

CSCI 5822 Assignment 1

Assigned: Jan 22, 2021. Due: Feb 5 (11:59pm, MST), on Canvas.

You have the option to solve the assignment in place and upload the notebook as a PDF (via "File" -> "Download as" -> "PDF via LaTeX (.pdf)") file. OR, use the notebook as a worksheet and upload your answers as a separated, typed PDF file. You may get partial points if you upload the notebook to that file.

Do not forget to edit your name here.

Name: Alex Book Collaborated With: Cole Sturza

The collaboration/anti-plagiarism policy for this course can be found in the syllabus statements.

The questions provided below will ask you to either write code or write answers in the form of markdown. Markdown syntax guide is here: click here. Using markdown you can typeset formulae using latex. This way you can write nice readable answers with formulae like this:

A Markov chain is of the form  $p(x_T) \prod_{i=1}^{T-1} p(x_i|x_{i+1})$  for some assignment of the variables to labels  $x_i$ .

We will follow the following grading rubrics.

- 10: correct answer, work is shown and clear
- 8: correct answer but work might be difficult to read
- 6: incorrect answer but good attempt
- 4: mediocre attempt or very difficult to read
- 2: problem is written down, no attempt to solve
- 0: problem is not written down

Double click anywhere on this box to find out how your instructor typeset it. Press Shift+Enter to go back.

The Tasks 1-5 are from Exercises 1.1-1.6 of Barber.

Task 1: (10 points). Prove

$$p(x,y,z) = p(x) p(y|x,z)$$

and also

$$p(x|y,z) = \frac{p(y|x,z)p(x|z)}{p(y|z)}.$$

your answer to task 1

$$p(x,y,z) = \frac{p(x,y,z)}{p(z)}$$

Conditional Probability (1)

$$= \frac{p(y,x,z)}{p(z)}$$

Commutativity of Probability (2)

$$= \frac{p(y|x,z)p(x|z)p(z)}{p(z)}$$

Chain Rule of Probability (3)

$$= \frac{p(y|x,z)p(x|z)}{p(y,z)}$$

Commutativity of Multiplication (4)

$$p(x|y,z) = \frac{p(x,y,z)}{p(y,z)}$$

Conditional Probability (1)

$$= \frac{p(x,y|z)p(z)}{p(y,z)}$$

Chain Rule of Probability (2)

$$= \frac{p(x,y|z)p(z)}{p(y|z)p(z)}$$

Chain Rule of Probability (3)

$$= \frac{p(x,y|z)}{p(y|z)}$$

(4)

$$= \frac{p(x|z)p(y|x,z)}{p(y|z)}$$

Proven in first proof of task 1 (5)

$$= \frac{p(y|x,z)p(x|z)}{p(y|z)}$$

Commutativity of Multiplication (4)

Task 2: (10 points). Prove the [Bonferroni inequality](https://en.wikipedia.org/wiki/Bonferroni\_inequality)  $p(a,b) \geq p(a) + p(b) - 1$ .

your answer to task 2

Note that  $p(a,b) = p(a \cap b)$ .

$$p(a \cup b) = p(a) + p(b) - p(a \cap b)$$

True by definition, any probability must be less than or equal to 1 (2)

$$1 \geq p(a \cup b)$$

Combined (1) and (2) (3)

$$p(a \cap b) \geq p(a) + p(b) - 1$$

(4)

$$p(a,b) \geq p(a) + p(b) - 1$$

(4)

Task 3: (10 points). Consider three variable distributions which admit the factorization  $p(a,b,c) = p(a)b|b|c)p(c)$ ,

your answer to task 3

In order to represent the distribution of  $p(a|b)$ , two parameters are needed. Say,  $p(a = 1|b = 0)$  and  $p(a = 1|b = 1)$ . Using complementation, all other variations of  $p(a|b)$  can be found. Similarly, two parameters are needed to represent the distribution of  $p(b|c)$ . Lastly, only one parameter is needed to represent the distribution of  $p(c)$ . So, five parameters in total are needed to represent the distributions of the given form.

Task 4: (20 points, 10 each). (a) There are two boxes. Box 1 contains three red and five white balls and box 2 contains two red and five white balls. A box is chosen at random  $p(b \oplus 1) = p(b \oplus 2) = 0.5$  and a ball chosen at random from this box turns out to be red. What is the posterior probability that the red ball came from box 1?

