

Name: _____

Homework 2: CSCI 347: Data Mining

Show your work. Include any code snippets you used to generate an answer, using comments in the code to clearly indicate which problem corresponds to which code.

Consider the following data matrix:

1. [5 points] Use matplotlib to create a bar plot for the counts of the variable X2. Make sure to label the axis.

Using Python:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
Data = [['red', 'yes', 'north'], ['blue', 'no', 'south'], ['yellow', 'no', 'east'], ['yellow', 'no', 'west'],
        ['red', 'yes', 'north'], ['yellow', 'yes', 'north'], ['blue', 'no', 'west']]
df = pd.DataFrame(Data, columns = ['X1', 'X2', 'X3'])
df['X2'].value_counts().plot(kind = 'bar', rot = 45)
plt.xlabel("Variable X_2")
plt.ylabel("Counts")
```

2. [2 points] Use one-hot encoding to transform **all** the categorical attributes to numerical values. Write down the transformed data matrix. Call this new matrix Y.

| | | X_{1R} | X_{1B} | X_{1Y} | X_{2Y} | X_{2N} | X_{3N} | X_{3S} | X_{3E} | X_{3W} |
|-------|-------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| $Y =$ | x_1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| | x_2 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| | x_3 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 |
| | x_4 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |
| | x_5 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| | x_6 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| | x_7 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |

3. [2 points] What is the Euclidean distance between data instance x_2 (second row) and data instance x_7 (seventh row) after applying one-hot encoding?

$$\|x_2 - x_7\|_2 = \sqrt{(0-0)^2 + (1-1)^2 + (0-0)^2 + (0-0)^2 + (1-1)^2 + (0-0)^2 + (1-0)^2 + (0-0)^2 + (0-1)^2}$$

$$= \sqrt{2}$$

4. [2 points] What is the cosine similarity (cosine of the angle) between data instance x_2 and data instance x_7 after applying one-hot encoding?

$$\cos(\theta) = \frac{0(0) + 1(1) + 0(0) + 0(0) + 1(1) + 0(0) + 1(0) + 0(0) + 0(1)}{\|x_2\|_2 \|x_7\|_2}$$

$$= \frac{2}{\sqrt{3}\sqrt{3}} = \frac{2}{3}$$

5. [2 points] What is the Hamming distance between data instance x_2 and data instance x_7 ?

$$\delta_H(x_2, x_7) = d - s = 3 - 2 = 1$$

where d is the number of categorical attributes and s is the number of matching values for categorical attributes in x_2 and x_7

6. [2 points] What is the Jaccard coefficient between data instance x_2 and data instance x_7 after applying one-hot encoding?

$$J(x_2, x_7) = \frac{s}{2d - s} = \frac{2}{6 - 2} = \frac{1}{2}$$

7. [2 points] What is the multivariate mean of Y

$$\left(\frac{2}{7} \quad \frac{2}{7} \quad \frac{3}{7} \quad \frac{4}{7} \quad \frac{4}{7} \quad \frac{3}{7} \quad \frac{1}{7} \quad \frac{1}{7} \quad \frac{2}{7} \right)$$

8. [2 points] What is the sample variance of the first column of Y (using the matrix written in the answer to (2))?

$$\begin{aligned} \hat{\sigma}_1^2 &= \left(\frac{5}{7}\right)^2 + \left(\frac{-2}{7}\right)^2 + \left(\frac{-2}{7}\right)^2 + \left(\frac{-2}{7}\right)^2 + \left(\frac{5}{7}\right)^2 + \left(\frac{-2}{7}\right)^2 + \left(\frac{-2}{7}\right)^2 \\ &= 1.429 \end{aligned}$$

9. [2 points] Write down the resulting matrix after applying standard (z-score) normalization to the matrix Y. Call this matrix Z.

Using Python:

```
from sklearn import preprocessing
import numpy as np
D = np.matrix('1 0 0 1 0 1 0 0 0; 0 1 0 0 1 0 1 0 0; 0 0 1 0 1 0 0 1 0; 0 0 1 0 1 0 0 0 0
1; 1 0 0 1 0 1 0 0 0; 0 0 1 1 0 1 0 0 0; 0 1 0 0 1 0 0 0 1')
standard_scaler = preprocessing.StandardScaler()
standard_normalized_D = standard_scaler.fit_transform(D)
```

$$Z = \begin{array}{c|ccccccccc} & X_{1R} & X_{1B} & X_{1Y} & X_{2Y} & X_{2N} & X_{3N} & X_{3S} & X_{3E} & X_{3W} \\ \hline x_1 & 1.58 & -0.63 & -0.87 & 1.15 & -1.15 & 1.15 & -0.41 & -0.41 & -0.63 \\ x_2 & -0.63 & 1.58 & -0.87 & -0.87 & 0.87 & -0.86 & 2.45 & -0.41 & -0.63 \\ x_3 & -0.63 & -0.63 & 1.15 & -0.87 & 0.87 & -0.87 & -0.41 & 2.45 & -0.63 \\ x_4 & -0.63 & -0.63 & 1.15 & -0.87 & 0.87 & -0.87 & -0.41 & -0.41 & 1.58 \\ x_5 & 1.58 & -0.63 & -0.87 & 1.15 & -1.15 & 1.15 & -0.41 & -0.41 & -0.63 \\ x_6 & -0.63 & -0.63 & 1.15 & 1.15 & -1.15 & 1.15 & -0.41 & -0.41 & -0.63 \\ x_7 & -0.63 & 1.58 & -0.87 & -0.87 & 0.87 & -0.87 & -0.41 & -0.41 & 1.58 \end{array}$$

[answers will vary depending on Y from question (2)]

10. [2 points] What is the (multivariate) mean of Z?

(0 0 0 0 0 0 0 0 0)

11. [2 points] Let z_i be the i th row of Z. What is the Euclidean distance between z_2 and z_7 ?

Using Python:

```
z2 = standard_normalized_D[1,:]
z7 = standard_normalized_D[6,:]
import numpy.linalg as LA
LA.norm(z2-z7)
```

3.16