# Machine Learning- Exercise 2
# Theoretical Part

Name:   Tomer Gill (תומר גיל)
ID:   318459450
U2 Username:   gilltom
Date:   31/03/2018

## Contents

## Multiclass Logistic Regression

(a) $P(Y = y|X = x) = softmax(\boldsymbol{W}\boldsymbol{x} + \boldsymbol{b})_{[y]} = \frac{e^{\boldsymbol{W}_{[y]}\boldsymbol{x} + b_{[y]}}}{\sum_i e^{\boldsymbol{W}_{[i]}\boldsymbol{x} + b_{[i]}}}$, where $[i]$ indicates the i-th index of the vector / matrix, $\boldsymbol{W}$ the weights matrix and $\boldsymbol{b}$ the bias vector.

(b) $\Theta^* = \underset{\Theta}{argmin} \frac{1}{m} \sum_{t=1}^{m} - \log(P(Y_t = y|X_t = x)) =$

$\underset{\Theta}{argmin} \sum_t - \log(softmax(\boldsymbol{W}\boldsymbol{x} + \boldsymbol{b})_{[y]}) \overset{w := W_{[y]}, b := b_{[y]}}{\triangleq}$

$\underset{\Theta}{argmin} \sum_t - \log\left(\frac{e^{wx+b}}{\sum_i e^{W_{[i]}x + b_{[i]}}}\right) = \underset{\Theta}{argmin} \sum_t -\boldsymbol{w}\boldsymbol{x} - b + \log(\sum_i \boldsymbol{W}_{[i]}\boldsymbol{x} + \boldsymbol{b}_{[i]})$,

where $\Theta$ is $\boldsymbol{W}$ and $\boldsymbol{b}$.

(c) I'll calculate the gradients of the average loss for all the training examples w.r.t each weight $(w_1, \dots, w_n - n \text{ is the number of classes})$ and the bias $(b)$:

**Weights:** $\nabla_{w_i} loss(\boldsymbol{W}; \boldsymbol{b}) = \nabla_{w_i} \frac{1}{m} \sum_{t=1}^{m} - \log(P(Y_t = y|X_t = x)) =$

$\frac{1}{m} \sum_{t=1}^{m} \nabla_{w_i} - \log(P(Y_t = y|X_t = x)) = \frac{1}{m} \sum_{t=1}^{m} \nabla_{w_i} - \log\left(\frac{e^{w_y x + b_y}}{\sum_{j=1}^{n} e^{W_j x + b_j}}\right) =$

$\frac{1}{m} \sum_{t=1}^{m} -\nabla_{w_i} \log(e^{w_y x + b_y}) + \nabla_{w_i} \log(\sum_{j=1}^{n} e^{W_j x + b_j}) =$

$\frac{1}{m} \sum_{t=1}^{m} -\nabla_{w_i}(\boldsymbol{w}_y \boldsymbol{x} + b_y) + \frac{1}{\sum_{j=1}^{n} e^{W_j x + b_j}} \nabla_{w_i}(\sum_{j=1}^{n} e^{W_j x + b_j})$

We'll call the derivative of the loss of example t as $\nabla_{w_i} loss_t(\boldsymbol{W}; \boldsymbol{b}) =$

$-\nabla_{w_i}(\boldsymbol{w}_y \boldsymbol{x} + b_y) + \frac{1}{\sum_{j=1}^{n} e^{W_j x + b_j}} \nabla_{w_i}(\sum_{j=1}^{n} e^{W_j x + b_j})$, so in total $\nabla_{w_i} loss(\boldsymbol{W}; \boldsymbol{b}) =$

$\frac{1}{m} \sum_{t=1}^{m} loss_t(\boldsymbol{W}; \boldsymbol{b})$. For $i = y$ (where y is the correct tag), $loss_t(\boldsymbol{W}; \boldsymbol{b}) =$

$-\nabla_{w_y}(\boldsymbol{w}_y \boldsymbol{x} + b_y) + \frac{1}{\sum_{j=1}^{n} e^{W_j x + b_j}} \nabla_{w_y}(\sum_{j=1}^{n} e^{W_j x + b_j}) = -\boldsymbol{x} +$

$\frac{1}{\sum_{j=1}^{n} e^{W_j x + b_j}} (\sum_{j=1}^{n} \nabla_{w_y} e^{W_j x + b_j}) = -\boldsymbol{x} + \frac{1}{\sum_{j=1}^{n} e^{W_j x + b_j}} \left(\nabla_{w_y} e^{W_0 x + b_0} + \right.$

