# CS680 – Assignment 1

## Ali ElSaid (20745892)

E1.Q1:



E1.1: Mistakes per epoch for perceptron on spambase

E1.Q2:

```
# Training errors: 139
# Testing errors: 145
```

E1.Q3:

```
# Training errors: 111
# Testing errors: 148
```

E1.Q4:

Consider a term from the summation:

$$\max\{0, -y_i(\langle x_i, w\rangle)\} = f(w)$$

this term is
$$\begin{cases} 0, & \text{when } y_i(\langle x_i, w\rangle) \text{ is positive, ie. the prediction is right} \\ \\ -y_i(\langle x_i, w\rangle), & \text{when the prediction is wrong} \end{cases}$$

- So this term can be seen as the penalty assigned when a wrong prediction is made.

- Taking its derivative by the rule for $\frac{d}{dw} f(w) = \max_k f_k(w)$

then $\frac{d}{dw} f(w) = g(w) = \begin{cases} 0, & \text{when prediction correct} \\ \\ -y_i x_i, & \text{when prediction wrong} \end{cases}$

$\overset{\text{gradient}}{\underset{\vee}{}}$
comparing that to the ~~state~~ update the algorithm makes

we find that they are exactly the same; i.e. when the prediction is correct, no update happens. But when the prediction is wrong, the algorithm updates $w \leftarrow w - (-y_i x_i)$.

∴ we see that the update is $w \leftarrow w - g(w)$.

E1.Q5:

**E1.Q5** when $c=2$: consider the two cases for $y_i$ & the corresponding term inside the summation

when $y_i = 1 \longrightarrow$ term is: $\text{Max}\left[\langle x_i, w_1\rangle - \langle x_i, w_1\rangle, \langle x_i, w_2\rangle - \langle x_i, w_1\rangle\right]$

$$= \text{Max}\left\{0, \langle x_i, w_2\rangle - \langle x_i, w_1\rangle\right\}$$

when $y_i = 2 \longrightarrow$ term is: $\text{Max}\left\{\langle x_i, w_2\rangle - \langle x_i, w_2\rangle, \langle x_i, w_1\rangle - \langle x_i, w_2\rangle\right\}$

$$= \text{Max}\left\{0, \langle x_i, w_1\rangle - \langle x_i, w_2\rangle\right\}$$

Since we

This term should be equal to the similar term in the binary case

So, equating both we get:

& if we map $y_i = 1 \longrightarrow y_i^{(2)} = -1$
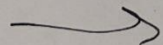
&

$y_i = 2 \longrightarrow y_i^{(2)} = 1$

$\underbrace{\phantom{xxxxx}}$   $\underbrace{\phantom{xxxxx}}$
when $c=2$          in binary
                    perceptron

then, when when $y_i = 1$ in our problem & $y_i^{(2)} = -1$ in the binary perceptron

$$\langle x_i, w_2\rangle - \langle x_i, w_1\rangle = -(\overset{y_i^{(2)}}{-1})(\langle x_i, w\rangle)$$

$\rightarrow \langle x_i, w_2 - w_1\rangle = \langle x_i, w\rangle$, which implies

$$W = W_2 - W_1$$

$\longrightarrow$

E1.Q6:

```
# Training errors: 126
# Testing errors: 132
```

E1.Q6 | Suppose we predict $\hat{y}_i$ ~~~~

So, by the definition of the prediction

$$\langle x_i, w_{\hat{y}_i} \rangle > \langle x_i, w_k \rangle \quad \forall k \neq \hat{y}_i$$

So, $\langle x_i, w_{\hat{y}_i} \rangle - \langle x_i, w_{y_i} \rangle > \langle x_i, w_k \rangle - \langle x_i, w_{y_i} \rangle \quad \forall k \neq \hat{y}_i$

Consider 2 cases:

Case 1: $\hat{y}_i ~~ = y_i$ , i.e. prediction is correct.

then $\langle x_i, w_{\hat{y}_i} \rangle - \langle x_i, w_{y_i} \rangle$

$= \langle x_i, w_{\hat{y}_i} \rangle - \langle x_i, w_{\hat{y}_i} \rangle = 0 > \langle x_i, w_k \rangle - \langle x_i, w_{\hat{y}_i} \rangle \quad \forall k \neq \hat{y}_i$

& So, the max term is equal to zero. ~~~

i.e. No penalty since we predicted correctly.

