

CSE512 Fall 2019 - Machine Learning - Homework 3

October 7, 2019

Name: Anmol Shukla

Solar ID: 112551470

NetID: anmshukla

Email address: Anmol.shukla@stonybrook.edu

Names of people whom I have discussed the homework with: None

1 Question 1

1.1.1

1.1.1. Bayes Risk is given by : $R(a_i/x) = \sum_j L(a_i/c_j) P(c_j|x)$

where L is the cost of taking action a_i given a class c_j

$$\therefore R(Y=1|x) = \underbrace{L(Y=1|\hat{Y}=1) P(\hat{Y}=1|x)}_{\text{this term is 0 as there is no risk/loss when prediction is correct}} + L(Y=1|\hat{Y}=0) P(\hat{Y}=0|x)$$

$$= \underbrace{L(Y=1|\hat{Y}=0) P(\hat{Y}=0|x)}_{\substack{\text{cost of false positive} \\ \rightarrow \text{prior of } \cancel{\text{positive}} \text{ negative}}} = \alpha \cdot (1 - \eta(x))$$

$$\begin{aligned} R(Y=0|x) &= L(Y=0|\hat{Y}=1) P(\hat{Y}=1|x) + \underbrace{L(Y=0|\hat{Y}=0) P(\hat{Y}=0|x)}_{\substack{0 \text{ as this is a right prediction}}} \\ &= L(Y=0|\hat{Y}=1) \cdot P(\hat{Y}=1|x) \\ &= 1 \cdot \eta(x) \Rightarrow \eta(x) \end{aligned}$$

Optimal bayes classifier minimizes the risk

$$\therefore \gamma^*(x) = \min(\eta(x), \alpha(1 - \eta(x)))$$



Scanned with
CamScanner

1.1.2

1.1.2. Following from 1.1.1

total risk for a data point x whose nearest neighbor is " z " = $r(x)$

$r(x)$ = error we make when we make a wrong decision

a wrong decision is made when $\text{label}(x) \neq \text{label}(z)$

$$r(x) = \underbrace{L(y=0|\hat{y}=1)}_1 P(\hat{y}=1|x) P(\hat{y}=0|z) + \underbrace{L(y=1|\hat{y}=1)}_\alpha P(\hat{y}=0|x) P(\hat{y}=1|z)$$

$$r(x) = \eta(x) (1 - \eta(z)) + \alpha (1 - \eta(x)) \eta(z)$$

but as $n \rightarrow \infty$ (no. of training data)
 $x \rightarrow z$ or $z \rightarrow x$

$$\therefore r(x) = \eta(x) (1 - \eta(x)) + \alpha \eta(x) (1 - \eta(x))$$

$$r(x) = (1 + \alpha) \eta(x) (1 - \eta(x))$$



Scanned with
CamScanner

1.1.3

1.1.3. To prove $r(x) \leq (1+\alpha) r^*(x) (1-r^*(x))$

from previous question, we know that

$$r(x) = (1+\alpha) \eta(x) (1-\eta(x))$$

$$\&, \quad r^*(x) = \min(\eta(x), \alpha(1-\eta(x)))$$

$$\therefore r^*(x) = \eta(x) \quad \text{if } \eta(x) \leq \alpha(1-\eta(x)) \quad \text{--- (1)}$$

$$= \alpha(1-\eta(x)) \quad \eta(x) > \alpha(1-\eta(x)) \quad \text{--- (2)}$$

We can prove the lower bound on $r^*(x)$ when $r^*(x) = \eta(x)$
from eq (1), $r^*(x) = \eta(x)$

$$\therefore r(x) = (1+\alpha) r^*(x) (1-r^*(x))$$

Now, we need to prove: $r(x) < (1+\alpha) r^*(x) (1-r^*(x))$

$$r(x) < (1+\alpha) \alpha(1-\eta(x)) (1-\alpha(1-\eta(x)))$$

subs. the value of $r(x)$

$$(1+\alpha) \eta(x) (1-\eta(x)) < (1+\alpha) \alpha(1-\eta(x)) (1-\alpha(1-\eta(x)))$$