(b) Two balls are placed in a box as follows. A fair coin is tossed and a white ball is placed in the box if a head occurs, otherwise a red ball is placed in the box. The coin is tossed again and a red ball is placed in the box if a tail occurs, otherwise a white ball is placed in the box. Balls are drawn from the box three times in succession (always with replacing the drawn ball back in the box). It is found that on all three occasions a red ball is drawn. What is the probability that both balls in the box are red?

your answer to task 4(a)

Let  $r$  represent a ball being red and  $b$  represent a ball coming from box 1.

$$p(b|r) = \frac{p(r|b)p(b)}{p(r)}$$

$$= \frac{\frac{2}{3} \times \frac{1}{2}}{(\frac{1}{3} + \frac{2}{3}) \times \frac{1}{2}} = \frac{\frac{1}{3}}{\frac{1}{2}} = \frac{2}{3} \approx .667$$

your answer to task 4(b)

Let any capitalized number in text represent that number of red balls drawn in a row (THREE = 3 red balls drawn in a row),  $RR$  represent two red balls in the box,  $WW$  represent two white balls in the box, and  $RW$  represent one ball of each color in the box.

$$p(RR|THREE) = \frac{p(THREE|RR)p(RR)}{p(THREE)}$$

$$= \frac{(1)p(RR)}{p(THREE)}$$

$$= \frac{(.5)(.5)}{.25}$$

$$= \frac{p(THREE|WW)p(WW) + p(THREE|RW)p(RW) + p(THREE|RR)p(RR)}{.25}$$

$$= \frac{0 \times .25 + .5^3 \times .5 + 1 \times .25}{.25}$$

$$= \frac{.0625 + .25}{.25}$$

$$= .8$$

Task 5: (10 points). A secret government agency has developed a scanner which determines whether a person is a terrorist. The scanner is fairly reliable; 95% of all scanned terrorists are identified as terrorists, and 95% of all upstanding citizens are identified as such. An informant tells the agency that exactly one passenger of 100 aboard an aeroplane in which you are seated is a terrorist. The police had off the plane the first person for which the scanner tests positive. What is the probability that this person is a terrorist?

your answer to task 5

Let  $T$  represent whether or not someone is a terrorist,  $t$  represent whether or not someone tested positive for terrorism, and  $H$  represent whether or not someone got hauled off the plane.

"Correct" Version

$$p(T = 1|H = 1) = \frac{\sum_{i=1}^{100} p(T_i = 1, t_i = 1, t_{1..i-1} = 0|H = 1)}{\sum_{i=1}^{100} p(H = 1|T_i = 1, t_i = 1, t_{1..i-1} = 0)p(T_i = 1, t_i = 1, t_{1..i-1} = 0)}$$

$$= \frac{p(H = 1|T_1 = 1, t_1 = 1, t_{1..i-1} = 0)p(T_1 = 1, t_1 = 1, t_{1..i-1} = 0)}{p(H = 1)}$$

$$= \frac{p(t_1 = 1|T_1 = 1, t_{1..i-1} = 0)p(T_1 = 1|p(T_1 = 1))}{p(H = 1)}$$

$$= \frac{p(t_1 = 1|T_1 = 1)p(T_1 = 1) \prod_{i=1}^{i-1} (t_i = 0|T_i = 0)}{p(H = 1)}$$

$$= \frac{\sum_{i=1}^{100} (.95)(.01) \prod_{i=1}^{i-1} (.95)}{p(H = 1)}$$

$$= \frac{p(H = 1)}{\sum_{i=1}^{100} (.95)(.01)(.95)^{i-1}}$$

$$= \frac{(.95)(.01) \sum_{i=1}^{100} (.95)^{i-1}}{p(H = 1)}$$

$$= \frac{(.95)(.01) \sum_{i=0}^{99} (.95)^i}{p(H = 1)}$$

$$= \frac{(.95)(.01) (1 - .95^{100})}{p(H = 1) (1 - .95^{100})}$$

$$= \frac{(.05)(.01)p(H = 1)}{p(H = 1)}$$

$$= \frac{(.19)(1 - .95^{100})}{p(H = 1)}$$

$$= \frac{1 - p(H = 0)}{p(H = 1)}$$

$$p(H = 0) = .95^{99} (.05), \text{ as it is equal to the probability that all passengers test negative}$$

