1}{2}\right).$$

Instead of the 0/1 error, if we consider the CIA error function, where a false positive (classifying a negative example as a positive one) is 1000 times more important than a false negative, the hidden target should be changed to

$$f_{CIA}(\mathbf{x}) = \operatorname{sign}(P(+1|\mathbf{x}) - \alpha).$$

Prove what the value of α should be.



2. (20 points) Consider a binary classification task, where God gives you some noiseless data i.i.d. from an unknown distribution $P(x)$ and an unknown target function $f(x)$ that maps from \mathcal{X} to $\{-1, +1\}$. After you use the data to obtain some $g(x)$ that suffers

$$E_{out}(g) = \mathbb{E}_{x \sim P(x)} [g(x) \neq f(x)] \text{ (here } \mathbb{E} \text{ means expectation, as shown in class slides)}$$

$$= \mathbb{E}_{x \sim P(x)} [g(x) \neq f(x)] \text{ (or if you like the more beautiful font E for expectation).}$$

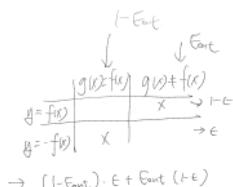
Now, assume that $g(x)$ is put in a noisy test environment where

$$\begin{aligned} P(y = +f(x)|x) &= 1 - \epsilon \\ P(y = -f(x)|x) &= \epsilon. \end{aligned}$$

Derive

$$\mathbb{E}_{(x,y) \sim P(x,y)} [g(x) \neq y]$$

as a function of $E_{out}(g)$ and ϵ .



$$\begin{aligned} \text{In ideal case } P[g(x) \neq f(x)] &= E_{out} \\ P[g(x) \neq f(x)] &= 1 - E_{out} \end{aligned}$$

$$\begin{aligned} \text{In noise } P[y \neq f(x) | x] &= 1 - \epsilon \\ P[y = -f(x) | x] &= \epsilon \end{aligned}$$

	$g(x) = f(x)$	$g(x) \neq f(x)$	
$y = f(x)$	s_1	s_2	$1 - \epsilon$
$y = -f(x)$	s_3	s_4	ϵ
P_2	$1 - E_{out}$	E_{out}	

$$\begin{aligned} \mathbb{E}_{(x,y) \sim P(x,y)} [g(x) \neq y] &= \mathbb{E}_{x \sim P(x)} [g(x) \neq f(x)] \cdot P(y = +f(x) | x) + \\ &\quad \mathbb{E}_{x \sim P(x)} [g(x) \neq f(x)] \cdot P(y = -f(x) | x) \\ &= (1 - E_{out}) \cdot \epsilon + E_{out} \cdot (1 - \epsilon) \end{aligned}$$

3. (20 points) Consider a hypothesis set that contains hypotheses of the form $h(x) = wx$ for $x \in \mathbb{R}$. Combine the hypothesis set with the squared error function to minimize

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

on a given data set $\{(x_n, y_n)\}_{n=1}^N$. Derive the optimal w_{LIN} in terms of (x_n, y_n) and express the result without using matrix/vector notations. You can assume all denominators to be non-zero.

(Hint: This is linear regression in \mathbb{R} without the added x_0 .)

$$h(x) = wx$$

$$E_{\text{lin}}(w) = \frac{1}{N} \sum_{n=1}^N (wx_n - y_n)^2$$

w_{LIN} exists when $\min E_{\text{lin}}(w)$ appears $\rightarrow \nabla E_{\text{lin}}(w) = 0$

$$E_{\text{lin}}(w) = \frac{1}{N} [(wx_1 - y_1)^2 + (wx_2 - y_2)^2 + (wx_3 - y_3)^2 + \dots + (wx_n - y_n)^2]$$

$$= \frac{1}{N} [w^2 x_1^2 - 2wx_1 y_1 + y_1^2 + w^2 x_2^2 - 2wx_2 y_2 + y_2^2 + \dots + w^2 x_n^2 - 2wx_n y_n + y_n^2]$$

$$\forall E_{\text{lin}}(w) = \frac{1}{N} (z_N x_1^2 - z_N x_1 y_1 + z_N x_2^2 + z_N x_2 y_2 + \dots + z_N x_n^2 + z_N x_n y_n)$$

