

CSE512 Spring 2021 - Machine Learning - Homework 2

Your Name: Sagar Onkar Toshniwal

Solar ID: 113260061

NetID email address: sagaronkar.toshniwal@stonybrook.edu

Names of people whom you discussed the homework with: None but I have taken reference from medium, towardsdatascience, Google, Stackoverflow, etc.

1) Question 1 – Parameter estimation (45 points)

1.1 Question 1.1 – MLE (15 points)

1. (5 points) Give the log-likelihood function of X given $\underline{\quad}$.
2. (5 points) Compute the MLE for $\underline{\quad}$ in the general case.
3. (5 point) Compute the MLE for $\underline{\quad}$ using the observed X .

ML HW 2

1.1] MLE (Log likelihood)

→ Poisson Distribution for the Random variable X such that $X = x_1, x_2, \dots, x_n$ for the expected outcome 1 is given by λ

$$P(X=k|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!} \quad \dots \textcircled{1} \quad k \in \{0, 1, 2, \dots, n\}$$

\therefore Log likelihood of X given λ .

$$\begin{aligned} P(x_1, \dots, x_n|\lambda) &= P(x_1|\lambda) \cdot P(x_2|\lambda) \dots \text{since they are conditionally independent} \\ &= \prod_{i=1}^n P(x_i|\lambda) \end{aligned}$$

Taking log on both sides

$$\begin{aligned} \log(P(x_1, \dots, x_n|\lambda)) &= \log \left(\prod_{i=1}^n P(x_i|\lambda) \right) \\ &= \prod_{i=1}^n \log \left(\frac{\lambda^{x_i} e^{-\lambda}}{x_i!} \right) \\ &= \prod_{i=1}^n \left[\log(\lambda^{x_i} \cdot e^{-\lambda}) - \log(x_i!) \right] \\ &= \log \lambda \sum_{i=1}^n x_i + \sum_{i=1}^n (-\lambda) \log e - \sum_{i=1}^n \log(x_i!) \\ &= \log \lambda \sum_{i=1}^n x_i + \sum_{i=1}^n (-\lambda) - \sum_{i=1}^n \log(x_i!) \end{aligned}$$

$\boxed{\log(P(x_1, \dots, x_n|\lambda)) = \log \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n \log(x_i!)}$
\textcircled{1}

1.1]

2) MLE for λ

→ MLE for given λ can be derived after taking its derivative and equating/setting it with 0.

$$\frac{d L(x)}{d \lambda} = 0$$

$$\therefore \hat{\lambda}_{MLE} = \arg \max_{\lambda} \ln(L(x/\lambda))$$

$$\therefore \frac{d}{d \lambda} \ln[L(x/\lambda)] = 0$$

$$\frac{d}{d \lambda} \left[\log \lambda \sum_{i=1}^n x_i - n \lambda - \sum_{i=1}^n \log(x_i) \right] = 0$$

$$\frac{1}{\lambda} \sum_{i=1}^n x_i - n - 0 = 0$$

$$\boxed{\lambda = \frac{1}{n} \sum_{i=1}^n x_i} \quad \therefore \textcircled{2}$$

3) MLE for λ for observed x .

→ From Eqⁿ (2)

$$\lambda = \frac{1}{7} [4 + 12 + 3 + 5 + 6 + 9 + 17]$$

$$= \frac{1}{7} [56]$$

$$= 8$$

1.2 Question 1.2 – MAP (15 points)

1.2] MAP

▷ Compute Posterior distribution over λ

→ Given that :-

$$P(\lambda|x, \alpha) = \frac{\beta^\lambda \lambda^{\alpha+1} e^{-\beta\lambda}}{\Gamma(\alpha)} \quad \dots \lambda > 0$$

To compute the Posterior Probability $P(\lambda|x_1, \dots, x_n)$, we need to multiply log likelihood with prior.

