# Cover page for answers.pdf
# CSE512 Fall 2018 - Machine Learning - Homework 3

Your Name: SAIF SULEMAN VAZIR

Solar ID: 112072061

NetID email address: saifsuleman.vazir@stonybrook.edu

Names of people whom you discussed the homework with: -

# HOMEWORK 3

1) 1.1   $X_1$ = boolean variable

$y$ = boolean variable

$X_2$ = continuous variable

As $X_1$ is a boolean variable, we can formulate it in terms of a bernoulli distribution (for 1 sample) or binomial distribution (>1 samples)

We require only one parameter for $X_1$ for 'each' class.

For $Y=1$, we have $\delta_{10}^{x_i}(1-\delta_{10})^{1-x_i}$   [$x_i$ is sample data point in $X_1$]

For $Y=0$, we have $\delta_{11}^{x_i}(1-\delta_{11})^{1-x_i}$   [$x_i$ is 1 or terms when $x_i = 1$ else 0 or false]

We need 1 parameter for $Y$, as we have only 2 classes which can be represented as $\gamma$ and $1-\gamma$.

For continuous variable $X_2$, we require a Gaussian distribution with $N \propto (\mu_{ik}, \sigma_{ik}^2)$.

As we have 2 classes, $k=2$, therefore we require 4 parameters in total.

Total no. of parameters required = $2+1+4 = \underline{7}$

Now, to find $P(Y|X)$ we use Bayes theorem:

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

Lets say we want to find $P(Y=1|X)$ [class = 1]

we have:
$$P(Y=1|X) = \frac{P(X|Y=1) \cdot P(Y=1)}{\sum_j P(Y=j) \cdot \prod_j P(X_j|Y=j)}$$

(using conditional independence for Naive Bayes)
$$P(Y=1|X) = \frac{P(Y=1) \cdot P(X_1|Y=1) \cdot P(X_2|Y=2)}{P(Y=1) \cdot P(X_1|Y=1) P(X_2|Y=1) + P(Y=0) P(X_1|Y=0) P(X_2|Y=0)}$$

$$= \frac{\gamma \cdot [\delta_1^{x_i}(1-\delta_1)^{x_i}] \left[\frac{1}{\sqrt{2\pi} \sigma_1} e^{-\frac{(z_i - \mu_{1k})^2}{2\sigma_1^2}}\right]}{\gamma[\delta_1^{x_i}(1-\delta_1)^{1-x_i}] \left(\frac{1}{\sqrt{2\pi}\sigma_1^2} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_1^2}}\right) + (1-\gamma)(\delta_0^{x_i}(1-\delta)^{x_i}) \cdot \left(\frac{1}{\sqrt{2\pi}\sigma_0^2} e^{-\frac{(x_i-\mu_0)^2}{2\sigma_0^2}}\right)}$$

Suppose we have total `N` data points
~~for each~~ $\therefore |X_1| = |X_2| = N$
we can write the above $f^n$ as:
(using only numerator as denominator
becomes independent of $Y=y$ as will
not affect ~~our minimization~~ our formula ]

$$P(Y=1|X) \propto \gamma \cdot [\delta_1^m (1-\delta_1)^{N-m}] \cdot \left[\prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x_i-\mu_1)^2}{2\sigma_1^2}}\right]$$

If we want to maximize our function
for any $Y=y_k$:

$$P(Y=y_k|X) \propto \underset{y_k}{\text{argmax}} \; P(Y=y_k) \prod P(X_1|Y=y_k) \cdot P(X_2|Y=y_k)$$

Taking log on both sides we get:
$$\log P(Y=y_k|X) = \underset{y_k}{\text{argmax}} \log P(Y=y_k) + \log P(X_1|Y=y_k) + \log P(X_2|Y=y_k)$$

$\log$ of prior $\qquad$ MLE for
$\qquad$ MLE for $\qquad$ binomial
$\qquad$ Gaussian distribution $\qquad$ distribution