$$\text{To prove: } \eta(x) < \alpha(1-\alpha(1-\eta(x))) \quad \text{--- (1)}$$

Since $r^*(x) = \alpha(1-\eta(x))$ in the above case
it follows that $\eta(x) > \alpha(1-\eta(x))$

$$\text{Also } \alpha > 1 \therefore -\eta(x) < -\alpha(1-\eta(x))$$

$$\therefore \alpha(1-\eta(x)) > \eta(x)$$

$$1-\eta(x) \leq 1-\alpha(1-\eta(x))$$

$$\alpha > 1$$

$$\therefore 1-\eta(x) < \alpha(1-\alpha(1-\eta(x)))$$

$$\therefore \eta(x) < \alpha(1-\alpha(1-\eta(x)))$$

\therefore We have proved eq (1)



1.1.4

1.1.4. from 1.1.3 we know that $r(x) \leq (1+\alpha)r^*(x)(1-r^*(x))$

taking expectation on both the sides

$$\begin{aligned} E[r(x)] &\leq (1+\alpha) E[r^*(x)(1-r^*(x))] \\ &\leq (1+\alpha) (E[r^*(x)] - E[r^{*2}(x)]) \quad \text{--- (1)} \end{aligned}$$

Now, we know that $E[x^2] = E[x]^2 + \text{Var}[x]$

$$\therefore E[r^{*2}(x)] = E[r^*(x)]^2 + \text{Var}[r^*(x)]$$

$$E[r^{*2}(x)] \geq E[r^*(x)]^2$$

$$\therefore R \leq (1+\alpha) E[r^*(x)] - (1+\alpha) E[r^*(x)]^2 \quad \text{from (1)}$$

$$R \leq (1+\alpha) R^* - (1+\alpha) R^{*2} \quad \left[\text{as } E[r^*(x)] = R^* \right]$$

$$R \leq (1+\alpha) R^* (1-R^*)$$

Hence Proved.



Scanned with
CamScanner

1.2.1

1.2.1. probability of a point being positive = η
 probability that at least $(K+1)/2$ out of K points are +ve.

$r(x)$ = asymptotic risk $r(x)$ for a point x in terms of $\eta(x)$
 & function $g(\cdot, \cdot)$

$$r(x) = \text{Prob}(x \text{ is positive but } \text{all } \frac{K+1}{2} \text{ points are not positive}) \\
+ \text{Prob}(x \text{ is negative but all } \frac{K+1}{2} \text{ points are +ve})$$

$$= \eta(x)(1 - g(\eta, K)) + (1 - \eta(x))g(\eta, K)$$

$$= \eta(x) - \eta(x)g(\eta, K) + g(\eta, K) - \eta(x)g(\eta, K)$$

$$= \eta(x) + g(\eta, K) - 2\eta(x)g(\eta, K)$$

$$= \eta(x) + (1 - 2\eta(x))g(\eta, K)$$

$$= \eta(x) + (1 - 2\eta(x))g(\eta, K)$$

2.1

2.1. To Prove: $\frac{\partial (P(y^i | \bar{x}^i; \theta))}{\partial \theta_c} = (\delta(c=y^i) - P(c | \bar{x}^i; \theta)) \bar{x}^i$

we can write

$$\log(P(y^i | \bar{x}^i; \theta)) = \delta(c=y^i) \log(P(y^i | \bar{x}^i; \theta)) + (1 - \delta(c=y^i)) \log(P(y^i | \bar{x}^i; \theta))$$

for simplicity, we'll write $\delta(c=y^i)$ as δ & \sum_j as a matrix multiplication

also,

$$P(y=i | x; \theta) = \frac{\exp(\theta_i^T \bar{x})}{1 + \sum_{j=1}^{K-1} \exp(\theta_j^T \bar{x})}$$

$$P(y=K | x; \theta) = \frac{1}{1 + \sum_{j=1}^{K-1} \exp(\theta_j^T \bar{x})}$$

$$\log(P(y | x; \theta)) = \delta \log(P(y=i | x; \theta)) + (1 - \delta) \log(P(y=K | x; \theta))$$

$$= \delta \log\left(\frac{\exp(\theta^T \bar{x})}{1 + \exp(\theta^T \bar{x})}\right) + (1 - \delta) \log\left(\frac{1}{1 + \exp(\theta^T \bar{x})}\right)$$