$$\therefore P(\lambda|x_1, \dots, x_n) = P(x_1, \dots, x_n|\lambda) \cdot P(\lambda)$$

we know from Eq^n ⑪

$$\text{Log}(P(x_1, \dots, x_n|\lambda)) = \log \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n \log(x_i!)$$

$$P(x_1, \dots, x_n|\lambda) = \frac{\lambda^{\sum_{i=1}^n x_i} e^{-n\lambda}}{\prod_{i=1}^n (x_i!)}$$

$$\therefore P(\lambda|x_1, \dots, x_n) = \frac{\lambda^{\sum_{i=1}^n x_i} e^{-n\lambda}}{\prod_{i=1}^n (x_i!)} \cdot \frac{\beta^\lambda}{\Gamma(\alpha)} \cdot \lambda^{\alpha+1} \cdot e^{-\beta\lambda}$$

$$\therefore P(\lambda|x_1, \dots, x_n) = \frac{e^{-\lambda(n+\beta)} \lambda^{\sum_{i=1}^n x_i + \alpha + 1} \cdot \beta^\alpha}{\prod_{i=1}^n (x_i!) \cdot \Gamma(\alpha)}$$

$$\boxed{P(\lambda|x_1, \dots, x_n) = \frac{e^{-\lambda(n+\beta)} \lambda^{\sum_{i=1}^n x_i + \alpha + 1} \cdot \beta^\alpha}{\prod_{i=1}^n (x_i!) \cdot \Gamma(\alpha)}}$$

1.2]

2) Derive an Analytical expression for MAP estimate of λ

$\rightarrow \lambda_{\text{map}} \Rightarrow$ Take log likelihood of the conditional probability
and then differentiate it to 0.

$$P(\lambda | x_1, \dots, x_n) = \frac{\beta^\alpha e^{-\lambda(n+\beta)} \lambda^{\sum_{i=1}^n x_i + \alpha - 1}}{\Gamma(\alpha) \cdot \prod_{i=1}^n (x_i!)}$$

Taking Log on both side

$$\begin{aligned} \text{Log}(P(\lambda | x_1, \dots, x_n)) &= \log \left(\frac{\beta^\alpha e^{-\lambda(n+\beta)} \lambda^{\sum_{i=1}^n x_i + \alpha - 1}}{\Gamma(\alpha) \cdot \prod_{i=1}^n (x_i!)} \right) \\ &= \arg \max_{\lambda} \left[\log \left(\frac{\beta^\alpha}{\Gamma(\alpha) \prod_{i=1}^n (x_i!)} \right) - \lambda(n+\beta) + \left(\sum_{i=1}^n x_i + \alpha - 1 \right) \log \lambda \right] \end{aligned}$$

Taking derivative & setting it to 0

$$\begin{aligned} \frac{d}{d\lambda} \text{Log}(P(\lambda | x_1, \dots, x_n)) &= 0 \\ 0 &= \frac{d}{d\lambda} \left[\log \left(\frac{\beta^\alpha}{\Gamma(\alpha) \prod_{i=1}^n (x_i!)} \right) - \lambda(n+\beta) + \left(\sum_{i=1}^n x_i + \alpha - 1 \right) \log \lambda \right] \\ 0 &= 0 - (n+\beta) + \frac{1}{\lambda} \sum_{i=1}^n x_i + \alpha - 1 \end{aligned}$$

$$\boxed{\lambda = \frac{\left(\sum_{i=1}^n x_i \right) + \alpha - 1}{n + \beta}}$$

1.3 Question 1.3 – Estimator Bias (15 points)

1.3] Estimator Bias

① $\hat{n} = e^{-2X}$, show that \hat{n} is max. likelihood estimate of n .

