

## Homework Assignment 1: Section A

Ankit Parekh

MS (Non-Thesis) Computer Engineering

ECE Department

ankitparekh@vt.edu

### Problem 1: Question 1)

If  $x_1, x_2, \dots, x_N$  are  $N$  observations sampled from a Bernoulli distribution with parameter  $p$ , the probability mass function for each  $x_i$  would be:

$$f(x_i; p) = p^{x_i} (1-p)^{1-x_i}$$

for  $x_i = 0$  and  $1$  &  $0 < p < 1$ .

$$\text{Likelihood Function } L(p) = \prod_{i=1}^n f(x_i; p) = p^{x_1} (1-p)^{1-x_1} \dots p^{x_n} (1-p)^{1-x_n}$$

$$L(p) = p^{\sum x_i} (1-p)^{n - \sum x_i}$$

To find maximum likelihood, one can maximize  $\ln L(p)$  instead of  $L(p)$ , since  $y = \ln(x)$  is a naturally increasing function.

So, we can find  $p$  at which  $\ln L(p)$  is maximum.  
or  $\log L(p)$

$$\log L(p) = (\sum x_i) \log(p) + (n - \sum x_i) \log(1-p)$$

$$\text{Taking derivative and setting it to zero: } \frac{\partial \log(L(p))}{\partial p} = \frac{\sum x_i}{p} - \frac{(n - \sum x_i)}{1-p} = 0 \quad \dots (1)$$

Solving (1): multiply  $p(1-p)$  on both sides:

$$(\sum x_i)(1-p) - (n - \sum x_i)p = 0$$

$$\text{To indicate } p \text{ as an estimate } \hat{p} = \frac{\sum_{i=1}^n x_i}{n} \quad \therefore \text{MLE of } p = \frac{\sum_{i=1}^n x_i}{n}$$

To show it is the maxima:

$$\frac{\partial^2 \log(L(p))}{\partial p^2} = - \frac{\sum_{i=1}^n x_i}{p^2} - \frac{\sum_{i=1}^n (1-x_i)}{(1-p)^2}$$

This term is negative  
indicating  $\hat{p}$  is a maxima.



**Problem 1: Question 2)**

①  $x_1, x_2, \dots, x_N$  are  $N$  observations sampled from a poisson distribution with probability mass function  $\Rightarrow$

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

Likelihood function is  $L(\lambda; x_1, \dots, x_N) = \prod_{i=1}^N \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$

Since  $y = \ln(x)$  is a naturally increasing function, we can try to find maxima of log-likelihood to get the maxima of likelihood.

Log-likelihood function is  $l(\lambda; x_1, \dots, x_N) = -N\lambda - \sum_{i=1}^N \ln(x_i!) + \ln(\lambda) \sum_{i=1}^N x_i$

We need to maximize log-likelihood and find the  $\lambda$  for the derivative is 0.

$$\frac{d}{d\lambda} l(\lambda; x_1, \dots, x_N) = -N + \frac{1}{\lambda} \sum_{i=1}^N x_i = 0$$

$\therefore \lambda = \frac{1}{N} \sum_{i=1}^N x_i$  MLE is  $\hat{\lambda} = \frac{1}{N} \sum_{i=1}^N x_i$  which is the mean of the observations.

2) By definition of expectation:  $E(Y) = \sum_{y \in Y} y P(Y=y)$

For Poisson Distribution:  $E(Y) = \sum_{y \geq 0} y \cdot \frac{\lambda^y e^{-\lambda}}{y!} \dots \textcircled{I}$

solving  $\textcircled{I}$ :  $E(Y) = \lambda e^{-\lambda} \sum_{y \geq 1} \frac{1}{(y-1)!} \lambda^{y-1}$  as  $y=0$  term vanishes

$$= \lambda e^{-\lambda} \sum_{z \geq 0} \frac{\lambda^z}{z!} \text{ putting } z = y-1$$

$$= \lambda e^{-\lambda} e^{\lambda} \quad \text{Taylor Series Expansion for } e^{\lambda}$$