$$\log(P(y | \bar{x}; \theta)) = \delta \theta^T \bar{x} - \delta \log(1 + \exp(\theta^T \bar{x})) + (1 - \delta) \log(1 + \exp(\theta^T \bar{x}))$$

$$\frac{\partial \log(P(y | \bar{x}; \theta))}{\partial \theta} = \delta \bar{x} - \frac{\delta \cdot \bar{x} \cdot \exp(\theta^T \bar{x})}{1 + \exp(\theta^T \bar{x})} + \frac{\delta \bar{x} \cdot \exp(\theta^T \bar{x})}{1 + \exp(\theta^T \bar{x})} - \frac{\bar{x} \cdot \exp(\theta^T \bar{x})}{1 + \exp(\theta^T \bar{x})}$$

$$= \delta \bar{x} - \frac{\bar{x} \cdot \exp(\theta^T \bar{x})}{1 + \exp(\theta^T \bar{x})}$$

$$\text{but } P(y=c | x; \theta) = \frac{\exp(\theta^T \bar{x})}{1 + \exp(\theta^T \bar{x})}$$

$$\therefore \frac{\partial \log(P(y | \bar{x}; \theta))}{\partial \theta} = \delta \bar{x} - \bar{x} P(c | \bar{x}; \theta)$$

$$= (\delta(c=y) - P(c | \bar{x}; \theta)) \bar{x}$$



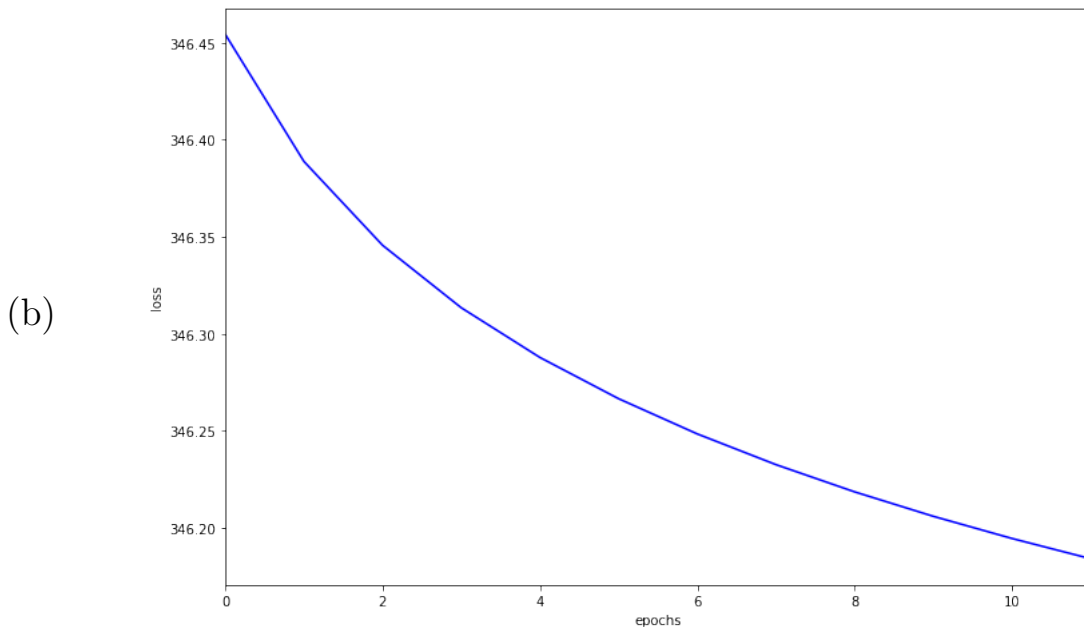
Scanned with
CamScanner

2 Question 2

2.3 Implement Logistic Regression with SGD

1. $\text{max_epoch} = 1000$, $m = 16$, $\eta_0 = 0.1$, $\eta_1 = 1$, $\delta = 0.00001$