→ we know that $n = e^{-2\lambda}$

$$\ln(n) = -2\lambda \ln e$$

$$\boxed{\lambda = -\frac{1}{2} \ln(n)} \quad \dots \textcircled{D}$$

from Poisson Distribution :-

$$P(X=k/\lambda) = \frac{e^{-\lambda} \lambda^k}{(k!)}$$

$$\lambda = -\frac{1}{2} (\ln(n)) \quad \dots \text{from } \textcircled{D}$$

$$\therefore P(X=k/\lambda) = \frac{e^{\frac{1}{2} \ln(n)}}{(k!)^k} \left(-\frac{1}{2} \ln(n) \right)^k$$

Taking log likelihood on both sides.

$$\ln(P(X=k/\lambda)) = \ln \left(\left(-\frac{1}{2} \ln(n) \right)^k + \frac{1}{2} \ln(n) \right) - \ln(k!)$$

$$= k \ln \left(-\frac{1}{2} \ln(n) \right) + \frac{\ln(n)}{2} - \ln(k!)$$

Taking derivative & setting it to 0.

$$\frac{d (\ln(\rho(\lambda)))}{d\lambda} = 0$$

$$0 = \frac{K}{n \ln(\bar{n})} + \frac{1}{2n}$$

$$\frac{1}{2n} = -\frac{K}{n \ln(\bar{n})}$$

$$\boxed{\hat{n}_{MLE} = e^{-2K}}$$

which is nothing but

$$\boxed{\hat{n}_{MLE} \approx \bar{n}}$$

② Bias of \hat{n}

$$\rightarrow \text{bias } (\hat{n}) = \text{Expected Value} - \text{True Value}$$
$$= E[\hat{n}] - n$$

$$E[\hat{n}] = E[e^{-2x}] = \sum_{k=0}^{\infty} e^{-2k} P(x=k)$$

$$P(x=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$\text{Substitute } P(X=K) = \frac{\sum_{K=0}^{\infty} e^{-2\lambda} \lambda^K e^{-\lambda}}{\sum_{K=0}^{\infty} K!}$$

$$\begin{aligned} E[e^{-2\lambda}] &= e^{-\lambda} \sum_{K=0}^{\infty} \left[\frac{e^{-2\lambda} \lambda^K}{K!} \right] \\ &= e^{-\lambda} \left[\frac{1}{1} + \frac{e^{-2\lambda} \lambda}{1} + \frac{e^{-4\lambda} \lambda^2}{2} + \dots \right] \\ &= e^{-\lambda} [e^{e^{-2\lambda}}] \quad \dots \text{ Taylor series} \\ &= e^{-\lambda} \left[1 + \frac{1}{e^2} \right] \end{aligned}$$

$$\boxed{\text{Bias } [\hat{n}] = e^{-\lambda} \left[1 + \frac{1}{e^2} \right] - e^{-2\lambda}}$$

$$\text{as } E[e^{-2\lambda}] = e^{-\lambda} \left[1 + \frac{1}{e^2} \right]$$

$$\hat{n} = e^{-2\lambda} \quad \dots \text{ proven from last (1.3.1)}$$

$$\therefore \text{Bias } [\hat{n}] = e^{-\lambda} \left[1 + \frac{1}{e^2} \right] - e^{-2\lambda}$$

③ $(-1)^X$ is unbiased & prove that it is the only unbiased of n . & briefly explain why it is unbiased.

→ Given, $\hat{n} = (-1)^X$

$$E[\hat{n}] = E[(-1)^X] = \sum_{k=0}^{\infty} \frac{(-1)^k k! e^{-\lambda}}{k!}$$
$$= e^{-\lambda} \cdot e^{-\lambda} = e^{-2\lambda}$$

$$\text{Bias} = E[\hat{n}] - n$$
$$= e^{-2\lambda} - e^{-2\lambda} = 0$$

When Bias is 0, then it is unbiased.

Hence the estimator is unbiased.

To prove only ~~$(-1)^X$~~ is ~~$E[-1]^X$~~ is the only unbiased

Reason why unbiased Estimator is a bad estimator to use :-

* It is a proven fact that Zero bias approach has a poor generalisability to new situation (in test data) and also assumes that the model has a precise knowledge of a true state of world (which definitely isn't the case in Machine Learning)

This results in wrong estimation of the expected value.

2.1 Question 2.1 – Logistic Regression Derivation (10 points)

* Also when such situation occurs, ~~error~~ increases.
 Noise \Rightarrow unavoidable noise
 error \Rightarrow Noise + Variance + (bias)²
 error has a drastic increase when the bias is 0.
 Because the variance increase a lot for any new data set.