$\boxed{E(Y) = \lambda}$   $\therefore$  Expectation of Poisson Distribution

## Problem 2: Question 1) Used Python for calculations

```
import numpy as np
```

```
1. vehicle_speed_dataset = {'number_of_wheels':[4, 4, 2, 8, 4, 3], 'cost':[15000,
25000, 5000, 40000, 22000, 17000]}
2. def standardize_feature(arr):
3.     feature = np.array(arr)
4.     mean = np.mean(feature)
5.     print("Mean: ", mean)
6.     std = np.std(feature)
7.     print("Standard Deviation: ", std)
8.     arr = [(x - mean)/std for x in arr]
9.     print("Standardized Values:", arr)
10. print("Number of Wheels : ", vehicle_speed_dataset['number_of_wheels'])
11. standardize_feature(vehicle_speed_dataset['number_of_wheels'])
12. print("Cost : ", vehicle_speed_dataset['cost'])
13. standardize_feature(vehicle_speed_dataset['cost'])
```

### Answer (Output):

Number of Wheels: [4, 4, 2, 8, 4, 3]

Mean: 4.167

Standard Deviation: 1.863

Standardized Values: [-0.0894, -0.0894, -1.1628, 2.057, -0.0894, -0.626]

Cost: [15000, 25000, 5000, 40000, 22000, 17000]

Mean: 20666.667

Standard Deviation: 10687.479

Standardized Values: [-0.530, 0.405, -1.466, 1.809, 0.125, -0.343]

## Problem 2: Question 2)

Q.2) Linear Model  $h_w(x) = w_0 + w_1x_1 + w_2x_2$

Mean squared error is the sum of squared differences between the predicted and the true values.

Cost function  $J(w)$  for least squares is given by:

Features:

$x_1$ : no. of wheels

$x_2$ : cost

$$J(w) = \frac{1}{n} \sum_{i=1}^n (\text{predicted}_i - \text{true}_i)^2$$

where  $n$  = number of samples

$$J(w) = \frac{1}{n} \sum_{i=1}^n (h_w(x_i^{(i)}) - y_i^{(i)})^2$$

$$J(w) = \frac{1}{n} \sum_{i=1}^n ((w_0 + w_1x_1^{(i)} + w_2x_2^{(i)}) - (y_i^{(i)}))^2$$

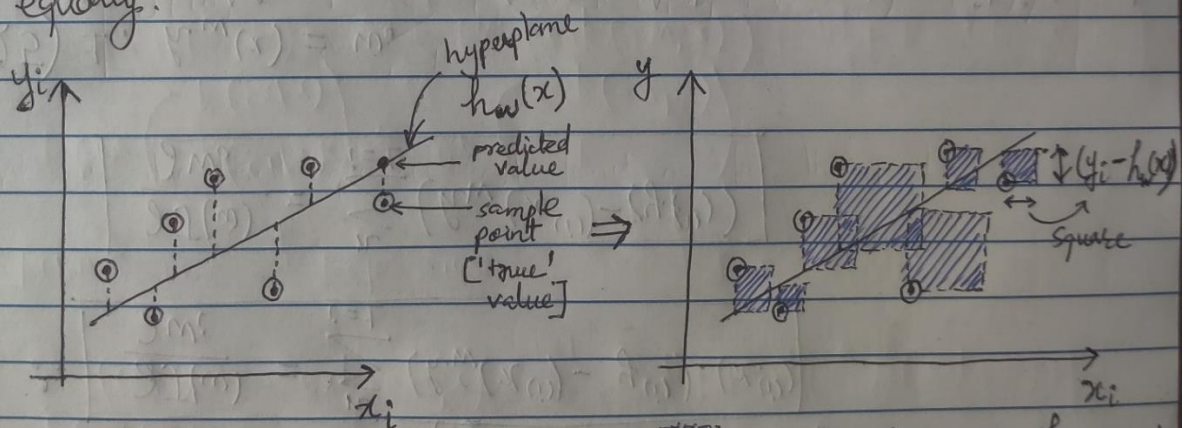


## Problem 2: Question 3)

Q.3) \* We have used Least Squares Cost function for this problem to find the parameters ( $w$ ) of a hyperparameter plane which best fits the linear regression dataset.

\* Computing the total squared error between the associated hyperparameter plane and the true values gives an appropriate measure of how well does the line fits the dataset.

\* It also takes into account that the difference between predicted values and true values might have different signs for different points, and hence squaring the individual terms so that both positive and negative values are treated equally.



\* We get the area marked in blue using Least squares framework and we want to minimize this area. We try to minimize this area using parameter ( $w$ ) tuning.

$$\underset{w}{\text{minimize}} J(w) = \underset{w}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n ((w_0 + w_1 x_1^{(i)} + w_2 x_2^{(i)}) - y^{(i)})^2$$

\* The residuals being normally distributed, L2 regularization performs better for linear regression.

