

**Machine learning – Assignment#2**  
**Aditya Gautam([agautam1@andrew.cmu.edu](mailto:agautam1@andrew.cmu.edu))**

Assignment #2 (Machine learning 601-B)

- ADITYA GAUTAM  
[agautam1@andrew.cmu.edu](mailto:agautam1@andrew.cmu.edu).

Problem 1) Independent events and Bayes theorem

- a) let  $A \cap B$  be two events.  
 probability that event  $B$  occurs given  $A$  has already occurred =  $P(B|A)$   
 $P(A \cap B) = P(A)P(B|A)$  [probability of occurrence of  $A$  given  $B$ ]

Similarly,  
 given event  $B$  has occurred, probability of occurrence of event  $A$   
 would be  $P(A)P(A|B) = P(A \cap B)$

From, above two points, we can say that

$$P(A \cap B) = P(A)P(B|A) = P(B)P(A|B)$$

$$\text{So, } P(A|B)P(B) = P(A)P(B|A)$$

$$\boxed{P(A|B) = \frac{P(A)P(B|A)}{P(B)}}$$

- b) Complete Sample Space  $\Sigma \rightarrow A_1, A_2, A_3, \dots, A_n$

Probability that event  $B$  occurs given  $A_i$  is the space

$$= P(A_i \text{ occurs})P(\text{B occurs given } A_i)$$

$$= P(A_i)P(B|A_i)$$

likewise for any  $A_i$   $P(B \cap A_i) = P(A_i)P(B|A_i)$

$$P(B) = P(B \cap \Sigma) = P(B \cap A_1) + P(B \cap A_2) + \dots + P(B \cap A_n)$$

complete space  $\Sigma$   $\left[ \text{as } \bigcup_i A_i = \Sigma \right]$

Applying ① in ②, we get.

$$P(B) = P(A_1)P(B|A_1) + P(A_2)P(B|A_2) + \dots + P(A_n)P(B|A_n)$$

$$\text{So, } \boxed{P(B) = \sum_{i=1}^n P(A_i)P(B|A_i)}$$

c).  $P(A_i) > 0 \text{ & } P(B) > 0$

Using the results from part a) & b), we can say the following about  $A_i$  and  $B$ ,

$$P(B)P(A_i^c|B) = P(B/A_i)P(A_i^c) \quad \text{--- (1)}$$

i.e probability that  $B$  has occurred given  $A_i^c$  is same as that of probability of  $B$  given  $A_i$  has occurred.

Also, from part b,

$$P(B) = \sum_i P(B/A_i)P(A_i) \text{ where } \cup A_i = \Omega \quad \text{--- (2)}$$

so, we can say that

$$P(A_i^c|B) = \frac{P(B/A_i)P(A_i^c)}{P(B)} \rightarrow \text{from (1)}$$

$$\Rightarrow P(A_i^c|B) = \frac{P(B/A_i)P(A_i)}{\sum_k P(B/A_k)P(A_k)} \rightarrow \text{applying (2)}$$

d)

a)  $P(A, B, C) = P(A|B, C)P(B|C)P(C)$

True by directly applying chain rule / ~~Bayes~~ theorem

b)  $P(A, B) = P(A|B)P(B|A)$

False.

Applying chain rule  $P(A|B) = P(A|B)P(B)$   
and  $P(B) = P(B|A)$  only when  $A$  &  $B$  are independent events

c)  $P(A, B, C) = P(B|A, C)P(C|A)$

True

Applying chain rule,

$$P(A, B, C) = P(B|A, C)P(C|A)$$

$$P(A, C) = P(C|A) = P(ANC) \rightarrow \text{proved in part (a)}$$

$$4) P(A_1B_1C_1) = P(B_1/A_1C_1) P(C_1/A_1) P(A_1) \rightarrow \boxed{\text{False}}$$

Applying chain rule of L.H.S

$$P(A_1B_1C_1) = P(B_1/A_1C_1) P(A_1C_1) = P(B_1/A_1C_1) P(A_1 \cap A_1) \rightarrow \text{proven in part (a)}$$

