

## Homework 3

Your Name: Srujan Shekar Shetty

Student ID: A20529733

---

### 1. [40] **Manually** solve the questions below by using Naïve Bayes Classifier

We conducted a survey to collect people's daily diets and try to build a model to predict whether their diets result in healthy conditions or not. The final results could be Yes, No. Note: using green rows as training, orange rows as testing.

Breakfast	Lunch	Dinner	Healthy?
Ham	Carnivorous	Beef	Y
Milk	Carnivorous	Beef	N
Bread	Veggie	Pork	N
Bread	Veggie	Veggie	Y
Ham	Veggie	Veggie	Y
Milk	Carnivorous	Pork	N
Bread	Carnivorous	Beef	N
Ham	Veggie	Pork	Y
Milk	Veggie	Pork	Y
Milk	Carnivorous	Veggie	N
Noddle	Carnivorous	Pork	?

1). [5 points] What is Laplace smoothing? And why we need it in the Naïve Bayesian classifier?

Ans Problem - Zero probability issue

Solution- Laplace smoothing.

How it works:

$$P(E | c_i) = P(e_1 \wedge e_2 \wedge \dots \wedge e_m | c_i) = \prod_{j=1}^m P(e_j | c_i)$$

From the above table if we are trying to predict that breakfast is noodles using a Naïve Bayesian classifier given the health. But in our training set, we have no breakfast as noodles. In this case, one among the  $e$  feature will be zero. Meaning the whole

$P(E | c_i)$  will be zero. This is called the Zero conditional probability problem.

To replace the zero value, we use the technology called Laplace smoothing.

Laplace smoothing:  $P(e_i = A | C = c_j) = (n_{c_j} + m \cdot p) / (n + m)$

- $A$  is a value in the  $i$ -th feature;  $c_j$  is a value in the label
- $n_{c_j}$  = # of training instances for which  $e_i = A$  and  $C = c_j$
- $n$  = # of training instances for which  $C = c_j$
- $m$  = a weight factor, usually  $m \geq 1$  and could be big value, for example, the size of your training data
- $p$  = an estimate or a probability value to decrease  $m$ , usually, we can set  $p$  as  $1/t$  and  $t$  is the number of unique values in  $e$

2). [15 points] Using the **Categorical Naive Bayesian Classification** to make predictions on the test sets, present confusion matrix, and calculate accuracy, precision, recall, F1 measure, specificity, by considering Y as positive label

$\therefore$  Breakfast-  $e_1$ , lunch  $e_2$ , dinner  $e_3$

(i) Bread (Carnivorous) Beef

$P(Y) = 3/6$   
 $P(N) = 3/6$

$P(Y|E) = \frac{P(Y) P(E|Y)}{P(E)} \quad \text{--- (1)}$

$P(N|E) = \frac{P(N) P(E|N)}{P(E)} \quad \text{--- (2)}$

We can only compare numerators of (1) & (2) To keep it simple as denominators are the same.

Form (1)  
 $P(Y|E) = P(Y) P(E|Y)$   
 $= \frac{3}{6} \times (P(e_1|Y) P(e_2|Y) P(e_3|Y))$   
 $= \frac{3}{6} \times (\frac{1}{3} \times \frac{1}{3} \times \frac{1}{3}) = \frac{1}{6} \times \frac{1}{3} \times \frac{1}{3} = 0.0185$

$P(N|E) = P(N) P(E|N) = (\frac{3}{6}) \times (\frac{1}{3} \times \frac{1}{3} \times \frac{1}{3}) = \frac{1}{27} = 0.037$

(ii) Ham Veggie Pork

Applying the above formula  
 $P(Y|E) = P(Y) P(E|Y)$   
 $= (\frac{3}{6}) \times (P.T.O)$

We have more confidence to say we can trust  $P(N|E) > P(Y|E)$  'N'  
 $(0.037) > (0.0185)$



(ii) Ham  $e_1$ , Veggie  $e_2$ , Pork  $e_3$

$$P(Y|E) = P(Y) (P(E|Y))$$

$$= \frac{3}{6} \times \left( \frac{2}{3} \times \frac{2}{3} \times \frac{0}{3} \right)$$