$$= \frac{2}{N} \left(w \sum_{k=1}^{N-1} x_k^2 - \sum_{k=1}^{N-1} x_k y_k \right) = 0$$

$$\rightarrow w = \frac{\sum_{k=1}^N x_k y_k}{\sum_{k=1}^N x_k^2} \#$$

4. (20 points) Consider the target function $f(x) = ax^2 + b$. Sample x uniformly from $[0, 1]$, and use all linear hypotheses $h(x) = w_0 + w_1 \cdot x$ to approximate the target function with respect to the squared error. For any given (a, b) , derive the weights (w_0^*, w_1^*) of the optimal hypothesis as a function of (a, b) .

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (f(x_n) - h(x_n))^2 = \frac{1}{N} \sum_{n=1}^N (ax_n^2 + b - w_0 - w_1 x_n)^2$$

$$\text{since } x \text{ uniformly from } [0, 1], E_{in}(w_0, w_1) = \int_0^1 (ax^2 + b - w_0 - w_1 x)^2 dx$$

$$\min E_{in}(w_0, w_1) \Leftrightarrow \nabla E_{in}(w_0, w_1) = 0$$

$$\begin{cases} \frac{\partial}{\partial w_0} = \int_0^1 (ax^2 - w_1 x + b - w_0) \cdot 2 \cdot (-1) dx = 0 \\ \frac{\partial}{\partial w_1} = \int_0^1 (ax^2 - w_1 x + b - w_0) \cdot 2 \cdot (-x) dx = 0 \end{cases}$$

$$\int_0^1 (ax^2 - w_1 x + b - w_0) dx = 0$$

$$\int_0^1 (ax^3 - w_1 x^2 + bx - w_0 x) dx = 0$$

$$\left\{ \begin{array}{l} \frac{1}{3}ax^3 - \frac{1}{2}w_1 x^2 + (b-w_0)x \Big|_0^1 = 0 \\ \frac{1}{4}ax^4 - \frac{1}{3}w_1 x^3 + \frac{1}{2}(b-w_0)x^2 \Big|_0^1 = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} \frac{1}{3}a - \frac{1}{2}w_1 + (b-w_0) = 0 \\ \frac{1}{4}a - \frac{1}{3}w_1 + \frac{1}{2}b - \frac{1}{2}w_0 = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} w_0 + \frac{1}{2}w_1 = \frac{1}{3}a + b \\ \frac{1}{2}w_0 + \frac{1}{3}w_1 = \frac{1}{4}a + \frac{1}{2}b, \quad w_0 + \frac{2}{3}w_1 = \frac{1}{4}a + b \end{array} \right.$$

$$\rightarrow \frac{1}{6}w_1 = \frac{1}{6}a, \quad w_1^* = a \#$$

$$w_0^* = -\frac{1}{6}a + b \#$$

5. (20 points) Consider running linear regression on $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, where \mathbf{x}_n includes the constant dimension $x_0 = 1$ as usual. For simplicity, you can assume that $\mathbf{X}^T \mathbf{X}$ is invertible. Assume that the unique (why \vdash) solution \mathbf{w}_{LIN} is obtained after running linear regression on the data above. Now, consider an output transformation

$$\mathbf{W}_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$y'_n = a y_n + b.$$

for some given constants (a, b) . Run linear regression on $\{(\mathbf{x}_n, y'_n)\}_{n=1}^N$ to obtain the unique solution \mathbf{w}'_{LIN} . Derive \mathbf{w}'_{LIN} as a function of \mathbf{w}_{LIN} and (a, b) .