Case 2: $\hat{y}_i ~~ \neq y_i$ , i.e. wrong prediction.

& $\langle x_i, w_{\hat{y}_i} \rangle - \langle x_i, w_{y_i} \rangle > 0$ & we incur penalty for wrong prediction.

To get the algorithm & update, we find the gradient of the term $f(w)$ inside the summation. From the definition of $f'(w)$ when $f(w) = \max_k f_k(w)$ we get $\longrightarrow$

if correct prediction

$$\nabla_w 0 = 0 \quad, \text{ so no update}$$

if incorrect prediction

$$\frac{\partial}{\partial w_k} \langle x_i, w_{\hat{y_i}} \rangle - \langle x_i, w_{y_i} \rangle = \begin{cases} 0 \quad, \text{when } k \neq \hat{y_i} \ \& \ k \neq y_i \\ x_i \quad, \text{when } k = \hat{y_i} \\ -x_i \quad, \text{when } k = y_i \end{cases}$$

So we do the following update:

$$w_{\hat{y_i}} \leftarrow w_{\hat{y_i}} - x_i \quad \& \quad w_{y_i} \leftarrow w_{y_i} + x_i \qquad \begin{pmatrix} \text{Padded } w \\ \& \ x_i \end{pmatrix}$$

---

Algorithm:

Input: $X \in \mathbb{R}^{d \times n}$, $y \in \{1, 2, -, c\}^{\wedge}$,

$w = [w_1, ---, w_c] \ \cancel{\phantom{xx}} = 0_{d \times c}$, i.e. each $w_j = 0_d \ \forall j \in \{1, ---, c\}$

$b = [b_1, ---, b_c] = 0_c$, i.e. each $b_j = 0 \ \forall j \in \{1, ---, c\}$

Max_Pass $\in \mathbb{N}$

Output: $W \in \mathbb{R}^{d \times c}$, $b \in \mathbb{R}^c$, mistake

for $t = 1, 2, ---,$ Max_Pass do
    mistake $(t) \leftarrow 0$
    for $i = 1, 2, ---, n$ do
        $\hat{y_i} = \underset{k=1,--,c}{\arg\max} \langle x_i, w_k \rangle + b_k$
        if $\hat{y_i} \neq y_i$ then:
            $w_{\hat{y_i}} \leftarrow w_{\hat{y_i}} - x_i, \quad b_{\hat{y_i}} \leftarrow b_{\hat{y_i}} - 1$
            $w_{y_i} \leftarrow w_{y_i} + x_i, \quad b_{y_i} \leftarrow b_{y_i} + 1$
            mistake $(t) \leftarrow$ mistake $(t) + 1$

E2.Q1:

likelihood function: $P(y_1, -, y_n | x_1, \cdots, x_n) = \prod_{i=1}^{\hat{n}} P(y_i | x_i)$

$$L(y|x)$$

since we assume indep.

$\ell_{\hat{}}(y|x) = -\log L(y|x) = -\log \prod_{i=1}^{\hat{n}} P(y_i | x_i)$

$= -\sum_{i=1}^{\hat{n}} \log P(y_i | x_i) = -\sum_{i=1}^{\hat{n}} \log \left( \exp(M(x_i) y_i - \lambda(x_i)) q(y_i) \right)$

$\boxed{= -\sum_{i=1}^{\hat{n}} M(x_i) y_i - \lambda(x_i) + \log q(y_i)}$   where $\ell_n(y|x)$ is the negative log-likelihood function

$\boxed{\text{E2.Q2}}$  Plugging into $\ell_n$

__we get__ $\ell_n(w) = -\sum\limits_{i=1}^{\hat{}} w^T x_i \; y_i - \lambda(x_i) + \log q(y_i)$

the gradient

$$\nabla_w \ell_n(w) = -\sum\limits_{i=1}^{\hat{}} x_i y_i - \nabla_w \lambda(x_i)$$