\* We want our approximation to penalize outliers more and capture the general trend. The Least Squares is also differentiable at all places and hence can be used as a cost function in gradient descent.



**Problem 2: Question 4)**

Q4) Partial differentiating the cost function  $J(w)$ :  $x^{(j)}$ th sample

$$\begin{aligned}\frac{\partial J(w)}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{n} \sum_{j=1}^n (h_w(x^{(j)}) - y^{(j)})^2 \\ &= \frac{1}{n} \cdot 2 \sum_{j=1}^n (h_w(x^{(j)}) - y^{(j)}) \cdot \frac{\partial (h_w(x^{(j)}))}{\partial w_i} \\ &= \frac{1}{n} \cdot 2 \sum_{j=1}^n (h_w(x^{(j)}) - y^{(j)}) \cdot \frac{\partial (w_0 + w_1 x_1 + w_2 x_2)}{\partial w_i}\end{aligned}$$

$\therefore \left[ \frac{\partial J(w)}{\partial w_i} = \sum_{j=1}^n (h_w(x^{(j)}) - y^{(j)}) \cdot (x_i^{(j)}) \right]$

$\left[ \frac{\partial J(w)}{\partial w_0} = \sum_{j=1}^n (h_w(x^{(j)}) - y^{(j)}) \right]$

**Problem 2: Question 5)**

- Q5) 1.  $h_w(x) = w_0$
- The model is only using a single parameter (bias weight). It does not have any flexibility (variance) and is heavily biased. It will underfit on the dataset.
2.  $h_w(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + w_5 x^5$
- The model has a large number of features resulting in low bias and high variance. On introducing regularization, it will most likely overfit on the dataset.

**Problem 2: Question 6)**

- Problem 2: Q6)
1. We can increase the number of features in the model to reduce bias. If we decrease the regularization it will increase variance, and the model might fit better to the distribution. Model complexity will increase.
  2. We can decrease the model complexity by decreasing the number of features and increasing regularization. This will introduce more bias and reduce variance respectively which should fix the overfitting issue.

**Problem 3: Question 1)**

Q1)  $P(C_0)$ : Probability of patient having infectious (class  $C_0$ ) disease.

$P(C_0/x)$ : Probability of patient having infectious (class  $C_0$ ) disease given he/she exhibits observed symptoms in data vector  $x$ .

$P(x/C_0)$ : Probability of patient exhibiting observed symptoms (in data vector  $x$ ) given that he/she has the infectious disease (class  $C_0$ ).

According to Baye's Theorem:

$$P(C_0/x) = \frac{\overbrace{P(x/C_0)}^{\text{Likelihood}} \cdot \overbrace{P(C_0)}^{\text{Prior}}}{P(x)}$$

[Posterior]

**Problem 3: Question 2)**

$$Q2) \quad P(C_0/x) = \frac{P(x/C_0) P(C_0)}{P(x)} \dots (1) \quad ; \quad P(C_1/x) = \frac{P(x/C_1) P(C_1)}{P(x)} \dots (2)$$

$$P(x) = P(x/C_1) P(C_1) + P(x/C_0) P(C_0) \dots (3)$$

Adding (1) & (2) and substituting (3):  
We get:

$$P(C_1/x) + P(C_0/x) = 1$$

$$P(C_0/x) = 1 - P(C_1/x) \dots (I)$$

We are given:  $P(C_1/x) = \frac{1}{1 + \exp(-w^T x)}$

Substituting this in (I), we get:



$$P(C_0/x) = 1 - \frac{1}{1 + \exp(-w^T x)} = \frac{\exp(-w^T x)}{1 + \exp(-w^T x)}$$

$$\frac{P(C_1/x)}{P(C_0/x)} = \left( \frac{1}{1 + \exp(-w^T x)} \right) / \left( \frac{\exp(-w^T x)}{1 + \exp(-w^T x)} \right)$$

$$\frac{P(C_1/x)}{P(C_0/x)} = \frac{1}{\exp(-w^T x)}$$

Taking log on both sides:

$$\begin{aligned} \ln \left( \frac{P(C_1/x)}{P(C_0/x)} \right) &= \ln \left( \frac{1}{\exp(-w^T x)} \right) \\ &= -\ln(\exp(-w^T x)) \\ &= -(-w^T x) \end{aligned}$$

$$\boxed{\ln \left( \frac{P(C_1/x)}{P(C_0/x)} \right) = w^T x}$$

**Problem 3: Question 3)**

Q3)1) Sigmoid function is mapping range of  $y_i$  values into a  $[0,1]$  interval.