So, R.H.S has additional term  $P(C_1)$ .

$$5) P(A_1B_1) = P(A_1)P(B_1) \rightarrow \boxed{\text{False}}$$

$$\text{Chain Rule} \Rightarrow P(A_1B_1) = P(A_1)P(B_1/A_1) = P(B_1)P(A_1/B_1)$$

& Bayes Theorem  
This will be true only when  $A_1 \cap B_1$  are independent event.

$$e) X = \begin{cases} 1 & \text{if } A \text{ occurs} \\ 0 & \text{otherwise} \end{cases}$$

$P(A)$  → probability of  $A$  occurs.

$$E[X] = P(A)X_1 + (1-P(A))X_0$$

$$\Rightarrow E[X] = P(A)$$

$$\Rightarrow [E[X] + P(A)] = 0$$

when  $X=1$  if event  $A$  occurs,

$$\text{so } X = \begin{cases} 1 & \text{if } A \text{ occurs} \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = P(A)X_1 + (1-P(A))X_0 = P(A)$$

$\Rightarrow [E[X] - P(A)]$  so, old result will not hold.

### Problem 2: Maximum Likelihood Estimation.

$$(i) P(X_i=k) = (1-\theta)^k \theta$$

$$L(\hat{\theta}) = P(X_1, \dots, X_n | \hat{\theta}) \quad X_i \text{ are iid}$$

$$\text{as } X_i \text{ are i.i.d} \Rightarrow P(X_1, X_2, \dots, X_n | \hat{\theta}) = P(X_1 | \hat{\theta}) P(X_2 | \hat{\theta}) \dots P(X_n | \hat{\theta})$$

$$\Rightarrow L(\hat{\theta}) = P(X_1 | \hat{\theta}) P(X_2 | \hat{\theta}) P(X_3 | \hat{\theta}) \dots P(X_n | \hat{\theta})$$

$$L(\hat{\theta}) = \prod_{i=1}^{n\theta} P(X_i | \hat{\theta})$$

$$\begin{aligned}
 l(\theta) &= \log L(\hat{\theta}) \\
 &= \log \prod_{i=1}^n P(x_i | \theta) = \sum_{i=1}^n \log P(x_i | \theta) \\
 &= \log(P(x_1 | \theta)) + \log(P(x_2 | \theta)) + \log(P(x_3 | \theta)) + \dots + \log(P(x_n | \theta)) \\
 &= \log((1-\theta)^{x_1} \theta) + \log((1-\theta)^{x_2} \theta) + \dots + \log((1-\theta)^{x_n} \theta) \\
 &= [x_1 \log(1-\theta) + \log \theta] + [x_2 \log(1-\theta) + \log \theta] + \dots + [x_n \log(1-\theta) + \log \theta] \\
 &= x_1 \log(1-\theta) + x_2 \log(1-\theta) + \dots + x_n \log(1-\theta) \\
 &\quad + x_1 \log \theta + \log \theta + \dots + \log \theta \\
 &= (\log(1-\theta)) [x_1 + x_2 + \dots + x_n] + n \log \theta \\
 l(\theta) &= \log(1-\theta) \sum_{i=1}^N x_i + N \log \theta
 \end{aligned}$$

This doesn't depend on the order of random variable.

b)  $x = (1, 0, 3, 5, 18, 14, 15, 7, 13, 9, 0, 17, 4, 24, 3)$

code, plots and figure submitted to autobb.  
code  $\rightarrow$  logMLE.m (file name).

c) let's say that log maximum likelihood is maximum at  $\hat{\theta}, \text{ so}$

$$\begin{aligned}
 \frac{d l(\hat{\theta})}{\theta} &= 0 \\
 \Rightarrow \frac{d(\sum x_i \log(1-\theta) + N \log \hat{\theta})}{d\theta} &= 0 \\
 \Rightarrow \frac{\sum x_i}{(1-\theta)} (-1) + \frac{N}{\theta} &= 0 \\
 \Rightarrow \frac{N}{\theta} - \frac{\sum x_i}{1-\theta} &= 0 \Rightarrow N - N\theta = (\sum x_i)\theta \\
 \Rightarrow N &= \theta(N + \sum x_i) \\
 \Rightarrow \theta &= \frac{N}{N + \sum x_i}
 \end{aligned}$$