& $\nabla_w \lambda(x_i) = \nabla_w \log \int_y \left(e^{w^T x_i \cdot y}\right) q(y) \, dy$ ,

$= \dfrac{\nabla_w \int_y \left(e^{w^T x_i \cdot y}\right) q(y) \, dy}{\int_y \exp(w^T x_i \cdot y) \, q(y) \, dy}$  , since $\frac{d}{dw} \log f(w) = \dfrac{\frac{d}{dw} f(w)}{f(w)}$

$= \dfrac{\int_y \nabla_w \exp(\cdots) q(y) dy}{\int_y \exp(\cdots) q(y) dy}$ , $(\cdots) = w^T x_i \cdot y$ ~~$\cancel{}$~~

$= \dfrac{x_i \int_y y \left(e^{w^T x_i \cdot y}\right) q(y) dy}{e^{\lambda(x_i)}}$ , since $\lambda(x_i) = \log \int_y e^{w^T x_i \cdot y} q(y) dy$

$= x_i \int_y y \, e^{w^T x_i y} \cdot e^{-\lambda(x_i)} \cdot q(y) \, dy = x_i \int_y y \, e^{w^T x_i y - \lambda(x_i)} q(y) dy$

$= x_i \int_y y \, P(y | x_i) \, dy = \boxed{x_i \cdot E(y | x_i)}$

So the gradient $\overset{\text{part V}}{\nabla_w \lambda(x_i)} = x_i \, E(y | x_i)$

& $\boxed{\nabla_w \ell_n(w) = -\sum\limits_{i=1}^{\hat{}} x_i (y_i - E(y | x_i))}$

E2.Q3:

## E2.Q3

$$y^2 - 2y\nu(x) + \nu(x)^2$$

$$\uparrow$$

$$P(y|x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - \nu(x))^2}{2}\right)$$

$$= \cancel{\exp(y\nu(x))} \frac{\exp\left(-\frac{y^2}{2}\right)}{\sqrt{2\pi}} \exp\left(y\underbrace{\nu(x)}_{\boxed{M(x)}} - \underbrace{\frac{\nu(x)^2}{2}}_{\boxed{A(x)}}\right)$$

$$\underbrace{\frac{\exp\left(-\frac{y^2}{2}\right)}{\sqrt{2\pi}}}_{\boxed{q(y)}}$$

Plugging into $\ell_n(w)$

$$\boxed{\ell_n(w)} = -\sum_{i=1}^{n} w^T x_i \cdot y_i - \frac{(w^T x_i)^2}{2} + \log\left(\frac{1}{\sqrt{2\pi}}\right) - \frac{y^2}{2}$$

$$\boxed{\nabla_w \ell_n(w)} = -\sum_{i=1}^{n} x_i (y_i) - x_i^T w x_i$$

$$= -\sum_{i=1}^{n} x_i (y_i - x_i^T w)$$

E2.Q4:

## E2.Q4

$$P(Y=y \mid x) = [\nu(x)]^{y}[1-\nu(x)]^{1-y}, \quad \text{taking } \exp(\log(\ ))$$

$$\rightarrow = \exp\Big(y \cdot \underbrace{\log\Big(\frac{\nu(x)}{1-\nu(x)}\Big)}_{M(x)} \; \underbrace{- \; -\log(1-\nu(x))}_{\lambda(x)}\Big)$$

$$\& \; q(y) = 1$$

---

$$M(x) = \log\Big(\frac{\nu(x)}{1-\nu(x)}\Big) \qquad \& \; M(x) = w^{T}x$$

$$\text{So } \nu(x) = \frac{e^{M(x)}}{1+e^{M(x)}} = \frac{e^{w^{T}x}}{1+e^{w^{T}x}}$$