$$\mathbf{W}_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad \text{--- ①}$$

$$E_{\text{lin}}(\mathbf{w}') = \frac{1}{N} \| \mathbf{X}^T \mathbf{X} \mathbf{w}' - (ay+b) \|_2^2 = \frac{1}{N} (\mathbf{w}'^T \mathbf{X}^T \mathbf{X} \mathbf{w}' - 2 \mathbf{w}'^T (\mathbf{X}^T \mathbf{y} + b) + (\mathbf{X}^T \mathbf{y} + b)^T (\mathbf{X}^T \mathbf{y} + b))$$

$$\nabla E_{\text{lin}}(\mathbf{w}') = \frac{1}{N} (2 \mathbf{X}^T \mathbf{X} \mathbf{w}' - 2 \mathbf{X}^T (\mathbf{X}^T \mathbf{y} + b))$$

$$\mathbf{W}'_{\text{LIN}} \text{ happens when } \nabla E_{\text{lin}} = 0 : 2 \mathbf{X}^T \mathbf{X} \mathbf{w}' = 2 \mathbf{X}^T (\mathbf{X}^T \mathbf{y} + b), \quad \mathbf{W}'_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}^T \mathbf{y} + b) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot \mathbf{ay} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot b$$

$$\begin{bmatrix} a (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot b \end{bmatrix}$$

$$\text{since } \mathbf{X} \cdot \mathbf{a} \mathbf{W}_{\text{LIN}} = \mathbf{ay}$$

$$\begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nn} \end{bmatrix} \times \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

all b 's comes from $(\mathbf{W}'_{\text{LIN}})$

$$\mathbf{X} \cdot \mathbf{W}'_{\text{LIN}} = \mathbf{ay} + b \quad \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nn} \end{bmatrix} \times \begin{bmatrix} aw_1 + b \\ aw_2 + b \\ \vdots \\ aw_n + b \end{bmatrix} = \begin{bmatrix} ay_1 + b \\ ay_2 + b \\ \vdots \\ ay_n + b \end{bmatrix}$$

$$\downarrow \quad a \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} + b \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$\Rightarrow \mathbf{W}'_{\text{LIN}} = a \mathbf{W}_{\text{LIN}} + b \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}, \text{ where } \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \text{ is a } n \times 1 \text{ matrix}.$$

6. (20 points) Let $E(\mathbf{w}) : \mathbb{R}^d \rightarrow \mathbb{R}$ be a function. Denote the gradient $\mathbf{b}_E(\mathbf{w})$ and the Hessian $\mathbf{A}_E(\mathbf{w})$ by

$$\mathbf{b}_E(\mathbf{w}) = \nabla E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E}{\partial w_1}(\mathbf{w}) \\ \vdots \\ \frac{\partial E}{\partial w_d}(\mathbf{w}) \end{bmatrix}_{d \times 1} \quad \text{and } \mathbf{A}_E(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 E}{\partial w_1^2}(\mathbf{w}) & \frac{\partial^2 E}{\partial w_1 \partial w_2}(\mathbf{w}) & \cdots & \frac{\partial^2 E}{\partial w_1 \partial w_d}(\mathbf{w}) \\ \frac{\partial^2 E}{\partial w_2 \partial w_1}(\mathbf{w}) & \frac{\partial^2 E}{\partial w_2^2}(\mathbf{w}) & \cdots & \frac{\partial^2 E}{\partial w_2 \partial w_d}(\mathbf{w}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 E}{\partial w_d \partial w_1}(\mathbf{w}) & \frac{\partial^2 E}{\partial w_d \partial w_2}(\mathbf{w}) & \cdots & \frac{\partial^2 E}{\partial w_d^2}(\mathbf{w}) \end{bmatrix}_{d \times d}.$$

Then, the second-order Taylor's expansion of $E(\mathbf{w})$ around \mathbf{u} is:

$$E(\mathbf{w}) \approx E(\mathbf{u}) + \mathbf{b}_E(\mathbf{u})^T(\mathbf{w} - \mathbf{u}) + \frac{1}{2}(\mathbf{w} - \mathbf{u})^T \mathbf{A}_E(\mathbf{u})(\mathbf{w} - \mathbf{u}).$$

Suppose $\mathbf{A}_E(\mathbf{u})$ is positive definite. The optimal direction \mathbf{v} such that $\mathbf{w} \leftarrow \mathbf{u} + \mathbf{v}$ minimizes the right-hand-side of the Taylor's expansion above is simply $-(\mathbf{A}_E(\mathbf{u}))^{-1}\mathbf{b}_E(\mathbf{u})$.