**Ques 2) Part b):**

**Matlab Code**

```
function [ ] = logMLE()

theta = 0.01:.01:0.5;

x = [1 0 3 5 18 14 5 7 13 9 0 17 4 24 3];

y5= sum(x(1:5))*log(1-theta) + 5*log(theta);
y10 = sum(x(1:10))*log(1-theta) + 10*log(theta);
y15 = sum(x(1:15))*log(1-theta) + 15*log(theta);

figure
plot(theta,y5)
xlabel('theta');
ylabel('y(log MLE)');
title('Log MLE (5 points)');

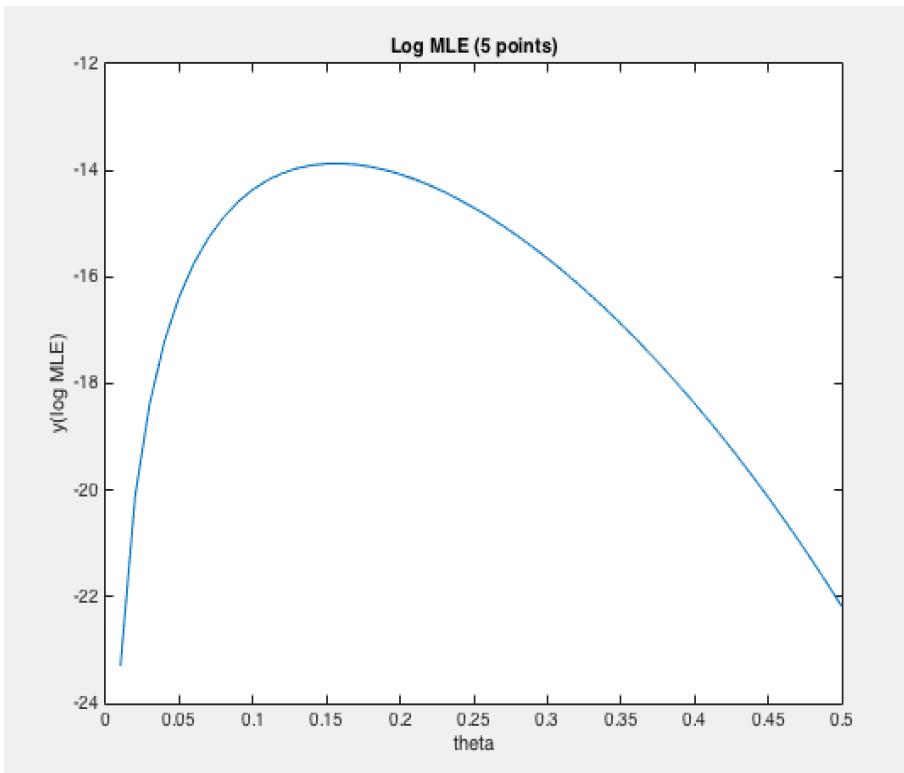
figure
plot(theta,y10)
xlabel('theta');
ylabel('y(log MLE)');
title('Log MLE (10 points)');

figure
plot(theta,y15)
xlabel('theta');
ylabel('y(log MLE)');
title('Log MLE (15 points)');

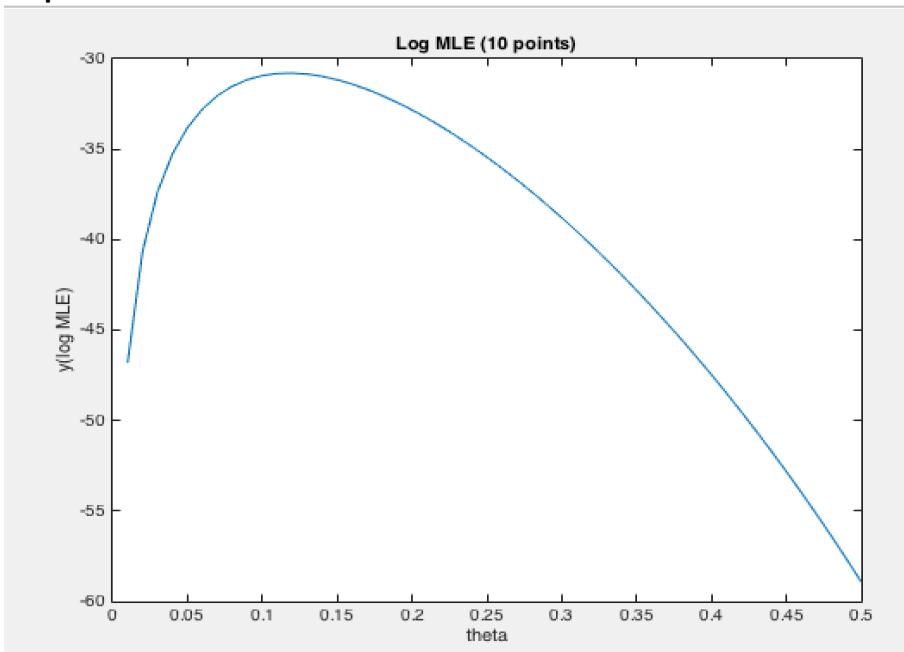
end
```