$$= \frac{(.19)(1 - .95^{100})}{1 - (.95^{99} (.05))}$$

$$\approx .189$$

Breast Cancer Version

$$p(T = 1|t = 1) = \frac{p(t = 1|T = 1)p(T = 1)}{p(t = 1)}$$

$$= \frac{.95 \times .01}{.95 \times .01 + .05 \times .99}$$

$$= \frac{.95 \times .01}{.95 \times .01 + .05 \times .99}$$

$$= .161$$

The goal of the following tasks is to give you a bit of practice manipulating data, using Bayes' rule, and constructing a naive Bayes classifier. Naive Bayes is described in 10.1 of Barber and understanding examples 10.1 and 10.2 of the text should help you do this assignment.

Dataset: The [titanic dataset](http://www.cs.toronto.edu/~dave/data/titanic/desc.html) gives the values of four categorical attributes for each of the 2201 people on board the Titanic when it struck an iceberg and sank. The attributes are social class (first class, second class, third class, crew member), age (adult or child), gender, and whether or not the person survived. The titanic dataset is available [here](https://home.cs.colorado.edu/~nozer/Teaching/sf/ProbabilisticModels/homeework/titanic.txt).

Preparation: Build a joint probability table, like the ones we discussed in class notes, that represents the joint distribution over all variables, i.e.  $P_r(\text{Gender}, \text{Age}, \text{Class}, \text{Outcome})$ . This table should have 32 entries because  $\text{Gender} \in \{\text{male}, \text{female}\}$ ,  $\text{Age} \in \{\text{child}, \text{adult}\}$ ,  $\text{Class} \in \{\text{1st}, \text{2nd}, \text{3rd}, \text{crew}\}$ , and  $\text{Outcome} \in \{\text{death}, \text{survival}\}$ . You will use the data in this table for the following tasks.

Task 6: Probability table (20 points). Build a probability table indicating  $P_r(\text{Death} | \text{Gender}, \text{Age}, \text{Class})$  for each combination of gender, age, and class. Display this table in the following way:

	Male		Female	
	Child	Adult	Child	Adult
First				
Second				
Third				
Crew				

The rows of each table represent the different classes and the columns the different ages and genders. In each cell of the table, insert the conditional probability. After you've built the probability table, come up with a rule that uses the probabilities to predict death or survival. Then make a second table, a classification table, which lets death or survival for each feature combination. Explain the rule you chose to classify.

Hint: \* Since it is [not possible](https://www.tablesgenerator.com/html\_tables) to create merged rows with Markdown, this "Tables Generator"([https://www.tablesgenerator.com/html\_tables]) may be useful. \* Be alert to the possibility of a cell containing no data.

In [2]:

# your answer to task 6

import numpy as np

import pandas as pd

titanic\_df = pd.read\_csv('titanic.txt', delim\_whitespace=True, header=None, names = ['Class', 'Age', 'Gender', 'Outcome'])

titanic\_df.value\_counts(['Class', 'Age', 'Gender', 'Outcome'])/titanic\_df.value\_counts(['Class', 'Age', 'Gender'])

Out [2]:

Class	Age	Gender	Outcome	
1st	adult	female	yes	0.972222
		female	no	0.027778
	male	no		0.674286
	child	female	yes	0.325714
		female	yes	1.000000
2nd	adult	male	yes	1.000000
		female	yes	0.000215
		no		0.139785
	male	no		0.918887
		yes		0.085333
	child	female	yes	1.000000
		male	yes	1.000000
3rd	adult	female	no	0.539394
		yes		0.460606
	male	no		0.837662
		yes		0.162338
	child	female	no	0.540387
		yes		0.451613
	male	no		0.729187
		yes		0.270813
crew	adult	female	yes	0.869565
		no		0.130435
	male	no		0.777262
		yes		0.222738
				dtype: float64

Probability table

	Male		Female	
	Child	Adult	Child	Adult
First	0	.674	0	.028
Second	0	.817	0	.140
Third	.729	.737	.548	.539
Crew	0	.778	0	.130

Classification table

Rule: Anyone with chances of death of greater than a coin flip (greater than .5) is given a 100% chance of death, while anyone with chances of death of less than a coin flip (less than .5) is given a 0% chance of death.

