

CSCI 347
Homework 02

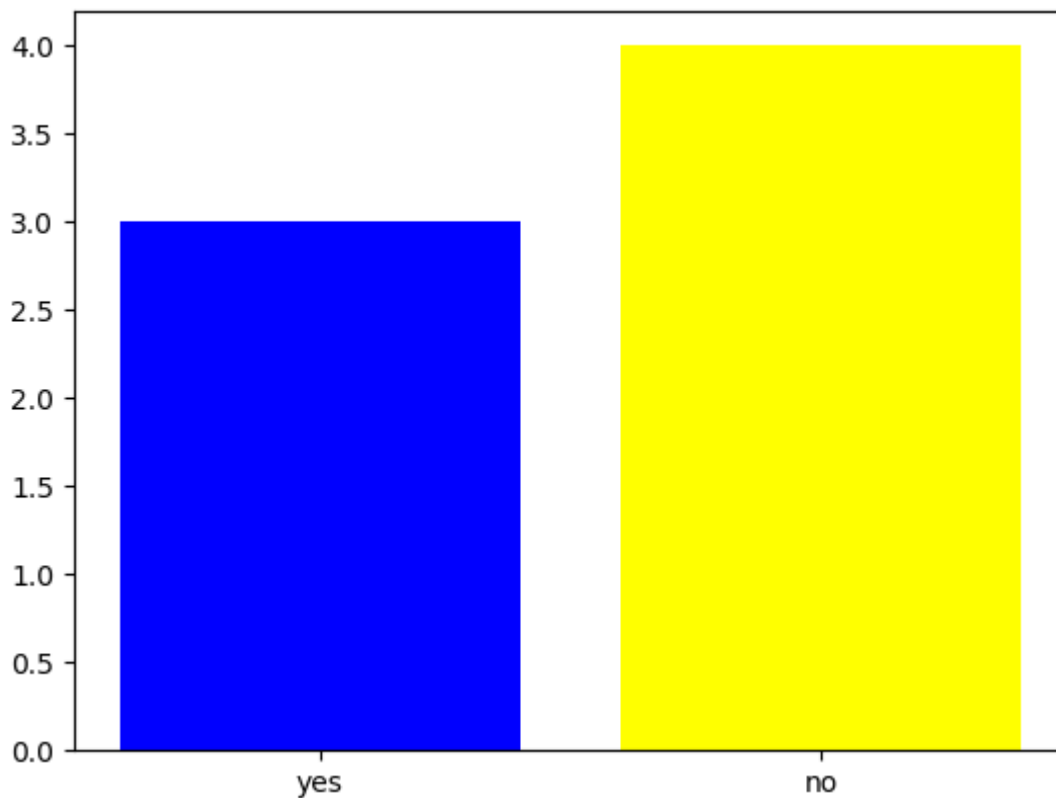
Show your work. Include any code snippets you used to generate an answer, using comments in the code to clearly indicate which problem corresponds to which code

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.



```

import numpy as np
import matplotlib.pyplot as plt

d = np.array([[ 'red ', 'yes ', 'north '], [ 'blue ', 'no ', 'south '],
               [ 'yellow ', 'no ', 'east '], [ 'yellow ', 'no ', 'west '],
               [ 'red ', 'yes ', 'north '], [ 'yellow ', 'yes ', 'north '],
               [ 'blue ', 'no ', 'west ']])

X2 = d[:,1]
countY = sum(X2 == 'yes')
countN = sum(X2 == 'no')

plt.bar(x = ('yes ', 'no '), height = (countY, countN),
        color = ('blue ', 'yellow '))
plt.show()

```

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 referred to the transformed data matrix as Y).

	X_{1r}	X_{1b}	X_{1y}	X_{2y}	X_{2n}	X_{3n}	X_{3e}	X_{3s}	X_{3w}
x_1	1	0	0	1	0	1	0	0	0
x_2	0	1	0	0	1	0	0	1	0
x_3	0	0	1	0	1	0	1	0	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 instance x_2 (second row) and x_7 (seventh row) after applying one-hot encoding.

$$\sqrt{(0-0)^2 + (1-1)^2 + (0-0)^2 + (0-0)^2 + (1-1)^2 + (0-0)^2 + (0-0)^2 + (1-0)^2 + (0-1)^2} = \sqrt{2} = \mathbf{1.4142}$$

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?