Plugging into $\lambda(x) = -\big(\log(1-\nu(x))\big) = -\Big(\log\Big(1- \frac{e^{w^{T}x}}{1+e^{w^{T}x}}\Big)\Big)$

$$= \log(1+e^{w^{T}x})$$

$$\& \; \ell_{n}(w) = -\sum_{i=1}^{\hat{n}} w^{T}x_{i}y_{i} - \log(1+e^{w^{T}x})$$

$$\& \; \nabla_{w}\ell_{n}(w) = -\sum_{i=1}^{\hat{n}} x_{i}y_{i} - \frac{e^{w^{T}x}}{1+e^{w^{T}x}}\cdot x_{i} = -\sum_{i=1}^{\hat{n}} x_{i}(y_{i} - \nu(x))$$

E2.Q5:

## E2.Q5

$$P(Y=y \mid x) = \frac{[v(x)]^y}{y!} e^{-v(x)} \quad , \quad \text{taking } \exp(\log(\ ))$$

$$\rightarrow = \exp\left( \underbrace{y \cdot \log v(x)}_{\mu(x)} - \underbrace{v(x)}_{\boxed{\lambda(x)}} \right) \cdot \underbrace{\frac{1}{y!}}_{\boxed{q(y)}}$$

E3.Q1:

<u>E3.Q1</u> If we treat it as normal multiclassification,
then we will count any wrong Prediction to be the same
no matter what. e.g. in the letter grade example
if the actual grade is A , we would like to Penalize more
for a prediction of C than B since B is "closer" to
A because of the ordering. However that will not happen
if we to normal multiclassification.

**E3.Q2** if we just encode the labels as $\{1, 2, \ldots c\}$ and run normal regression, we will basically be telling the model that the quantitative difference between label 1 & 2 is the same as between 2 & 3 for example. i.e. the predictors have to change equivalenty to go from 1 to 2 & from 2 to 3.

And this assumption will not be a sound one often in real life, since the quantitative difference between labels might actually be different.

### E3.Q3

replacing the max with the dual variable

$$\to \underset{w,\, b_1 \le b_2 \le \cdots \le b_{c-1},\; 0 \le \alpha \le C}{\text{Min}} \; \text{Max} \quad \overbrace{\frac{1}{2}\|w\|_2^2 + \sum_{k=1}^{c-1}\sum_{i=1}^{\hat{}} \alpha_{ki}\left[1 - \left(\llbracket y_i \ge k \rrbracket - \llbracket y_i \ge k+1 \rrbracket\right)\left(\langle x_i, w\rangle + b_k\right)\right]}^{f(w,b)}$$

~~differentiating~~ switching Min with Max & differentiating we get.

$$\frac{\partial}{\partial w} f(w,b) = w - \overbrace{\sum_{k=1}^{c-1}\sum_{i=1}^{\hat{}} \alpha_{ki}\left(\llbracket y_i \ge k \rrbracket - \llbracket y_i \ge k+1 \rrbracket\right) x_i}^{} = 0$$

$$\to w = \underline{\qquad\qquad} \downarrow$$

$$\frac{\partial}{\partial b_j} = 0 + \sum_{i=1}^{n} \alpha_{ji}\left(\llbracket\;\rrbracket - \llbracket\;\rrbracket\right) = 0, \; \forall j \in \{1, \ldots, c-1\}$$

Plugging back into Primal Problem. we get the dual Problem

**E3.Q4** In the original formulation the penalty is

1 if $y_i \neq k$ and $y_i \neq k+1$ for a given point

however a better formulation would be ~~to just~~ ~~all the points~~, for a given point, to punish all the planes represented by ~~the~~ $b_k$'s that are supposed to be

"above" it if $k \leq y_i - 1$, since all planes that

have $y_i \geq k+1 \Leftrightarrow k \leq y_i - 1$ should be "below" the point

& vice versa for $y_i \leq k$, the planes should be "above" it

So the formulation would be:

$$\min_{w, b_1 \leq b_2 \leq \cdots \leq b_{c-1}} \frac{\lambda}{2} \|w\|_2^2 + \sum_{k=1}^{c-1} \sum_{i=1}^{n} \max\left\{0, 1 - \left(\left[\!\left[y_i \leq k\right]\!\right] - \left[\!\left[y \geq k+1\right]\!\right]\right)\left(\langle x_i, w \rangle + b_k\right)\right\}$$

, this will punish all the planes in the wrong place

for any given point.