↓

Use Laplace smoothing to overcome this problem (Zero Conditional Probability problem)

$$P(e_i = A | C = G) = (n_i + m \cdot p) / (n + m)$$

for our case,

$$= \frac{3}{6} \times \left( \frac{2}{3} \times \frac{2}{3} \times \frac{0 + (m \cdot p)}{m + n} \right)$$

where  $m = 6$   
 $p = 1/3 = 1/3$   
 $n = 3$

$$= \frac{3}{6} \times \left( \frac{2}{3} \times \frac{2}{3} \times \frac{0 + (6 \cdot 1/3)}{6 + 3} \right)$$

$$= \frac{1}{2} \times \left( \frac{2}{3} \times \frac{2}{3} \times \frac{2}{3} \right)$$

$$P(Y|E) = \frac{1}{6} \times \left( \frac{2}{3} \times \frac{2}{3} \times \frac{2}{3} \right) = 0.4938$$

$$P(N|E) = \frac{1}{6} \times \left( \frac{0 + (6 \cdot 1/3)}{6 + 3} \right) \times \frac{1}{3} \times \frac{2}{3} = 0.246$$

$$\frac{3}{6} \times \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3}$$

∴  $P(Y|E) > P(N|E)$   
 So we trust 'Y' as Healthy label

(ignoring denominator as we compare it to  $P(N|E)$  having same denominator)

$$\frac{2 \times 2 \times 1}{6 \times 3} = 2/9$$

$$\frac{10}{3 \times 3} = \frac{10}{9}$$



(iii) Milk Veggie

Label	1	2
Milk	1	1
Veggie	1	0

$$P(Y|E) = P(Y) P(E|Y) = \frac{3}{6} \times \left( \frac{6 \times \frac{1}{3}}{6+3} \right) \times \frac{2}{3} \times \frac{1}{3} = \frac{6 \times \frac{1}{3}}{6+3}$$

$$= \frac{2}{3} \times \frac{2}{9} \times \frac{2}{3} \times \frac{2}{9} = 0.01646$$

$$P(N|E) = \frac{3}{6} \times \frac{2}{3} \times \frac{1}{3} \times \frac{2}{3} = 0.7407$$

(iv) Milk Carnivorous Veggie

$P(N|E) > P(Y|E)$   
 $\therefore$  We can trust 'N' as health label.

$$P(Y|E) = \frac{3}{6} \times \left( \frac{2}{9} \times \frac{1}{3} \times \frac{2}{3} \right) = 0.0246$$

$$P(N|E) = \frac{3}{6} \times \left( \frac{2}{3} \times \frac{2}{3} \times \frac{2}{9} \right) = 0.04938$$

$P(N|E) > P(Y|E)$   
 $\therefore$  We can trust 'N' as healthy label



Confusion Matrix

		Predicted labels	
		Y	N
Actual labels	Y	True positives (TP) 1	False Negatives (FN) 1
	N	False positives (FP) 0	True Negatives (TN) 2

$$\therefore \text{Accuracy} = (TP + TN) / \text{All} \\ = (1 + 2) / 4 = 3/4 = 75\%$$

$$\therefore \text{Precision} = (TP) / (TP + FP) \\ = (1 / 1 + 0) = 100\%$$

$$\therefore \text{Recall} = (TP) / (TP + FN) = 1 / 1 + 1 = 50\%$$

$$\therefore \text{F1 Measure} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{precision} + \text{Recall}} \\ = \frac{2 \times 100 \times 50}{100 + 50} = \frac{2 \times 1 \times 1/2}{1 + 1/2} \\ = \frac{1}{1 + 1/2} = \frac{1}{3/2} = 2/3 = 0.666 \\ = \underline{\underline{66\%}}$$

$\therefore$  Specificity considering Y as positive label

$$\text{Specificity} = TN / (FP + TN) \\ = 2 / (0 + 2) = 100\%$$

NOTE:- The last observation with '?' label is not used for the laplace calculation. ONLY training & testing dataset is used.