* for ex :- $(-1)^x$ $x \Rightarrow$ even $\Rightarrow 1$
 $x \Rightarrow$ odd $\Rightarrow -1$
 A zero bias model will assume either 1/-1 for all cases, making it a bad estimator.

2.1] Logistic Regression :-

\rightarrow we know,

$$P(Y=1/x; \alpha) = \frac{\exp(\alpha^T \bar{x})}{1 + \sum_{j=1}^{K-1} \exp(\alpha_j^T \bar{x})} \quad \dots (1)$$

Considering two cases only $K=2$.

$$\therefore P(Y=1/x; \alpha) = \frac{\exp(\alpha^T \bar{x})}{1 + \exp(\alpha^T \bar{x})} \quad \dots (2)$$

$$P(Y=0/x; \alpha) = \frac{1}{1 + \exp(\alpha^T \bar{x})} \quad \dots (3)$$

Taking Log likelihood for eqn ③

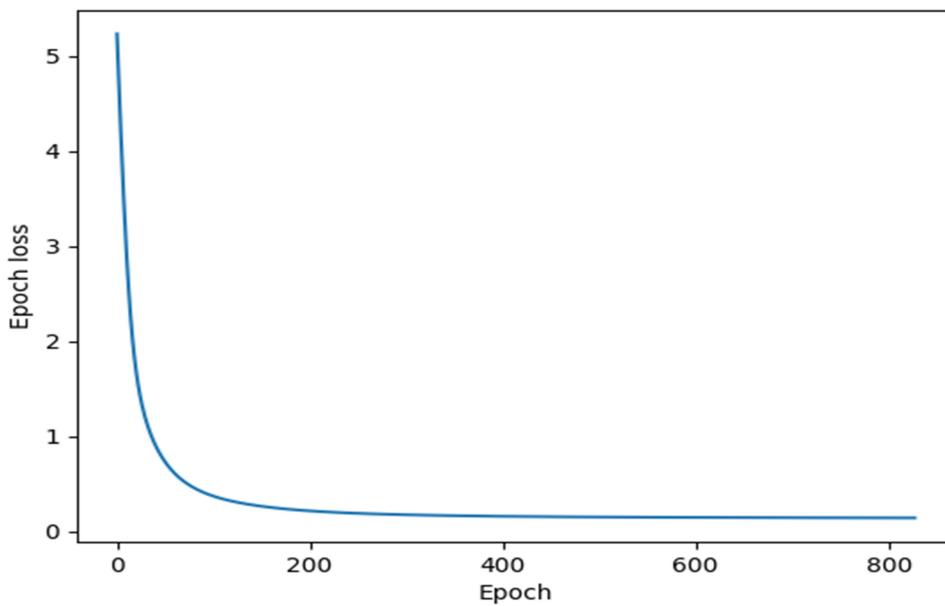
$$\begin{aligned} \underbrace{\log(P(y^i/x_i; \phi))}_{(1-y^i)^{1-y^i}} &= y^i \log(P(y=1/x_i; \phi)) \\ &\quad + (1-y^i) \log(P(y=0/x_i; \phi)) \\ &= y^i \log\left(\frac{\exp(\phi^T \bar{x})}{1+\exp(\phi^T \bar{x})}\right) + (1-y^i) \left(\log\left(\frac{1}{1+\exp(\phi^T \bar{x})}\right)\right) \\ &= y^i \phi^T \bar{x} - \log(1+\exp(1+\phi^T \bar{x})) \end{aligned}$$

Taking derivative & setting it on both sides.

$$\begin{aligned} \frac{d \log(P(y \neq x_i; \phi_c))}{d \phi_c} &= \frac{d}{d \phi_c} \left((y^i \phi_c^T - \log(1+\exp(1+\phi_c^T \bar{x}))) \bar{x}^i \right) \\ &= \frac{d}{d \phi_c} \left((y^i \phi_c^T - \log(1+\exp(1+\phi_c^T \bar{x}))) \bar{x}^i \right) \\ &= \left(y^i - \frac{\exp(\phi_c^T \bar{x}^i)}{1+\exp(\phi_c^T \bar{x}^i)} \right) \bar{x}^i \\ &= (y^i - P(y=1/\bar{x}, \phi_c)) \bar{x}^i \\ &= (\delta(c=y^i) - P(c=1/\bar{x}; \phi_c)) \bar{x}^i \end{aligned}$$

where $\delta(c=y^i)$ is indicator funn. & equating ground truth.