	Male		Female	
	Child	Adult	Child	Adult
First	0	1	0	0
Second	0	1	0	0
Third	1	1	1	1
Crew	0	1	0	0

Task 7: Build a Naive Bayes classifier (20 points). To build the classifier, you must first construct six one-dimensional tables:

- $P_r(\text{Class} | \text{death})$
- $P_r(\text{Age} | \text{death})$
- $P_r(\text{Gender} | \text{death})$
- $P_r(\text{Class} | \text{survival})$
- $P_r(\text{Age} | \text{survival})$
- $P_r(\text{Gender} | \text{survival})$

To be clear on this notation, for  $P_r(\text{Age} | \text{death})$ , your table should have two rows, one for adult and one for child, and you should compute, for each age group, the probability of the deceased being in that age group. Also compute the unconditional probabilities,  $P_r(\text{death})$  and  $P_r(\text{survival})$ , with  $P_r(\text{death}) + P_r(\text{survival}) = 1$ . From this information, compute  $P_r(\text{death} | \text{Gender}, \text{Age}, \text{Class})$  using the Naive Bayes assumption. In addition to the probability table, build the classification table as well.

Naive Bayes Equation

$$p(\text{death} | \text{Gender}, \text{Age}, \text{Class}) = \frac{p(\text{Gender}, \text{Age}, \text{Class} | \text{death}) p(\text{death})}{p(\text{Gender}, \text{Age}, \text{Class})}$$

$$= \frac{p(\text{Gender} | \text{death}) p(\text{Age} | \text{death}) p(\text{Class} | \text{death}) p(\text{death})}{p(\text{Gender}, \text{Age}, \text{Class})}$$

$$= \frac{p(\text{Gender} | \text{death}) p(\text{Age} | \text{death}) p(\text{Class} | \text{death}) p(\text{death})}{\sum_{\text{Outcome}} p(\text{Gender} | \text{Outcome}) p(\text{Age} | \text{Outcome}) p(\text{Class} | \text{Outcome}) p(\text{Outcome})}$$

In [58]:

# your answer to task 7

num\_death = len(titanic\_df[titanic\_df['Outcome'] == 'no'])

num\_survive = len(titanic\_df[titanic\_df['Outcome'] == 'yes'])

pr\_death = num\_death/(num\_death+num\_survive)

pr\_survive = num\_survive/(num\_death+num\_survive)

classes = titanic\_df[['Class']].unique()

class\_given\_death = np.zeros(len(classes))

class\_given\_survive = np.zeros(len(classes))

for i in range(len(classes)):

class\_given\_death[i] = len(titanic\_df[titanic\_df['Outcome'] == 'no' & (titanic\_df['Class'] == classes[i])]) / num\_death

class\_given\_survive[i] = len(titanic\_df[titanic\_df['Outcome'] == 'yes' & (titanic\_df['Class'] == classes[i])]) / num\_survive

d\_class = {'Pr(Class | death)': class\_given\_death, 'Pr(Class | survive)': class\_given\_survive}

class\_df = pd.DataFrame(data=d\_class, index=classes)

#####

ages = titanic\_df[['Age']].unique()

age\_given\_death = np.zeros(len(ages))

age\_given\_survive = np.zeros(len(ages))

for i in range(len(ages)):

age\_given\_death[i] = len(titanic\_df[titanic\_df['Outcome'] == 'no' & (titanic\_df['Age'] == ages[i])]) / num\_death

age\_given\_survive[i] = len(titanic\_df[titanic\_df['Outcome'] == 'yes' & (titanic\_df['Age'] == ages[i])]) / num\_survive

d\_age = {'Pr(Age | death)': age\_given\_death, 'Pr(Age | survive)': age\_given\_survive}

age\_df = pd.DataFrame(data=d\_age, index=ages)

#####

genders = titanic\_df[['Gender']].unique()

gender\_given\_death = np.zeros(len(genders))

gender\_given\_survive = np.zeros(len(genders))

for i in range(len(genders)):

gender\_given\_death[i] = len(titanic\_df[titanic\_df['Outcome'] == 'no' & (titanic\_df['Gender'] == genders[i])]) / num\_death

gender\_given\_survive[i] = len(titanic\_df[titanic\_df['Outcome'] == 'yes' & (titanic\_df['Gender'] == genders[i])]) / num\_survive

d\_gender = {'Pr(Gender | death)': gender\_given\_death, 'Pr(Gender | survive)': gender\_given\_survive}

gender\_df = pd.DataFrame(data=d\_gender, index=genders)