**Graphs : ( 2b)**

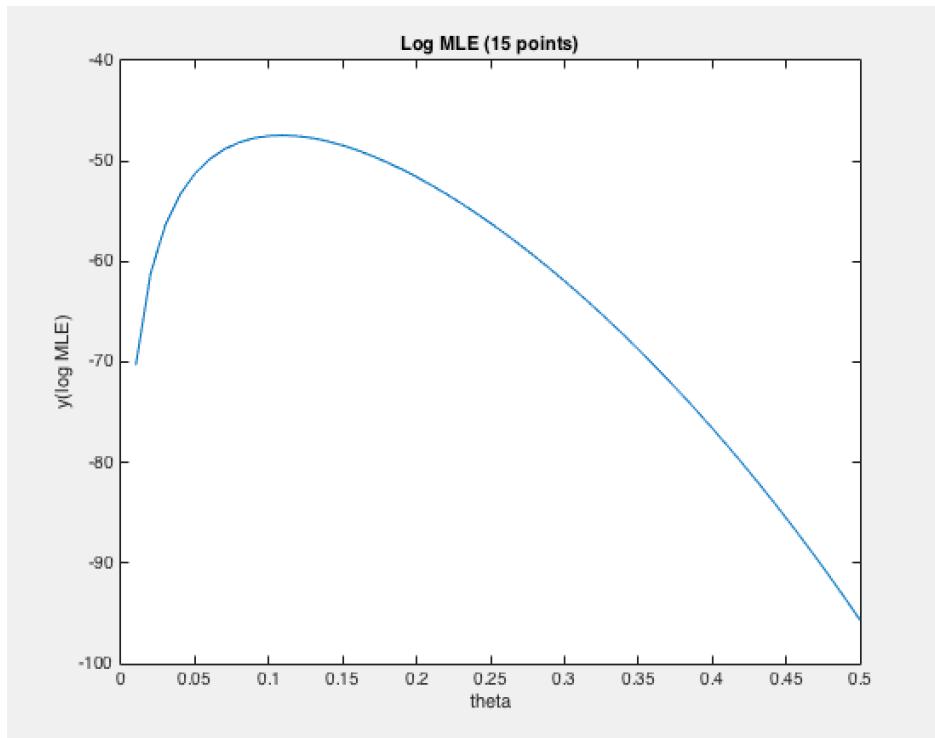
**5 Points**



**10 points**



**15 points**



Yes, the Closed form expression matches with the plots as the maximum of plot ( $\hat{\theta}$ ) is same as one derived from the previous equation.