1.2) Consider X to be a set of boolean variables, we can say that each $X_i$ has a binomial distribution.
We use the same notations from previous answer:

$P(Y=1) = \gamma \quad P(Y=0) = 1-\gamma$

Now, as $X = <X_1, X_2 \ldots X_n>$ all $X_i$ are boolean R.V's themselves, we have:

$P(X_i | Y=1) = \delta_{i1}^{X_i} (1-\delta_{i1})^{1-X_i}$

$P(X_i | Y=0) = \delta_{i0}^{X_i} (1-\delta_{i0})^{1-X_i}$

Now we use the eq$^n$:

$$P(Y=1|X) = \frac{P(Y=1)\cdot P(X|Y=1)}{P(Y=1)P(X=1|Y=1) + P(Y=0)P(X|Y=0)}$$

[dividing by $P(Y=1)P(X|Y=1)$]

$\therefore P(Y=1|X) = \dfrac{1}{1 + \dfrac{P(Y=0)\,P(X|Y=0)}{P(Y=1)\,P(X|Y=1)}}$

$= \dfrac{1}{1 + e^{\left[\ln\left(\frac{P(Y=0)\,P(X|Y=0)}{P(Y=1)\,P(X|Y=1)}\right)\right]}}$

$= \dfrac{1}{1 + e^{\left[\ln\left(\frac{1-\gamma}{\gamma}\cdot\frac{\prod_{i=1}^{n} P(X_i|Y=0)}{\prod_{i=1}^{n} P(X_i|Y=1)}\right)\right]}}$

$= \dfrac{1}{1 + e^{\left[\ln\frac{1-\gamma}{\gamma} + \ln\prod_{i=1}^{n}\frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right]}}$

$= \dfrac{1}{1 + e^{\left[\ln\frac{1-\gamma}{\gamma} + \sum_{i=1}^{n}\ln\frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right]}}$

$= \dfrac{1}{1 + e^{\left[\ln\frac{1-\gamma}{\gamma} + \sum_{i=1}^{n}\ln\left[\frac{\delta_{i0}^{X_i}(1-\delta_{i0})^{1-X_i}}{\delta_{i1}^{X_i}(1-\delta_{i1})^{1-X_i}}\right]\right]}}$

$$= \frac{1}{1+e^{\left[\ln\frac{1-\gamma}{\gamma} + \sum_{j=1}^{n} \ln\left(\frac{\delta_{io}}{\delta_{i1}}\right)x_i + \ln\left(\frac{1-\delta_{io}}{1-\delta_{i1}}\right)^{(1-x_i)}\right]}}$$

$$= \frac{1}{1+e^{\left[\ln\frac{1-\gamma}{\gamma} + \sum_{j=1}^{n} x_i \ln\frac{\delta_{io}}{\delta_{i1}} + (1-x_i)\ln\frac{1-\delta_{io}}{1-\delta_{i1}}\right]}}$$

[split last term]

$$= \frac{1}{1+e^{\left[\ln\frac{1-\gamma}{\gamma} + \sum_{i=1}^{n} x_i\ln\frac{\delta_{io}}{\delta_{i1}} + \ln\frac{1-\delta_{io}}{1-\delta_{i1}} - x_i\ln\frac{1-\delta_{io}}{1-\delta_{i1}}\right]}}$$

$$= \frac{1}{1+e^{\left[\ln\frac{1-\gamma}{\gamma} + \sum_{i=1}^{n}\ln\frac{1-\delta_{io}}{1-\delta_{i1}} + \sum_{i=1}^{n} x_i\left[\ln\frac{\delta_{io}}{\delta_{i1}} - \ln\frac{1-\delta_{io}}{1-\delta_{i1}}\right]\right]}}$$

$$= \frac{1}{1+e^{\left[\ln\frac{1-\gamma}{\gamma} + \sum_{i=1}^{n}\ln\frac{1-\delta_{io}}{1-\delta_{i1}} + \sum_{i=1}^{n} x_i\ln\frac{\delta_{io}\cdot(1-\delta_{i1})}{\delta_{i1}(1-\delta_{io})}\right]}}$$

