

CSCI 347  
Homework 02

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

	$X_1$	$X_2$	$X_3$
$x_1$	<i>red</i>	<i>yes</i>	<i>north</i>
$x_2$	<i>blue</i>	<i>no</i>	<i>south</i>
$x_3$	<i>yellow</i>	<i>no</i>	<i>east</i>
$x_4$	<i>yellow</i>	<i>no</i>	<i>west</i>
$x_5$	<i>red</i>	<i>yes</i>	<i>north</i>
$x_6$	<i>yellow</i>	<i>yes</i>	<i>north</i>
$x_7$	<i>blue</i>	<i>no</i>	<i>west</i>

Answer the following:

1. (5 points) Use matplotlib to create a bar plot for the counts of the variable  $X_2$ . Make sure to label the axis.
2. (2 points) Use one-hot encoding to transform all the categorical attributes to numerical values. Write down the transformed data matrix. (In what follows, we will refer to the transformed data matrix as  $Y$ ).
3. (2 points) What is the Euclidean distance between instance  $x_2$  (second row) and  $x_7$  (seventh row) after applying one-hot encoding.
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?
5. (2 points) What is the Hamming distance between data instance  $x_2$  and data instance  $x_7$  after applying one-hot encoding?
6. (2 points) What is the Jaccard similarity between data instance  $x_2$  and  $x_7$  after applying one-hot encoding?
7. (2 points) What is the multi-dimensional mean of  $Y$ ?
8. (2 points) What is the estimated variance of the first column of  $Y$ ?
9. (2 points) What is the resulting matrix after applying standard (z-score) normalization to the matrix  $Y$ . In the following, we will call this matrix  $Z$ .
10. (2 points) What is the multi-dimensional mean of  $Z$ ?
11. (2 points) Let  $z_i$  be the  $i$ -th row of  $Z$ . What is Euclidean distance between  $z_2$  and  $z_7$ ?

**Acknowledgements:** Homework problems adapted from assignments of Veronika Strnadova-Neeley.