Like for  $N=5 \quad \sum x_i = 27$ .

$$\hat{\theta} = \frac{5}{5+27} = \frac{5}{32} = 0.15$$

We are getting maximum log likelihood around the same point in the matlab plot so it is consistent.

Same is true with other two plots having 10 and 15 points.

for 10 points

$$\hat{\theta} = \frac{10}{10+75} = \frac{10}{85} = 0.11$$

for 15 points

$$\hat{\theta} = \frac{15}{15+123} = \frac{15}{138} = 0.10$$

All above values are consistent with matlab plots submitted.

d)  $L(\hat{\theta}) = \log(1-\theta) \sum_{i=1}^N x_i + N \log \theta$

as  $\theta < 1 \Rightarrow \log \theta & \log(1-\theta)$  both will be negative values.

as no. of sample increases,  $N$  increases or  $N \log \theta$  will

~~becomes~~ more negative.

also  $\sum_{i=1}^N x_i < \sum_{i=1}^N 1$ , and  $\log(1-\theta)$  is negative as  $(1-\theta) < 1$  &

~~base of log is greater than zero.~~

base of log is greater than zero.

so, if number of sample increases, log likelihood function became more negative.

### Problem 3) Implementing Naive Bayes

a)  $\hat{y} = \arg \max_y P(y|x)$

If  $x_i$ 's are conditionally independent given a class  $y$ ,

$$\Rightarrow P(x_1, x_2, y) = P(x_1|y) P(x_2|y)$$

Likewise for  $N$  features.

$$P(x_1, x_2, x_3, \dots, x_N, y) = P(x_1|y) P(x_2|y) P(x_3|y) \dots P(x_N|y) \quad \text{--- (1)}$$

$$\Rightarrow P(x_1, x_2, \dots, x_N|y) = \prod_{i=1}^N P(x_i|y) \quad \text{--- (2)}$$

$$P(y|x) P(x) = P(x|y) P(y)$$

$$\Rightarrow P(y|x) = \frac{P(x|y) P(y)}{P(x)} \xrightarrow{\text{normalization}}$$

$$\Rightarrow \hat{y} = \arg \max_y P(x|y) P(y)$$

$$\hat{y} = \arg \max_y P(x_1, x_2, \dots, x_N|y) P(y)$$

$$\boxed{\hat{y} = \prod_{i=1}^N P(x_i|y) P(y)} \quad \text{--- from (2)}$$

b) No. of parameters with naive Bayes assumption of conditionally independence would be  $[2n+1]$ , where  $n$  is the number of features.

If we don't assume Naive Bayes assumption, then total number of parameters would be  $(2^n - 1) \times 2$  as we would have 2 possibilities for every  $x_i$  and  $y$  i.e. 0 or 1, thus.

Explain it in the case of Naive Bayes i.e. conditional independence.

c), d), e)  $\rightarrow$  functions submitted on autolab.

f) Training error = .0034483 (2 mismatch out of 580)  
Test error = 0.027586 (2 mismatch out of 145)

On a new collection of set, we should expect test error as the best representative.

Since Naive Bayes is based on the training data, it will give comparatively good results for training data set when tested as test data matches with training data.

However, when we try different type of data (test data) we can expect high error rate since the distribution of test might be different from training data.

So, Naive Bayes attempt to minimize the training error.

g)