#####

death\_given\_class\_age\_gender = np.zeros((len(classes), len(ages), len(genders)))

for i in range(len(classes)):

for j in range(len(ages)):

for k in range(len(genders)):

numerator = class\_given\_death[i]\*age\_given\_death[j]\*gender\_given\_death[k]\*pr\_death

denominator = numerator + class\_given\_survive[i]\*age\_given\_survive[j]\*gender\_given\_survive[k]\*pr\_survive

death\_given\_class\_age\_gender[i, j, k] = numerator/denominator

# classes, ages, genders

# death\_given\_class\_age\_gender

first\_adult\_male = death\_given\_class\_age\_gender[0,0,0]

first\_adult\_female = death\_given\_class\_age\_gender[0,0,1]

first\_child\_male = death\_given\_class\_age\_gender[0,1,0]

first\_child\_female = death\_given\_class\_age\_gender[0,1,1]

second\_adult\_male = death\_given\_class\_age\_gender[1,0,0]

second\_adult\_female = death\_given\_class\_age\_gender[1,0,1]

second\_child\_male = death\_given\_class\_age\_gender[1,1,0]

second\_child\_female = death\_given\_class\_age\_gender[1,1,1]

third\_adult\_male = death\_given\_class\_age\_gender[2,0,0]

third\_adult\_female = death\_given\_class\_age\_gender[2,0,1]

third\_child\_male = death\_given\_class\_age\_gender[2,1,0]

third\_child\_female = death\_given\_class\_age\_gender[2,1,1]

crew\_adult\_male = death\_given\_class\_age\_gender[3,0,0]

crew\_adult\_female = death\_given\_class\_age\_gender[3,0,1]

crew\_child\_male = death\_given\_class\_age\_gender[3,1,0]

crew\_child\_female = death\_given\_class\_age\_gender[3,1,1]

# print('first\_adult\_male:', first\_adult\_male)

# print('first\_adult\_female:', first\_adult\_female)

# print('first\_child\_male:', first\_child\_male)

# print('first\_child\_female:', first\_child\_female)

# print('second\_adult\_male:', second\_adult\_male)

# print('second\_adult\_female:', second\_adult\_female)

# print('second\_child\_male:', second\_child\_male)

# print('second\_child\_female:', second\_child\_female)

# print('third\_adult\_male:', third\_adult\_male)

# print('third\_adult\_female:', third\_adult\_female)

# print('third\_child\_male:', third\_child\_male)

# print('third\_child\_female:', third\_child\_female)

# print('crew\_adult\_male:', crew\_adult\_male)

# print('crew\_adult\_female:', crew\_adult\_female)

# print('crew\_child\_male:', crew\_child\_male)

# print('crew\_child\_female:', crew\_child\_female)

Probability table

	Male		Female	
	Child	Adult	Child	Adult
First	0	.528	.047	.099
Second	522	.725	.094	.206
Third	.696	.647	.184	.352
Crew	.710	.855	.195	.368

Classification table

Rule: Anyone with chances of death of greater than a coin flip (greater than .5) is given a 100% chance of death, while anyone with chances of death of less than a coin flip (less than .5) is given a 0% chance of death.

	Male		Female	
	Child	Adult	Child	Adult
First	0	1	0	0
Second	1	1	0	0
Third	1	1	0	0
Crew	1	1	0	0

Task 8: Discussions (10 points, bonus) The classification tables you built in Tasks 6 and 7 are not identical. Discusses the advantages/disadvantages of each table for making predictions in case of another disaster like the Titanic (assuming it occurred at the same time in history). Under what circumstances would you expect an empirical table to provide better predictions? Under what circumstances would you expect the naive Bayes table to provide better predictions?

your answer to task 8

An empirical table may provide better predictions when there are similar passenger distributions as there were on Titanic (which seems to have been a fairly typical distribution). For example, it works better when there are no child crew members (which one would hope to be true in all cases).

A naive Bayes table may be of more use when there are more even (and likely unorthodox) distributions of passengers and crew. For example, it would work better if the probability of a crew member being a child were equal to that of a crew member being an adult. In such a case, the Naive Bayes assumption that all given attributes have independent probabilities would be correct, and thus wouldn't harm the accuracy of its predictions.