2). [20 points] Using the **Categorical Naive Bayesian Classification** to make prediction on the unseen data (note: building the model by using both the green and orange rows, and predicting the label for unseen data/last row)

2) (b) all the rows considering as training data

$E = \text{Noodle, Carnivorous, Pork}$   
 $C_1 = Y \quad C_2 = N$

$$P(\text{Noodle} | C_1) = \frac{m_c + m \times p}{m + n} = \frac{0 + 10 \times \frac{1}{4}}{10 + 5} = \frac{2.5}{15} = \frac{5}{30} = \frac{1}{6} = 0.1667$$

$m = \text{Total Rows}$   
 $p = 1/t$   
 $n = \text{no of 't' in labels}$

$$P(\text{Carnivorous} | C_1) = \frac{1}{5}$$

$$P(\text{Pork} | C_1) = \frac{2}{5}$$

$$P(E | C_1) = 0.1667 \times \frac{1}{5} \times \frac{2}{5} = 0.013336$$

$$P(\text{Noodle} | C_2) = \frac{0 + 10 \times \frac{1}{4}}{10 + 5} = 0.1667$$

$$P(\text{Pork} | C_2) = 2/5$$

$$P(\text{Carnivorous} | C_2) = 4/5$$

$$P(E | C_2) = 0.1667 \times 4/5 \times 2/5 = 0.05334$$

$$P(C_1) = 5/10 = 1/2 \quad ; \quad P(C_2) = 5/10 = 1/2$$



$$P(E) = \frac{1}{2} \times 0.01336 + \frac{1}{2} \times 0.053344$$

$$= \frac{1}{2} (0.01336 + 0.053344) = \underline{\underline{0.03334}}$$

$$P(C_1|E) = \frac{P(C_1)P(E|C_1)}{P(E)} = \frac{\frac{1}{2} \times 0.01336}{0.03334} = \underline{\underline{0.200}}$$

$$P(C_2|E) = \frac{P(C_2)P(E|C_2)}{P(E)} = \frac{\frac{1}{2} \times 0.053344}{0.03334} = \underline{\underline{0.800}}$$

$$P(C_2|E) > P(C_1|E)$$

∴ We can ~~simply~~ say we have more confidence that the label is 'N'.  
Healthy - 'N'



## 2. (60 points) Python practice for Naïve Bayes classification

Use the `Malware_MultiClass.csv` data (predicting the column “classification”), and run 5 Naïve Bayes techniques by using hold-out evaluations (80% as training)

Note:

- You need to change different/multiple parameters to find the best NB model.
- You should evaluate the models by using accuracy, micro-precision and micro-recall, micro-F1 and micro-AUC
- Give conclusions about the best model by comparing the evaluation metrics above

### Submission

- The ipynb and saved html files
- The comparison and conclusions of different models

Ans- Comparison

Model

Results of Categorical Naive Bayes using `df_binary`.

Hold-out Evaluation: Accuracy = 0.7458

Micro Precision = 0.7648960114630021

Micro Recall = 0.7460409171725716

Micro F1 = 0.7412574511263368

Micro AUC = 0.8057403055985901

Results of Categorical Naive Bayes using df\_binary.

Hold-out Evaluation: Accuracy = 0.7458

Micro Precision = 0.7648960114630021

Micro Recall = 0.7460409171725716

Micro F1 = 0.7412574511263368

Micro AUC = 0.8057403055985901

Results of Multinomial Naive Bayes using df\_num.

Hold-out Evaluation: Accuracy = 0.62225

Micro Precision = 0.6265632168368146

Micro Recall = 0.6222222310000007

Micro F1 = 0.6189721844020803

Micro AUC = 0.6038926643503398

Results of Gaussian Naive Bayes using df\_num.

Hold-out Evaluation: Accuracy = 0.6174

Micro Precision = 0.6203667501296704

Micro Recall = 0.6170361858205575

Micro F1 = 0.6145891369097447

Micro AUC = 0.7118165207678421

Results of Gaussian Naive Bayes using df\_num\_std.



Hold-out Evaluation: Accuracy = 0.49495

Micro Precision = 0.706251630151278

Micro Recall = 0.7024574526847485

Micro F1 = 0.7005634219234647

Micro AUC = 0.5

Results of Bernoulli Naive Bayes using df\_binary.

Hold-out Evaluation: Accuracy = 0.74935

Micro Precision = 0.7689996022120577

Micro Recall = 0.7500736021462386

Micro F1 = 0.7450445669649486

Micro AUC = 0.8097013358909546

Comparing the above metrics

Results of Bernoulli Naive Bayes using df\_binary.

Hold-out Evaluation: Accuracy = 0.74935

Gives more accuracy Precision recall and we can conclude this model is better comparatively.