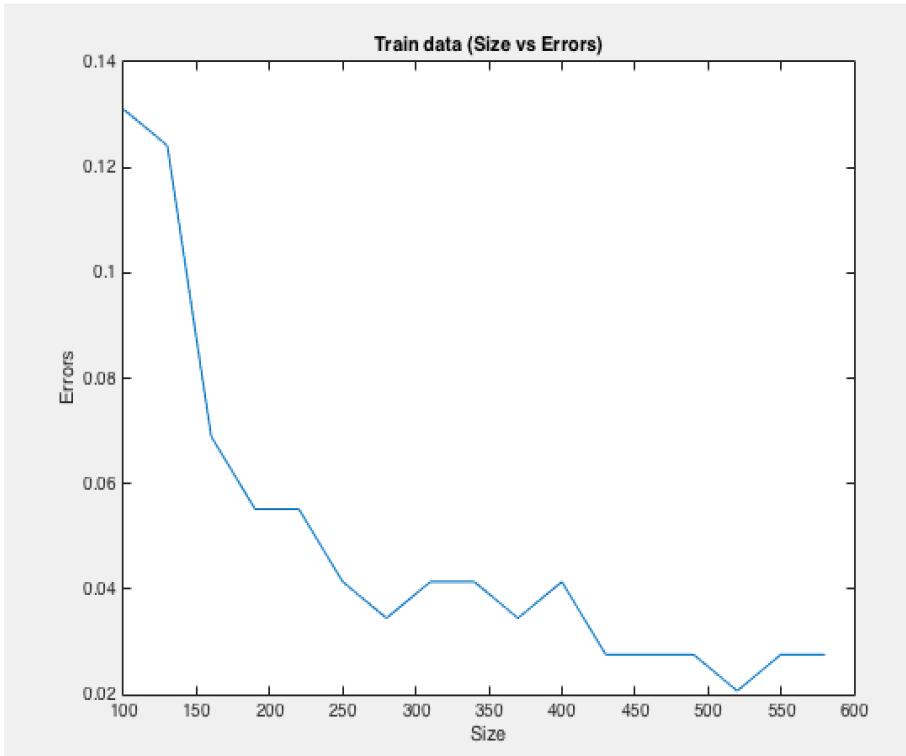
plots and Matlab code present in folder "Q3-g".  
Contains two subfolders i.e Test-Data-Error and Training-Data-Error.

Observation: In general we can say that as training data increase the classification error reduces or prediction gets better.

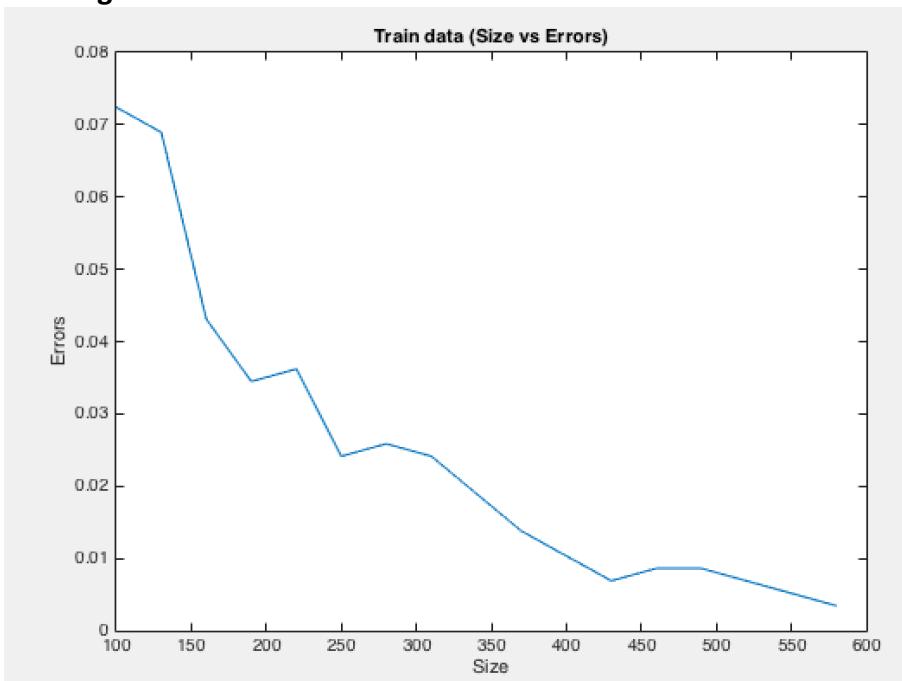
The error in the training data experiment is much lower than testing data for some no. of samples (documents). However, there are some glitches in the error i.e sometimes it goes up when I increased the samples by 36. This is observed both in training and test samples. Some spikes with increase in sample set could be because additional document might not have a distribution as existing sample thus causing distortion.