[Hint: Homework 07 -]

Now, consider minimizing $E_{\text{in}}(\mathbf{w})$ in logistic regression problem with Newton's method on a data set $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ with the cross-entropy error function for E_{in} :

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)).$$

$$\subset -\frac{1}{N} \sum_{n=1}^N \ln \left(1 + \frac{1}{h(\mathbf{x}_n)} \right)$$

For any given \mathbf{w}_t , let

$$h_t(\mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}_t^T \mathbf{x})}.$$

Express the Hessian $\mathbf{A}_E(\mathbf{w}_t)$ with $E = E_{\text{in}}$ as $\mathbf{X}^T \mathbf{D} \mathbf{X}$, where \mathbf{D} is an N by N diagonal matrix. Derive what \mathbf{D} should be in terms of h_t , \mathbf{w}_t , \mathbf{x}_n , and y_n .

$\frac{\partial^2}{\partial x_i^2}$

$$\begin{aligned} h_t(x) &= \frac{1}{1 + \exp(\mathbf{w}_t^T \mathbf{x})}, & E_{\text{in}} &= \frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y_n \mathbf{w}_t^T \mathbf{x}_n)) \\ h_t(-x) &= \frac{1}{1 + \exp(-\mathbf{w}_t^T \mathbf{x})} = 1 - h_t(x), & & \\ b_E(w) &= \nabla E_{\text{in}} = \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{x}_n)} (-y_n \mathbf{x}_n) = \frac{1}{N} \sum_{n=1}^N (1 - h_t(y_n \mathbf{x}_n)) (-y_n \mathbf{x}_n), & & \\ A_E(w) &= \nabla b_E(w) = \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n)(-y_n \mathbf{x}_n)(1 + \exp(-y_n \mathbf{x}_n)) - \exp^2(-y_n \mathbf{x}_n)(-y_n \mathbf{x}_n)}{(1 + \exp(-y_n \mathbf{x}_n))^2} (-y_n \mathbf{x}_n), & & (-y_n \mathbf{x}_n) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n)(-y_n \mathbf{x}_n)(-y_n \mathbf{x}_n) [1 + \exp(-y_n \mathbf{x}_n) - \exp(-y_n \mathbf{x}_n)]}{(1 + \exp(-y_n \mathbf{x}_n))^2} (-y_n \mathbf{x}_n), & & (-y_n \mathbf{x}_n) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n)(-y_n \mathbf{x}_n)(-y_n \mathbf{x}_n)}{(1 + \exp(-y_n \mathbf{x}_n))^2} (-y_n \mathbf{x}_n) = \mathbf{x}^T \mathbf{p} \mathbf{x}. & & \end{aligned}$$

$$D_{\bar{i}, \bar{j}} = \frac{1}{N} h_t(\mathbf{x}_{\bar{i}} y_{\bar{j}}) (1 - h_t(\mathbf{x}_{\bar{i}} y_{\bar{j}})) \cdot y_{\bar{i}, \bar{j}} \cdot y_{\bar{i}, \bar{j}}$$

$$= \frac{1}{N} (1 - h_t(\mathbf{x}_{\bar{i}})) h_t(\mathbf{x}_{\bar{i}}) y_{\bar{i}, \bar{j}} \cdot y_{\bar{i}, \bar{j}}$$

$$\text{since } y_{\bar{i}, \bar{i}} \cdot y_{\bar{i}, \bar{i}} = 1 \quad (\text{D diagonal}), \quad D_{\bar{i}, \bar{j}} = \frac{1}{N} \begin{bmatrix} h_t(\mathbf{x}_{\bar{i}}) (1 - h_t(\mathbf{x}_{\bar{i}})) & & \\ & h_t(\mathbf{x}_{\bar{i}}) (1 - h_t(\mathbf{x}_{\bar{i}})) & \\ & & h_t(\mathbf{x}_{\bar{i}}) (1 - h_t(\mathbf{x}_{\bar{i}})) \end{bmatrix}$$

7. (20 points) The truncated squared loss

$$\text{err}(s, y) = (\max(0, 1 - ys))^2$$

can be easily shown to be an upper bound on the 0/1 error. Assume that s is generated from a linear scoring function $s = w^T x$ like Page 3/25 of Lecture 11. Derive a "perceptron learning algorithm" by applying SGD on the truncated squared loss. Compare the resulting algorithm with the original PLA. Discuss the similarities and differences using 5 to 10 sentences.

