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:

| | $X_1$ | $X_2$ | $X_3$ |
|---|---|---|---|
| $x_1$ | red | yes | North |
| $x_2$ | blue | no | South |
| $x_3$ | yellow | no | East |
| $x_4$ | yellow | no | West |
| $x_5$ | red | yes | North |
| $x_6$ | yellow | yes | North |
| $x_7$ | blue | no | West |

1.  [4 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.

2. [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?

3. [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?

4. [2 points] What is the Hamming distance between data instance $x_2$ and data instance $x_7$?

5. [2 points] What is the Jaccard coefficient between data instance $x_2$ and data instance $x_7$ after applying one-hot encoding?

6. [2 points] What is the (multivariate) mean of Y?

7. [2 points] What is the sample variance of the first column of Y (using the matrix written in the answer to (1) ) ?

8. [4 points] Write down the resulting matrix after applying standard (z-score) normalization to the matrix Y. Call this matrix Z.

9. [2 points] What is the (multivariate) mean of Z?

10. [2 points] Let $z_i$ be the $i$th row of Z. What is the Euclidean distance between $z_2$ and $z_7$?