(a) Number of epochs taken before existing = 12

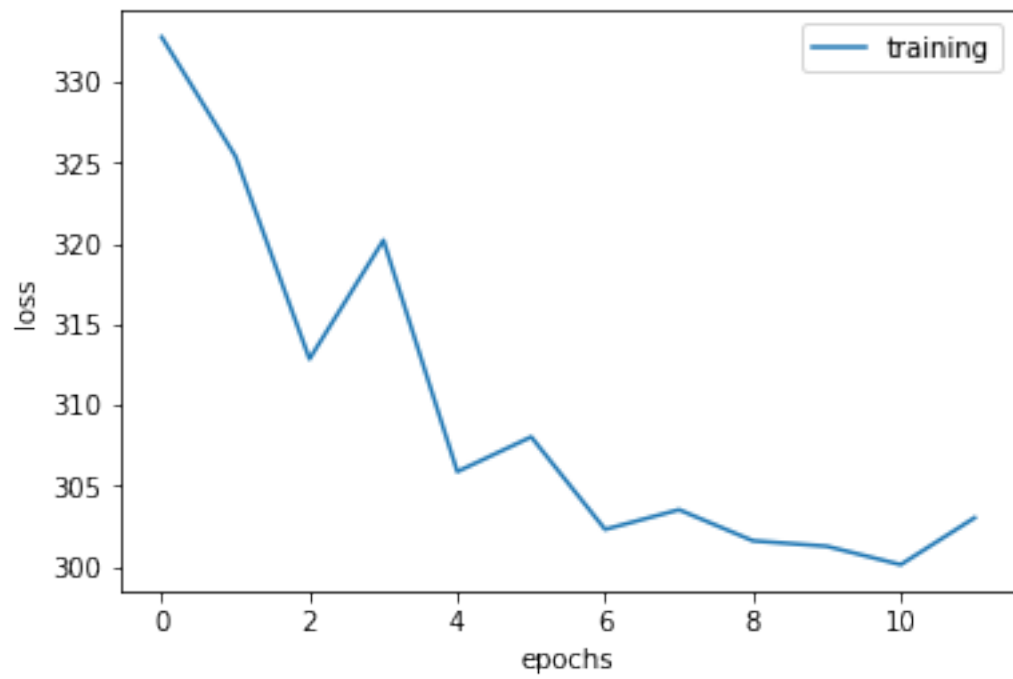


(c) $L(\theta) = 346.14301852200146$

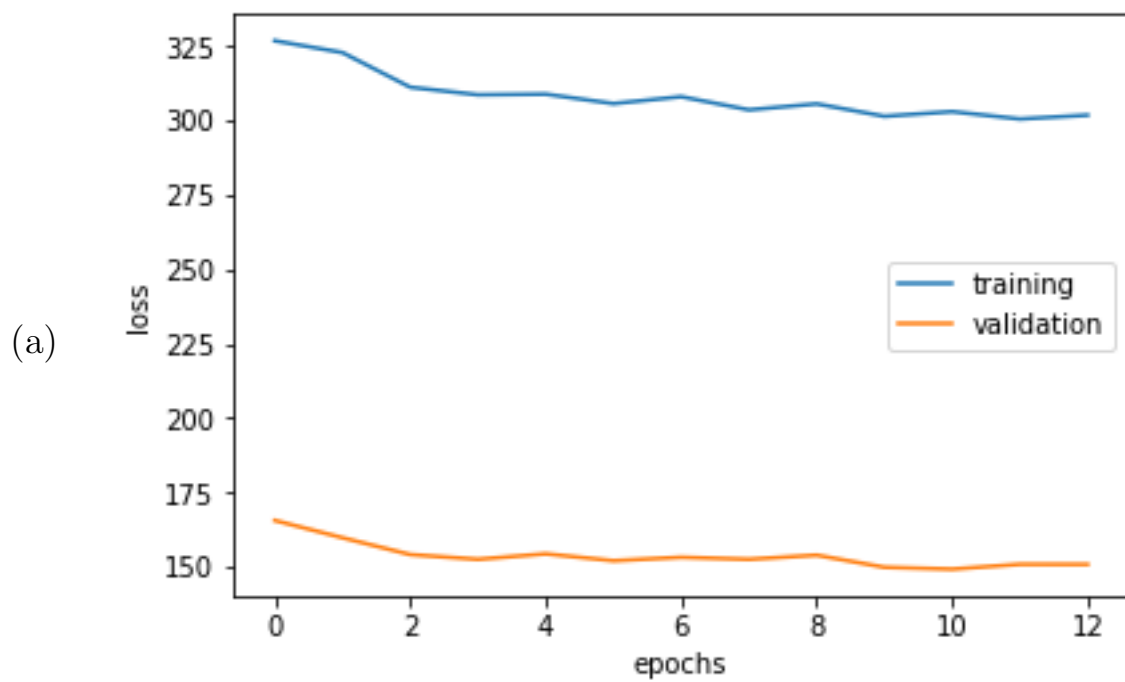
2. Experimenting with values of η_0, η_1