$$\begin{aligned} \text{PLA: } & \text{error}_{\text{PLA}}(s, y) = [\text{sign}(s) \neq y] \\ & W_{t+1} = W_t + \eta \cdot [y_n \neq \text{sign}(s)] (y_n x_n) \end{aligned}$$

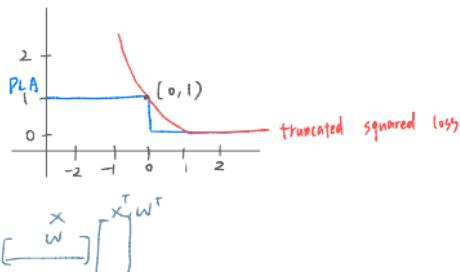
$$\begin{aligned} \text{SGD: } & \text{error}(s, y) = (\max(0, 1 - ys))^2 \\ & W_{t+1} = W_t + \eta \cdot (-\nabla \text{error}(y, w^T x)) \\ & = W_t + \eta \cdot (-\nabla (\max(0, 1 - ys))^2) \end{aligned}$$

• If $ys > 0$, $\text{error}(s, y) = 0$
 $0 < ys < 1$, $\text{error}(s, y) > \text{error}_{\text{PLA}}(s, y)$

$ys = 1$, $\text{error}(s, y) = \text{error}_{\text{PLA}}(s, y)$

$$\begin{aligned} (\nabla \max(0, 1 - ys))^2 &= \nabla (-(1 - w^T x_n))^2 = \nabla (w^T x_n - 1)^2 \\ &= \nabla (x \cdot w^T w \cdot x^T - 2x^T w y_n + 1)^2 \\ &= 2(x^T x \cdot w - y_n w) = 2(s - y_n)x = (1 - ys)x y_n \end{aligned}$$

$$W_{t+1} = W_t + 2\eta (1 - ys)x_n y_n$$



PLA always updates with $\eta = 1$ when $|y_n \neq \text{sign}(s)|$; while SGD updates with altered η

when $\max(0, 1 - ys) > 0$. When $ys > 1$, both PLA and SGD won't update ($\text{error} = 0$) .

when $0 < ys < 1$, SGD updates (PLA doesn't), when $ys < 0$, $W_{t+1}^{\text{PLA}} = W^{\text{PLA}} + \eta y_n x_n$,

$W_{t+1}^{\text{SGD}} = W_t^{\text{SGD}} + 2\eta (1 - ys)x_n y_n$ (updates depend on sign(y_n), alternate).

If $ys \ll 0$, updates larger; $ys \gg 0$, updates less.

However, eventually how much to update depend on η , which is usually smaller than 0.1

Consequently, SGD usually updates with flexible and smaller changes, PLA offers constant and unstable steps. In error measure, $\text{Err(SGD)} > \text{Err(PLA)}$, that is, SGD can also be viewed as upper bound of PLA.

Multinomial Logistic Regression

8. In Lecture 11, we solve multiclass classification by OVA or OVO decompositions. One alternative to deal with multiclass classification is to extend the original logistic regression model to Multinomial Logistic Regression (MLR). For a K -class classification problem, we will denote the output space $\mathcal{Y} = \{1, 2, \dots, K\}$. The hypotheses considered by MLR can be indexed by a matrix

$$W = \begin{bmatrix} | & | & \cdots & | & \cdots & | \\ w_1 & w_2 & \cdots & w_k & \cdots & w_K \\ | & | & \cdots & | & \cdots & | \\ & & & & & |_{(d+1) \times K} \end{bmatrix},$$

that contains weight vectors (w_1, \dots, w_K) , each of length $d+1$. The matrix represents a hypothesis