$$x_2 \cdot x_7 = 0 * 0 + 1 * 1 + 0 * 0 + 0 * 0 + 1 * 1 + 0 * 0 + 0 * 0 + 1 * 0 + 0 * 1 = 2$$

$$||x_2|| = \sqrt{(0)^2 + (1)^2 + (0)^2 + (0)^2 + (1)^2 + (0)^2 + (0)^2 + (1)^2 + (0)^2} = \sqrt{3}$$

$$||x_7|| = \sqrt{(0)^2 + (1)^2 + (0)^2 + (0)^2 + (1)^2 + (0)^2 + (0)^2 + (0)^2 + (1)^2} = \sqrt{3}$$

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

5. (2 points) What is the Hamming distance between data instance x_2 and data instance x_7 after applying one-hot encoding?

$$0 \oplus 0 + 1 \oplus 1 + 0 \oplus 0 + 0 \oplus 0 + 1 \oplus 1 + 0 \oplus 0 + 0 \oplus 0 + 1 \oplus 0 + 0 \oplus 1 \\ = 0 + 0 + 0 + 0 + 0 + 0 + 0 + 0 + 1 + 1 = \mathbf{2}$$

6. (2 points) What is the Jaccard similarity between data instance x_2 and x_7 after applying one-hot encoding?

$$\frac{0 \wedge 0 + 1 \wedge 1 + 0 \wedge 0 + 0 \wedge 0 + 1 \wedge 1 + 0 \wedge 0 + 0 \wedge 0 + 1 \wedge 0 + 0 \wedge 0 + 1 \wedge 1}{0 \vee 0 + 1 \vee 1 + 0 \vee 0 + 0 \vee 0 + 1 \vee 1 + 0 \vee 0 + 0 \vee 0 + 1 \vee 0 + 0 \vee 0 + 1 \vee 1} \\ = \frac{0+1+0+0+1+0+0+0+0}{0+1+0+0+1+0+0+1+1} = \frac{1}{2}$$

7. (2 points) What is the multivariate mean of Y ?

```
import pandas as pd
```

```
data = [[1, 0, 0, 1, 0, 1, 0, 0, 0],
        [0, 1, 0, 0, 1, 0, 0, 1, 0],
        [0, 0, 1, 0, 1, 0, 1, 0, 0],
        [0, 0, 1, 0, 1, 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]]
```

```
df = pd.DataFrame(data)
```

```
x = df.mean()
print(x)
```

multivariate mean = **0.285714, 0.285714, 0.428571, 0.428571, 0.571429, 0.428571, 0.142857, 0.142857, 0.285714**

8. (2 points) What is the estimated variance of the first column of Y ?

$$\frac{1}{7}(1 + 0 + 0 + 0 + 1 + 0 + 0) = 0.286 \\ \frac{1}{6}(1 - 0.286)^2 + (0 - 0.286)^2 + (0 - 0.286)^2 + (0 - 0.286)^2 + (1 - 0.286)^2 + (0 - 0.286)^2 + (0 - 0.286)^2) \\ = \mathbf{0.238}$$

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 .

```
import numpy as np

data = [[1, 0, 0, 1, 0, 1, 0, 0, 0],
        [0, 1, 0, 0, 1, 0, 0, 1, 0],
        [0, 0, 1, 0, 1, 0, 1, 0, 0],
        [0, 0, 1, 0, 1, 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]]
```

```
myArray = np.array(data)
```

```
mean = np.mean(myArray, axis = 0)
```

```
std = np.std(myArray, axis = 0)
```

```
print(mean)
```

```
print("-----")
```

```
print(std)
```

```
print("-----")
```

```
Z = ((myArray-mean)/std)
```

```
print(Z)
```

	X_{1r}	X_{1b}	X_{1y}	X_{2y}	X_{2n}	X_{3n}	X_{3e}	X_{3s}	X_{3w}
x_1	1.58	-0.63	-0.86	1.15	-1.15	1.15	-0.41	-0.41	-0.63
x_2	-0.63	1.58	-0.86	-0.86	0.86	-0.86	-0.41	2.45	-0.63
x_3	-0.63	-0.63	1.15	-0.86	0.86	-0.86	2.45	-0.41	-0.63
x_4	-0.63	-0.63	1.15	-0.86	0.86	-0.86	-0.41	-0.41	1.58
x_5	1.58	-0.63	-0.86	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	1.15	-0.86	-0.86	0.86	-0.86	-0.41	-0.41	1.58

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

```
Zmean = np.mean(Z, axis = 0)
```

```
print(Zmean)
```

```
[3.17206578e-17 3.17206578e-17 3.17206578e-17 1.58603289e-17 4.75809868e-17 1.58603289e-17
 6.34413157e-17 3.17206578e-17 3.17206578e-17] = [0 0 0 0 0 0 0]
```

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

```
temp = Z[1] - Z[6]
sumsq = np.dot(temp.T, temp)
print(np.sqrt(sumsq))
```

= 3.614784

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