Comparing with the original eqⁿ:
$$P(Y=i\mid X) = \frac{1}{1+e^{-\left(\sum_{i=1}^{d}\theta_i x_i + \theta_{d+1}\right)}}$$

we have:
$$P(Y=1\mid X) = \frac{1}{1+e^{-\left[\ln\frac{\gamma}{1-\gamma} + \sum_{j=1}^{n}\ln\frac{1-\delta_{i1}}{1-\delta_{io}} + \sum_{i=1}^{n} x_i\ln\frac{\delta_{i1}(1-\delta_{io})}{\delta_{io}(1-\delta_{i1})}\right]}}$$

[note that (-1) has been multiplied ~~2 times~~ as
we have reversed the numerator and
denominator in ln]

$$\theta_{d+1} = \ln\frac{\gamma}{1-\gamma} + \sum_{i=1}^{n}\ln\frac{1-\delta_{i1}}{1-\delta_{io}}$$

$$\theta_i = \ln\frac{\delta_{i1}(1-\delta_{io})}{\delta_{io}(1-\delta_{i1})}$$

we got: $P(Y=1\mid X) = \dfrac{1}{1+e^{-\left[\sum_{i=1}^{d}\theta_i x_i + \theta_{d+1}\right]}}$

2-1 To prove:

$$\frac{\partial}{\partial\Theta} \log P(y^i \mid \bar{x}^i; \Theta) = [y^i - P(y=1 \mid \bar{x}^i; \Theta)]\bar{x}^i$$

we can write:

$$\log P(y^i \mid \bar{x}^i; \Theta) = \underset{(\text{when } y=1)}{y^i \log P(y=1 \mid \bar{x}^i; \Theta)} + \underset{(\text{when } y=0)}{(1-y^i)\log P(y=0 \mid x^i; \Theta)}$$

⇒ we know that $P(y=1 \mid \bar{x}^i; \Theta) = \dfrac{1}{1+e^{-\Theta^T \bar{x}^i}} = \dfrac{e^{\Theta^T \bar{x}^i}}{1+e^{\Theta^T \bar{x}^i}}$

$$P(y=0 \mid x^i; \Theta) = 1 - P(y=1 \mid \bar{x}^i; \Theta)$$
$$= 1 - \frac{e^{\Theta^T \bar{x}^i}}{1+e^{\Theta^T \bar{x}^i}} = \frac{1}{1+e^{\Theta^T \bar{x}^i}}$$

we get:

$$y^i \log \frac{e^{\Theta^T \bar{x}^i}}{1+e^{\Theta^T \bar{x}^i}} + (1-y^i)\log \frac{1}{1+e^{\Theta^T \bar{x}^i}}$$

$$y^i \log e^{\Theta^T \bar{x}^i} + y^i \log \frac{1}{1+e^{\Theta^T \bar{x}^i}} + \log \frac{1}{1+e^{\Theta^T \bar{x}^i}} - y^i \log \frac{1}{1+e^{\Theta^T \bar{x}^i}}$$

⇒ $y^i \log e^{\Theta^T \bar{x}^i} - \log 1 + e^{\Theta^T \bar{x}^i}$

⇒ $y^i \Theta^T \bar{x}^i - \log 1 + e^{\Theta^T \bar{x}^i}$

Differentiating w.r.t $\Theta$ we get:

$$\frac{\partial P(y^i \mid \bar{x}^i; \Theta)}{\partial\Theta} = y^i \bar{x}^i + \frac{1}{1+e^{\Theta^T \bar{x}^i}} \cdot e^{\Theta^T \bar{x}^i} \cdot \bar{x}^i$$
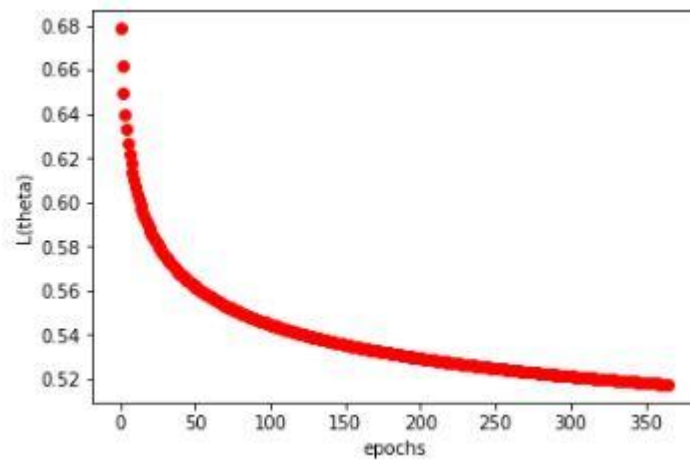
$$= \left[ y^i + \frac{e^{\Theta^T \bar{x}^i}}{1+e^{\Theta^T \bar{x}^i}} \right]\bar{x}^i$$

But $P(y=1 \mid \bar{x}^i; \Theta) = \dfrac{e^{\Theta^T \bar{x}^i}}{1+e^{\Theta^T \bar{x}^i}}$

∴ $\dfrac{\partial}{\partial\Theta} P(y^i \mid \bar{x}^i; \Theta) = [y^i - P(y=1 \mid \bar{x}^i; \Theta)]\bar{x}^i$

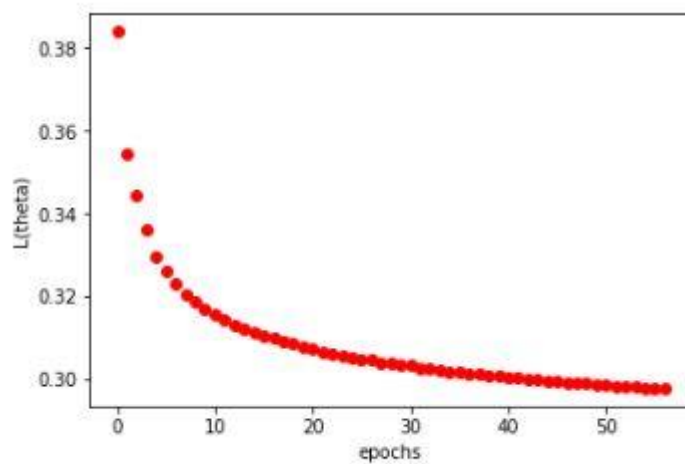2.3 1) a) Number of epochs till termination: 365

b) Plot:



c) Final value of $L(\Theta) = 0.5175504365325073$
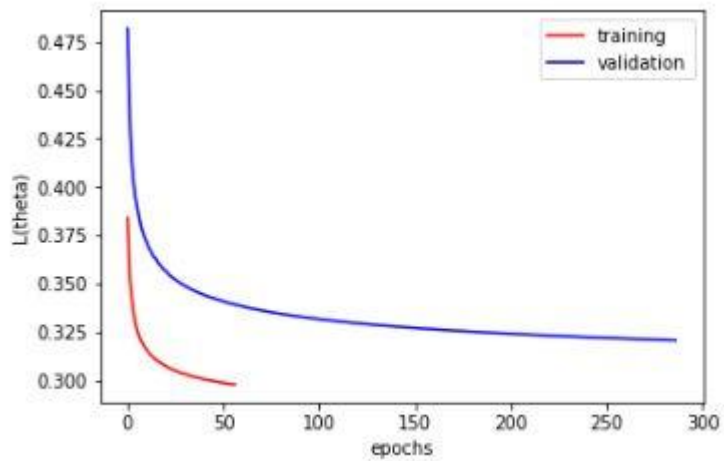
2) a) Best value of $\eta_0$, $\eta_1 = (2.5, 0.1)$