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

that can be used to approximate the target distribution $P(y|\mathbf{x})$ for any (\mathbf{x}, y) . MLR then seeks for the maximum likelihood solution over all such hypotheses. For a given data set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ generated i.i.d. from some $P(\mathbf{x})$ and target distribution $P(y|\mathbf{x})$, the likelihood of $h_y(\mathbf{x})$ is proportional to $\prod_{n=1}^N h_{y_n}(\mathbf{x}_n)$. That is, minimizing the negative log likelihood is equivalent to minimizing an $E_{\text{in}}(W)$ that is composed of the following error function

$$\text{err}(W, \mathbf{x}, y) = -\ln h_y(\mathbf{x}) = -\sum_{k=1}^K \|y = k\| \ln h_k(\mathbf{x}).$$

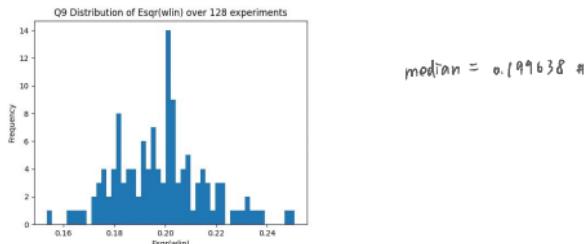
Consider minimizing $E_{\text{in}}(W) = \frac{1}{N} \sum_{n=1}^N \text{err}(W, \mathbf{x}_n, y_n)$ with gradient descent. Derive $\nabla E_{\text{in}}(W)$. Your result should simply be a matrix with the same size as W . (Note: the hypothesis that transforms the scores $(w_i^T \mathbf{x})_{i=1}^K$ to $h_y(\mathbf{x})$ is often called a softmax function in (multiclass) deep learning.)

$$\begin{aligned} E_{\text{in}}(W) &= \frac{1}{N} \sum_{n=1}^N \ln h_y(\mathbf{x}_n) \\ \nabla E_{\text{in}}(W) &= -\frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial W} \ln \frac{\exp(w_y^T \mathbf{x}_n)}{\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)} = -\frac{1}{N} \sum_{n=1}^N \frac{\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)}{\exp(w_y^T \mathbf{x}_n)} \times \frac{\exp(w_y^T \mathbf{x}_n) \cdot \frac{1}{N} \sum_{k=1}^K \exp(w_k^T \mathbf{x}_n) - \exp(w_y^T \mathbf{x}_n) \exp(w_k^T \mathbf{x}_n)}{\left(\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)\right)^2} \\ &= -\frac{1}{N} \sum_{n=1}^N \frac{\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)}{\left(\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)\right)^2} \times \frac{\exp(w_y^T \mathbf{x}_n) \cdot \frac{1}{N} \sum_{k=1}^K \exp(w_k^T \mathbf{x}_n) - \exp(w_y^T \mathbf{x}_n) \exp(w_k^T \mathbf{x}_n)}{\exp(w_y^T \mathbf{x}_n)} \\ &= -\frac{1}{N} \sum_{n=1}^N \frac{\left[\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n) - \exp(w_y^T \mathbf{x}_n) \right]}{\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)} \cdot \mathbf{x}_n = \frac{1}{N} \sum_{n=1}^N (h_t(x) - 1) \mathbf{x}_n \end{aligned}$$

$$\text{if } j \neq k, -\frac{1}{N} \sum_{n=1}^N \frac{\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)}{\exp(w_j^T \mathbf{x}_n)} \times \frac{-\exp(w_j^T \mathbf{x}_n) \exp(w_k^T \mathbf{x}_n)}{\sum_{k=1}^K \exp(w_k^T \mathbf{x}_n)^2} \cdot \mathbf{x}_n = \frac{1}{N} \sum_{n=1}^N h_t(x) \cdot \mathbf{x}_n$$

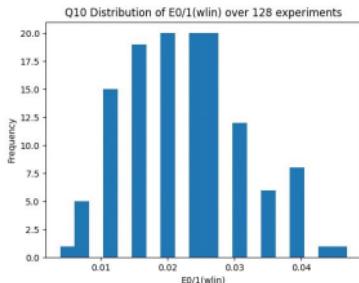
$$\rightarrow \nabla E_{\text{in}}(W) = -\frac{1}{N} \begin{bmatrix} | & | & \cdots & | & | \\ \sum_{n=1}^N (h_1(x_n) \mathbb{I}[y=1]) \mathbf{x}_n & \sum_{n=1}^N (h_2(x_n) - \mathbb{I}[y=2]) \mathbf{x}_n & \cdots & | \\ | & | & \cdots & | \\ & & & | \\ & & & | \\ & & & | \end{bmatrix}$$