Class Mapping:  $\sigma(y_i) \geq 0.5 \Rightarrow \text{Class 1}$   
 $\sigma(y_i) < 0.5 \Rightarrow \text{Class 2}$  }  $c_i \in \{0,1\}$

$y_i = w^T x_i + w_0$  } Logistic Regression for Binary Classification

P.M.F. for  $\sigma(y_i) \Rightarrow f(x)^{c_i} (1-f(x))^{1-c_i}$

$$\text{Likelihood Function} = \prod_{i=1}^N f(x)^{c_i} (1-f(x))^{1-c_i}$$

$$\text{Likelihood Function} = L(w, w_0) = \prod_{i=1}^N \sigma(y_i)^{c_i} (1 - \sigma(y_i))^{1-c_i}$$

2) Taking log of likelihood function:

$$\log(L(w, w_0)) = \log \prod_{i=1}^N \sigma(y_i)^{c_i} (1 - \sigma(y_i))^{1-c_i}$$

$$= \sum_{i=1}^N c_i \log(\sigma(y_i)) + \sum_{i=1}^N (1-c_i) \log(1 - \sigma(y_i))$$

$$= \sum_{i=1}^N c_i \log \left[ \frac{\sigma(y_i)}{1 - \sigma(y_i)} \right] + \sum_{i=1}^N \log(1 - \sigma(y_i)) \dots \textcircled{I}$$

$$\frac{\sigma(y_i)}{1 - \sigma(y_i)} = \frac{1/1 + \exp(-y_i)}{1 - 1/1 + \exp(y_i)} = \frac{1}{\exp(-y_i)} \Rightarrow \log \left[ \frac{\sigma(y_i)}{1 - \sigma(y_i)} \right] = -\log(\exp(-y_i)) = y_i \dots \textcircled{II}$$

$$\text{Using } \textcircled{II} \text{ in } \textcircled{I}: \log(L(w, w_0)) = \sum_{i=1}^N c_i y_i + \sum_{i=1}^N \log(1 - \sigma(y_i))$$



$$\begin{aligned}
 \log(L(w, w_0)) &= \sum_{i=1}^N c_i y_i - \sum_{i=1}^N \log \left( \frac{1}{1 - \sigma(y_i)} \right) \\
 &= \sum_{i=1}^N c_i y_i - \sum_{i=1}^N \log \left( \frac{1}{1 - \frac{1}{1 + \exp(-y_i)}} \right) \\
 &= \sum_{i=1}^N c_i y_i - \sum_{i=1}^N \log \left( \frac{1 + \exp(-y_i)}{\exp(-y_i)} \right) \\
 &= \sum_{i=1}^N c_i y_i - \sum_{i=1}^N \log [(1 + \exp(-y_i)) (\exp(y_i))] \\
 &= \sum_{i=1}^N c_i y_i - \sum_{i=1}^N \log [1 + \exp(y_i)]
 \end{aligned}$$

Putting  $y_i = w^T x_i + w_0$

$$\log(L(w, w_0)) = \sum_{i=1}^N c_i (w^T x_i + w_0) - \log [1 + \exp(w^T x_i + w_0)]$$

$$3) \text{ i) } \frac{\partial \log(L(w, w_0))}{\partial w} = \sum_{i=1}^N c_i x_i - \frac{1}{1 + \exp(w^T x_i + w_0)} \exp(w^T x_i + w_0) (x_i)$$

$$\Rightarrow \frac{\partial \log(L(w, w_0))}{\partial w} = \sum_{i=1}^N c_i x_i - \frac{x_i \exp(w^T x_i + w_0)}{1 + \exp(w^T x_i + w_0)}$$

$$\Rightarrow \left[ \frac{\partial \log(L(w, w_0))}{\partial w} = \sum_{i=1}^N x_i \left[ c_i - \frac{\exp(w^T x_i + w_0)}{1 + \exp(w^T x_i + w_0)} \right] \right] - \text{①}$$

$$\text{ii) } \frac{\partial \log(L(w, w_0))}{\partial w_0} = \sum_{i=1}^N c_i - \frac{1}{1 + \exp(w^T x_i + w_0)} \cdot \exp(w^T x_i + w_0)$$

$$\Rightarrow \left[ \frac{\partial \log(L(w, w_0))}{\partial w_0} = \sum_{i=1}^N c_i - \frac{\exp(w^T x_i + w_0)}{1 + \exp(w^T x_i + w_0)} \right]$$