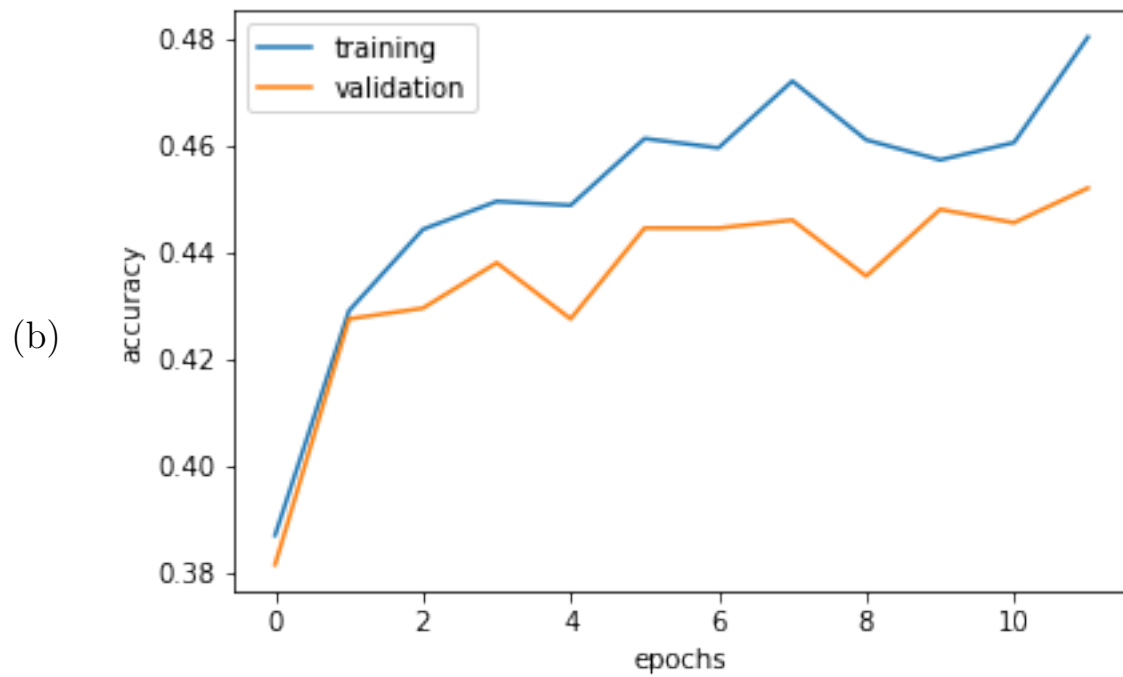
(a) $\eta_0 = 10$, $\eta_1 = 5$, $epochs = 12$, $L(\theta) = 303.0168711323694$

(b) The below graph shows the variation of loss vs epochs. Please note that due to high variance of loss in the few initial epochs, I have kept the minimum epochs required to check for the stopping condition is 10 epochs. Therefore, the algorithm doesn't exit at epoch = 2 as shown in the below graph.