9. (20 points, *) Implement the linear regression algorithm taught in the lecture. Run the algorithm for 128 times, each with a different random seed for generating the two data sets above. Plot a histogram to visualize the distribution of $E_{\text{in}}^{\text{sqr}}(\mathbf{w}_{\text{lin}})$, where $E_{\text{in}}^{\text{sqr}}$ denotes the averaged squared error over N examples. What is the median E_{in} over the 128 experiments?



10. (20 points, *) Following the previous problem, plot a histogram to visualize the distribution of $E_{\text{in}}^{0/1}(\mathbf{w}_{\text{LIN}})$, where $E_{\text{in}}^{0/1}$ denotes the averaged 0/1 error over N examples (i.e. using \mathbf{w}_{LIN} for binary classification). What is the median $E_{\text{in}}^{0/1}$ over the 128 experiments?

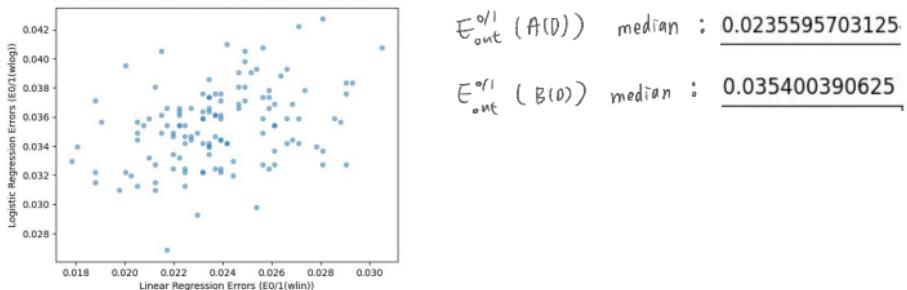
(Note: You can choose to run 128 new experiments in this problem, or just re-use the 128 hypotheses \mathbf{w}_{LIN} and test data sets obtained from the previous problem.)



median = 0.02343 #

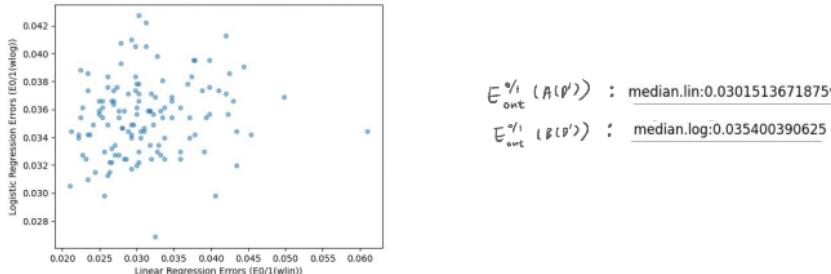
11. (20 points, *) Consider two algorithms. The first one, \mathcal{A} , is the linear regression algorithm above. The second one \mathcal{B} is logistic regression, trained with fixed learning rate gradient descent with $\eta = 0.1$ for $T = 500$ iterations, starting from $\mathbf{w}_0 = \mathbf{0}$. Run the algorithms on the same \mathcal{D} , and record $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D})), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}))]$. Repeat the process for 128 times, each with a different random seed for generating the training and test data sets above. Plot a scatter plot for $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D})), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}))]$. What is the median of $E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}))$ and what is the median of $E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}))$?

Q11 Scatter Plot of Linear vs. Logistic Regression Errors



12. (20 points, *) Following the previous problem, in addition to the 256 examples in \mathcal{D} , add 16 outlier examples generated from the following process to your training data (but not to your test data). All outlier examples will be labeled $y = +1$ and $\mathbf{x} = [1, x_1, x_2]$ where (x_1, x_2) comes from a normal distribution of mean $[0, 6]$ and covariance $\begin{bmatrix} 0.1 & 0 \\ 0 & 0.3 \end{bmatrix}$. Name the new training data set \mathcal{D}' . Run the algorithms on the same \mathcal{D}' , and record $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}')), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}'))]$. Repeat the process for 128 times, each with a different random seed for generating the training and test data sets above. Plot a scatter plot for $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}')), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}'))]$. What is the median of $E_{\text{out}}^0(\mathcal{A}(\mathcal{D}'))$ and what is the median of $E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}'))$? Compare your results to the previous problem. Describe your findings.