2.3.1) a) LOSS FUNCTION



Initially the loss is high as the model has considered random weights and its getting updated in the initial 200 epoch. Post 200 epoch the model has learned optimal weights, thus it doesn't update and that's why we see a saturation line

2.3.1) b) Accuracy , Confusion Matrix and Accuracy based on matrix

```
[82]: from sklearn.metrics import confusion_matrix
conf_mat = confusion_matrix(y_train, y_p)
conf_mat1 = confusion_matrix(y_test, y_ptest)
print("Confusion Matrix :")
print (confusion_matrix(y_test, y_ptest))

Confusion Matrix :
[[611  24]
 [ 43  13]]

[79]: from sklearn.metrics import accuracy_score,average_precision_score,precision_recall_curve
print ("Accuracy on training Data :",accuracy_score(y_train, y_p))
print ("Accuracy on testing Data :",accuracy_score(y_test, y_ptest))

Accuracy on training Data : 0.9424746743849494
Accuracy on testing Data : 0.9030390738060782

[80]: acc_conf = np.trace(conf_mat)/np.sum(conf_mat)
acc_conf1 = np.trace(conf_mat1)/np.sum(conf_mat1)
print ()

print ("Accuracy on training Confusion Matrix : ",acc_conf)
print ("Accuracy on testing Confusion Matrix : ",acc_conf1)

Accuracy on training Confusion Matrix : 0.9424746743849494
Accuracy on testing Confusion Matrix : 0.9030390738060782
```

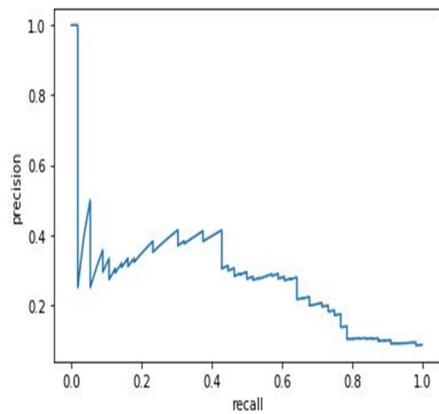
2.3.1) c) Average Precision and Recall

```
[111]: precision, recall, threshold = precision_recall_curve(y_test, Prob[:,1])
print ("Average Precision Score : ",average_precision_score(y_test, Prob[:,1]))
```

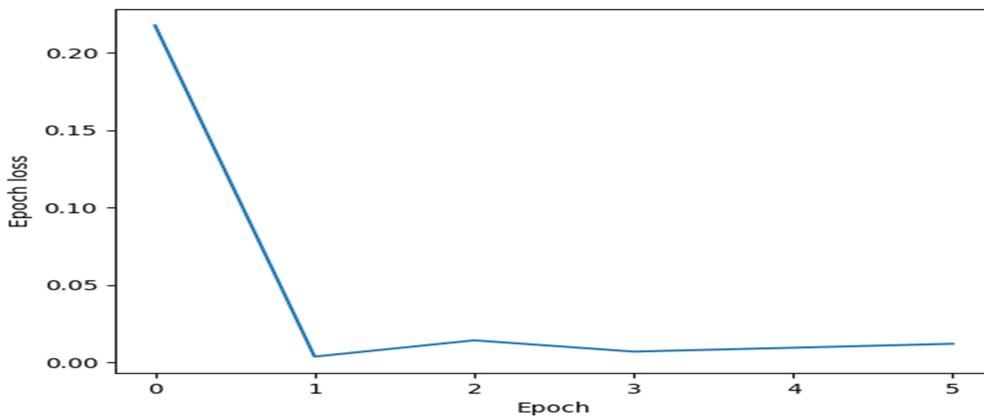
```
Average Precision Score :  0.2825801776059559
```

```
[110]: plt.plot (recall,precision)
plt.xlabel("recall")
plt.ylabel("precision")
plt.show
```

```
[110]: <function matplotlib.pyplot.show(*args, **kw)>
```



2.3.2) a) Loss Function without Standardization



For non-standardized data, The presence of feature value X in the formula will affect the step size of the gradient descent. The difference in ranges of features will cause different step sizes for each feature. Because of this reason the model learns only for few epoch and then comes out of the loop.