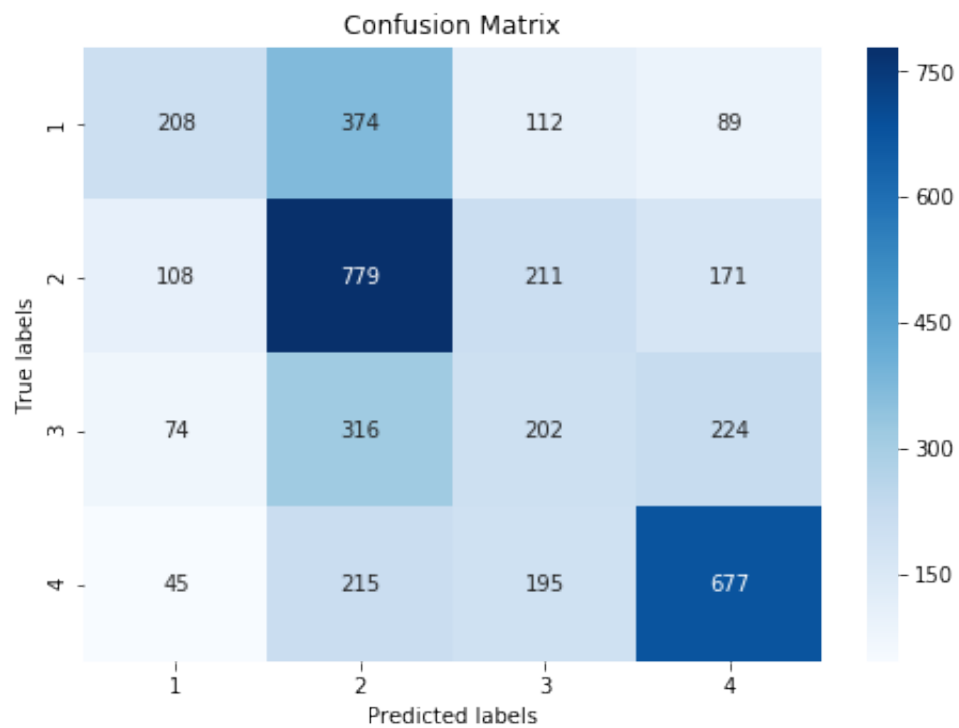
3. Evaluation on validation set



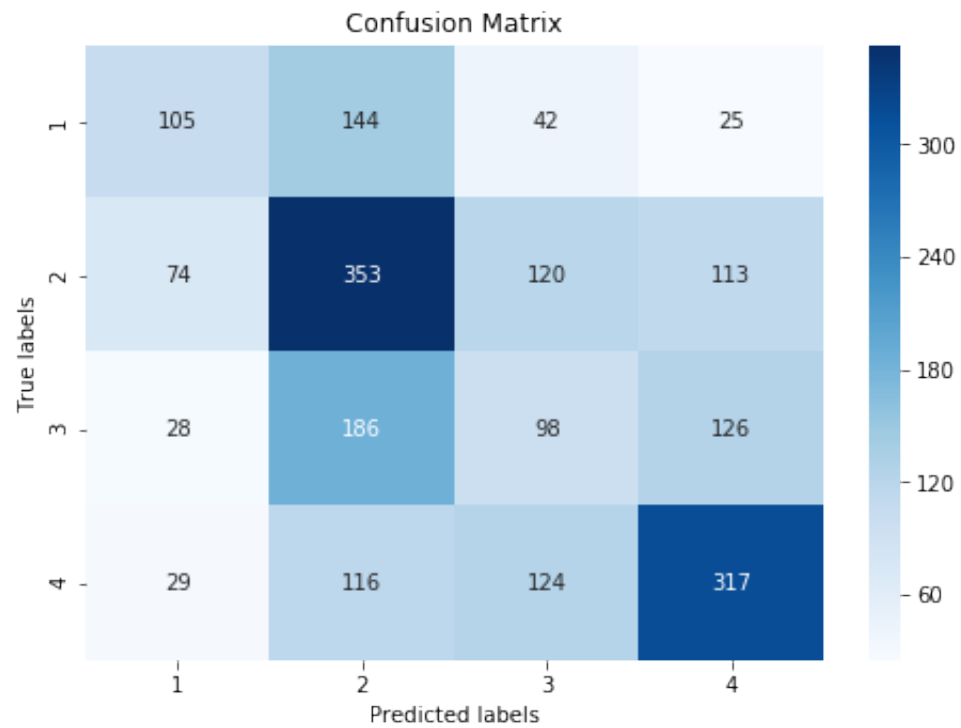


4. Confusion Matrices

(a) Training data



(b) Validation data



2.4 Kaggle Challenge

Accuracy from kaggle - 0.47333