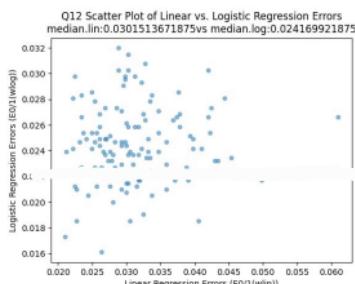
Q12 Scatter Plot of Linear vs. Logistic Regression Errors



Median of linear regression algorithm becomes larger than training without outlier data, while median of logistic algorithms doesn't change much.

Since linear regression make error $\hat{\sigma}^2$, when we use noise in training data, it may sacrifice original best g, modify itself to reduce the error from noise when training.

The plot also shows that $E_{\text{out}}(\text{log})$ doesn't increase with $E_{\text{out}}(\text{Lin})$, but in most of cases $E_{\text{out}}(\text{log})$ is bigger than $E_{\text{out}}(\text{Lin})$, so I run another plot to see if iteration=20000(which means it walks more steps to modify), $E_{\text{out}}'(\text{log}) : 0.0241$ the consequence indicates $E_{\text{out}}(\text{log})$ becomes smaller than $E_{\text{out}}(\text{Lin})$. That is, if we increase iteration, logistic have the ability to perform better than linear.



13. (Bonus 20 points) When using Newton's method for solving logistic regression, as discussed in Problem 6, each update \mathbf{v} is calculated by

$$\mathbf{v} = -(\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \nabla E_{\text{in}}(\mathbf{w}_t)$$

when $(\mathbf{X}^T \mathbf{D} \mathbf{X})$ is invertible. In linear regression, when $\mathbf{X}^T \mathbf{X}$ is invertible, the optimal

$$\mathbf{w}_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

If we can express $-\nabla E_{\text{in}}(\mathbf{w}_t)$ as some $\tilde{\mathbf{X}}^T \tilde{\mathbf{y}}$, and $\mathbf{X}^T \mathbf{D} \mathbf{X}$ as some $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$, then each iteration of Newton-solving logistic regression is performing an internal linear regression! State the internal linear regression problem—in particular, what are $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{y}}$?

$$13. \quad h_t(x) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}, \quad \mathbf{D} = \begin{bmatrix} h_t(x_1)(1-h_t(x_1)) & & & \\ & h_t(x_2)(1-h_t(x_2)) & & \\ & & \ddots & \\ & & & h_t(x_n)(1-h_t(x_n)) \end{bmatrix}$$

$$\textcircled{1} \quad \mathbf{X}^T \mathbf{D} \mathbf{X} = \mathbf{X}^T \sqrt{\mathbf{D}} \cdot \sqrt{\mathbf{D}}^T \mathbf{X} \quad , \quad \tilde{\mathbf{X}} = \sqrt{\mathbf{D}} \mathbf{X} \quad \#$$

$$\textcircled{2} \quad -\nabla E_{\text{in}}(\mathbf{w}_t) = \frac{1}{N} \sum_{n=1}^N y_n \mathbf{x}_n (1-h_t(x_n)) = \mathbf{X}^T [y(1-h_t(x))]$$

$$\frac{1}{N} (\mathbf{X}^T \mathbf{D}) (\sqrt{\mathbf{D}}^{-1} y (1-h_t(x))) = \tilde{\mathbf{X}}^T \tilde{y} \quad , \quad \tilde{\mathbf{X}}^T \tilde{y} = \mathbf{x}^T \sqrt{\mathbf{D}}^{-1} \tilde{y} = \mathbf{X}^T \sqrt{\mathbf{D}} \tilde{y}$$

$$\tilde{y} = \frac{1}{N} \sqrt{\mathbf{D}}^{-1} y (1-h_t(x)) \quad \#$$