$\left. \dots \nabla_{w_y} e^{W_y x + b_y} + \dots \nabla_{w_y} e^{W_n x + b_n}\right) = -\boldsymbol{x} + \frac{1}{\sum_{j=1}^{n} e^{W_j x + b_j}} (e^{W_y x + b_y} \boldsymbol{x}) =$

$$x\left(\frac{e^{W_y x+b_y}}{\sum_{j=1}^n e^{W_j x+b_j}}-1\right)=x\left(softmax_{[y]}(Wx+b)-1\right).$$ For $i\neq y$: $loss_t(W;b)=$

$$-\nabla_{w_i}(w_y x+b_y)+\frac{1}{\sum_{j=1}^n e^{W_j x+b_j}}\nabla_{w_i}\left(\sum_{j=1}^n e^{W_j x+b_j}\right)=-0+$$

$$\frac{1}{\sum_{j=1}^n e^{W_j x+b_j}}\left(\sum_{j=1}^n \nabla_{w_i}e^{W_j x+b_j}\right)=\frac{e^{W_i x+b_i}}{\sum_{j=1}^n e^{W_j x+b_j}}x=softmax(Wx+b)_{[i]}*x.$$

In conclusion, $\nabla_{w_i}loss(W;b)=\frac{1}{m}\sum_{t=1}^m \nabla_{w_i}loss_t(W;b),$

where $\nabla_{w_i}loss_t(W;b)=\begin{cases}-x+\dfrac{e^{W_y x+b_y}}{\sum_{j=1}^n e^{W_j x+b_j}}x\ ,i=y\\[2ex]\dfrac{e^{W_i x+b_i}}{\sum_{j=1}^n e^{W_j x+b_j}}x\qquad,i\neq y\end{cases}.$

**Bias**: $\nabla_{b_i}loss(W;b)=\nabla_{b_i}\frac{1}{m}\sum_{t=1}^m -\log\left(P(Y_t=y|X_t=x)\right)=\frac{1}{m}\sum_{t=1}^m \nabla_{b_i}-$

$$\log\left(P(Y_t=y|X_t=x)\right)=\frac{1}{m}\sum_{t=1}^m \nabla_{b_i}-\log\left(\frac{e^{W_y x+b_y}}{\sum_{j=1}^n e^{W_j x+b_j}}\right)=$$

$$\frac{1}{m}\sum_{t=1}^m -\nabla_{b_i}\log\left(e^{W_y x+b_y}\right)+\nabla_{b_i}\log\left(\sum_{j=1}^n e^{W_j x+b_j}\right)=$$

$$\frac{1}{m}\sum_{t=1}^m -\nabla_{b_i}(w_y x+b_y)+\frac{1}{\sum_{j=1}^n e^{W_j x+b_j}}\nabla_{b_i}\left(\sum_{j=1}^n e^{W_j x+b_j}\right),$$ so like before we'll

look at the case $i=y$ for a single example t:

$$-\nabla_{b_y}(w_y x+b_y)+\frac{1}{\sum_{j=1}^n e^{W_j x+b_j}}\nabla_{b_y}\left(\sum_{j=1}^n e^{W_j x+b_j}\right)=-1+$$

$$\frac{1}{\sum_{j=1}^n e^{W_j x+b_j}}\left(\nabla_{b_y}e^{W_0 x+b_0}+\cdots\nabla_{b_y}e^{W_y x+b_y}+\cdots\nabla_{b_y}e^{W_n x+b_n}\right)=-1+$$

$$\frac{e^{W_y x+b_y}*1}{\sum_{j=1}^n e^{W_j x+b_j}}=softmax(Wx+b)_{[y]}-1,$$ and for $i\neq y$ it's the same without the -1.

For conclusion, $\nabla_{b_i}loss(W;b)=\frac{1}{m}\sum_{t=1}^m \nabla_{b_i}loss_t(W;b),$

where $\nabla_{b_i}loss_t(W;b)=\begin{cases}\dfrac{e^{W_y x+b_y}}{\sum_{j=1}^n e^{W_j x+b_j}}-1\ ,i=y\\[2ex]\dfrac{e^{W_i x+b_i}}{\sum_{j=1}^n e^{W_j x+b_j}}\qquad,i\neq y\end{cases}.$

Now, for the update rules: the general SGD update rule is $w_i=\eta\nabla_{w_i}loss(W;b)$ for each row $w_i$ in the weight matrix $W$, and $b_i=\eta\nabla_{b_i}loss(W;b)$ for each element $b_i$ in the bias vector $b$. So the update rules are:

$$w_i=\begin{cases}-\eta x+\eta\dfrac{e^{W_y x+b_y}}{\sum_{j=1}^n e^{W_j x+b_j}}x\ ,i=y\\[2ex]\eta\dfrac{e^{W_i x+b_i}}{\sum_{j=1}^n e^{W_j x+b_j}}x\qquad,i\neq y\end{cases}$$

$$b_i=\begin{cases}\eta\dfrac{e^{W_y x+b_y}}{\sum_{j=1}^n e^{W_j x+b_j}}-\eta\ ,i=y\\[2ex]\eta\dfrac{e^{W_i x+b_i}}{\sum_{j=1}^n e^{W_j x+b_j}}\qquad,i\neq y\end{cases}$$

# Practical Part – Graph



Learned P(Y=1|X) VS. Probability by Density