3 g)

Test data errors :



### Training Data Error:



h) Top five words.

A)  $y=1$  (Economist)

$$\text{Word Indexes} = \underline{32, 36, 46, 55, 12}$$

$y=2$  (Union)

$$\text{Word Indexes} = \underline{2, 24, 32, 36, 46} \rightarrow 'a', 'and', 'the', 'to', \cancel{'is'}$$

words  $\Rightarrow$  'the', 'to', 'of', 'in', 'a'

$$P(x=1 | y=1)$$

B) Top five words for  $P(x=1 | y=1) \rightarrow \text{economist}$

$$\frac{P(x=1 | y=1)}{P(x=1 | y=2)} \rightarrow \text{union}$$

$$P(x=1 | y=1)$$

$$P(x=1 | y=2)$$

$$\text{Word Indexes} = \underline{827, 1450, 3550, 199, 1516}$$

Top five words for  $\frac{P(x=1 | y=2)}{P(x=1 | y=1)}$   $\rightarrow$  origin, recession, favour, excess, labour

$$\text{Word Indexes} = \underline{20846, 20899, 21003, 20849, 2175}$$

4enlang, 5enlang, rethrift, realiz, coach

Words which has the largest ratio value for two given classes  
would have the highest contribution in classification or  
are most informative about class  $y$ .

These words are:

4enlang, 5enlang, rethrift, realiz, coach.

**Code for 3(g) :**

```
function [] = NB_Size_Error(XTest,yTest,XTrain,yTrain)
load('HW2Data.mat');
error = zeros([1 length(100:30:580)]);
count =1 ;
for m= 100:30:580
D = NB_XGivenY_Size(XTrain,yTrain,m);
p = sum(yTrain(1:m)==1)/length(yTrain);
fprintf('p = %d \n',p);
[yHat] = NB_Classify(D,p,XTest);

error(count)=sum((yTest~=yHat))/length(yTest);
fprintf('size = %d, error = %d \n',m,error(count));
count = count+1;
end

data_size = 100:30:580;

figure
plot(data_size,error)
xlabel('Size');
ylabel('Errors');
title('Train data (Size vs Errors)');
end
```

**Code for 3(h):**

```
function [] = NB_FreqWords(XTrain,yTrain)

load('HW2Data.mat');
[n, v] = size(XTrain);
D = zeros([2 v]);

word_count = zeros([2 v]);

[rows,col,value]= find(XTrain);

num_eco_docs = sum(yTrain==1);
num_oni_docs = sum(yTrain==2);
```

```

length_spare_matrix = length(rows);

for i = 1:length_spare_matrix
    if yTrain(rows(i),1)==1
        word_count(1,col(i)) = word_count(1,col(i))+1;
    else
        word_count(2,col(i)) = word_count(2,col(i))+1;
    end
end

D(1,:) = (word_count(1,:)+ .001)/(num_eco_docs + .901);
D(2,:) = (word_count(2,:)+ .001)/(num_oni_docs + .901);

B = D(1,:);
for i=1:5
    [M,I] = max(B);
    fprintf(' Economist : val = %d, idx = %d \n',M,I);
    B(1,I)=0;
end

B = D(2,:);
for i=1:5
    [M,I] = max(B);
    fprintf(' Onion : val = %d, idx = %d \n',M,I);
    B(1,I)=0;
end

B = D(1,:)./ D(2,:);
for i=1:5
    [M,I] = max(B);
    fprintf(' Eco/Onion : val = %d, idx = %d \n',M,I);
    B(1,I)=0;
end

B = D(2,:)./ D(1,:);
for i=1:5
    [M,I] = max(B);
    fprintf(' Onion/Eco : val = %d, idx = %d \n',M,I);
    B(1,I)=0;
end

end

```