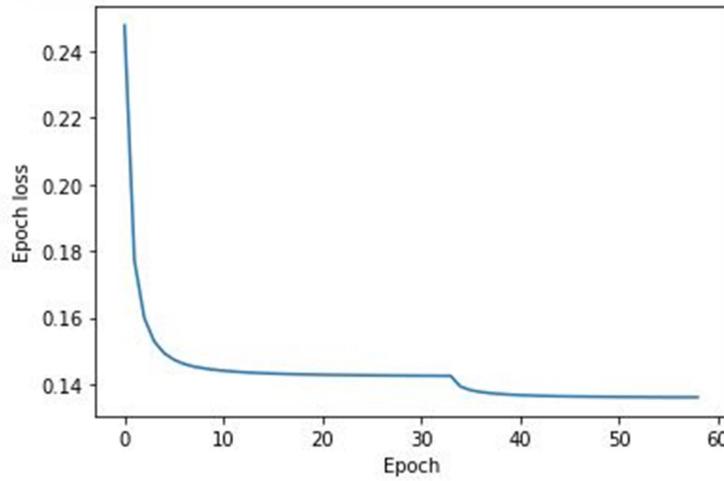
It does not converge, infact it diverges. The number of epoch are very few. In our case just below 10. Since the parameter are not learned well, accuracy is very low. Thus the final performance is very bad as compared to standardized data.

2.3.2) b) Tuning Hyperparameter

- When eta_start increases it converges faster but if it is increased to a lot then it may also diverges. So an optimum value would be between (0.01,0.1). When eta_start decreases it tends to be closer to eta_end and thus it doesn't learn efficiently and may or may not give optimal result.
- When eta_end decreases it converges slow but with optimal value. If eta_end decreases then it dosent provide optimal solution . it approximates early.
- For max_epoch I found 1000 was optimal. Post 1000 it showed the same result with more computational time. When below 1000 accuracy was not as good as it was with 1000 (when other parameters were given in HW2)
- Decreasing batch size take more training time but it converges faster (Also no. of epoch decreases). Increasing batch size will take less training time, more epoch and converge slow.

2.3.2) c) different parameter

- Based on above description(2.3.2.b) and multiple combination, I decreased batch size to 16, increased eta_start to 0.1, decreased epoch to 250. I found that the accuracy increased, precision increased, and computationally it worked faster and converge faster with minimum loss.



```
[121]: from sklearn.metrics import confusion_matrix
conf_mat = confusion_matrix(y_train, y_p)
conf_mat1 = confusion_matrix(y_test, y_ptest)
print("Confusion Matrix :")
print(confusion_matrix(y_test, y_ptest))

Confusion Matrix :
[[609  26]
 [ 36  20]]

[122]: from sklearn.metrics import accuracy_score,average_precision_score,precision_recall_curve
print ("Accuracy on training Data :",accuracy_score(y_train, y_p))
print ("Accuracy on testing Data :",accuracy_score(y_test, y_ptest))

Accuracy on training Data : 0.9424746743849494
Accuracy on testing Data : 0.91027496382055

[123]: acc_conf = np.trace(conf_mat)/np.sum(conf_mat)
acc_conf1 = np.trace(conf_mat1)/np.sum(conf_mat1)
print ()

print ("Accuracy on training Confusion Matrix :",acc_conf)
print ("Accuracy on testing Confusion Matrix :",acc_conf1)

Accuracy on training Confusion Matrix : 0.9424746743849494
Accuracy on testing Confusion Matrix : 0.91027496382055
```