Number of epochs= 69

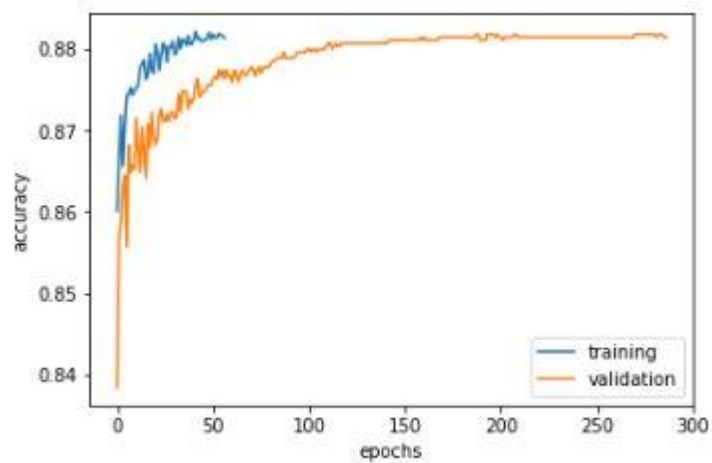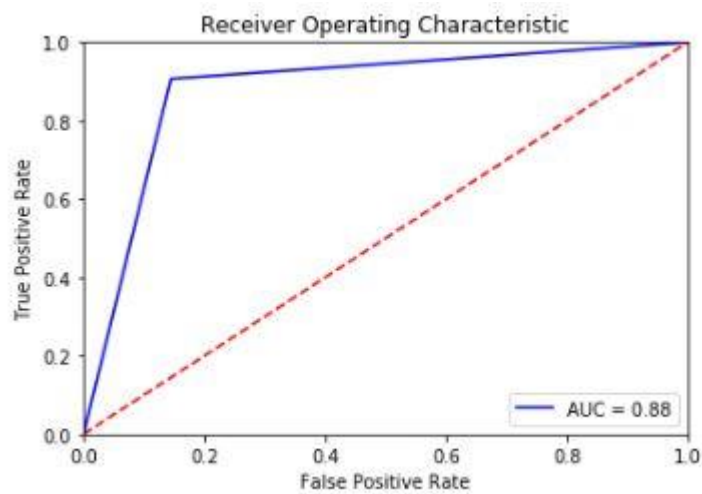Final value of $L(\Theta)= 0.29680946550912407$

b) Plot:

3) a) Plot for validation+training:

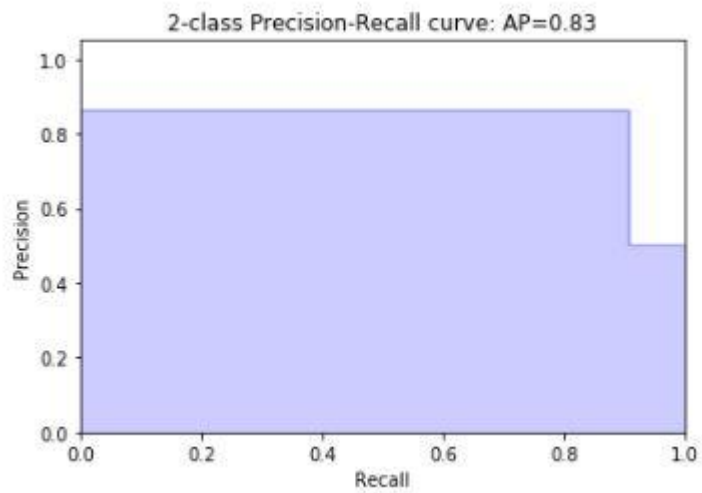

b) Plot for accuracy of model on training and validation:



4) a) Plot for ROC curve on validation data:



Area under the curve= 0.88

b) Plot for precision-recall curve on validation data:



2-class Precision-Recall curve: AP=0.83

Average precision= 0.83

2.4 1) Best accuracy for Kaggle submission = 87.333%