

AI 3603 ARTIFICIAL INTELLIGENCE: PRINCIPLES AND TECHNIQUES

By: LiuQiaoan (520030910220)

HW#: 3

December 7, 2022

Contents

I. Introduction	2
II. Coding Part	2
A. <code>joinFactors(Factor1,Factor2)</code>	2
B. <code>marginalizeFactor(factorTable, hiddenVar)</code>	3
C. <code>marginalizeNetworkVariables(bayesNet, hiddenVar)</code>	3
D. <code>evidenceUpdateFactor(factorTable, evidenceVars, evidenceVals)</code>	4
E. <code>evidenceUpdateNet(bayesnet, evidenceVars, evidenceVals)</code>	4
F. <code>inference(bayesnet, hiddenVar, evidenceVars, evidenceVals)</code>	4
III. Written Part	5
A. sub-problem 1	5
B. sub-problem 2	5
(a)	7
(b)	9
C. sub-problem 3	11
D. sub-problem 4	12
(a)	13
(b)	15
conclusion of sub-problem 4	16
E. sub-problem 5	17
F. sub-problem 6	18
IV. Conclusion	19
References	19

I. INTRODUCTION

In this homework, I implement code for computing exact inferences in Bayesian networks of discrete random variables. After that, I use the implementation of Bayes Network to make query, which gives hidden facts behind the data.

II. CODING PART

In this section, I implement the code for computing exact inferences in Bayesian networks of discrete random variables.

A. `joinFactors(Factor1,Factor2)`

In this function, I use `pandas.merge` to join to factors together. There are two situation, both of which can be dealt with by `pandas.merge`.

- **Factor1** and **Factor2** have no common columns. For this situation, we need to extent the table and then crossly multiply the probability.
- **Factor1** and **Factor2** have common columns. For this situation, we need to concat the different columns together.

After merge two factors to one factor, we need to multiply their probability together to get new probability.

```

1  ## Join of two factors
2  ## Factor1, Factor2: two factor tables
3  ##
4  ## Should return a factor table that is the join of factor 1 and 2.
5  ## You can assume that the join of two factors is a valid operation.
6  ## Hint: You can look up pd.merge for merging two factors
7  def joinFactors(Factor1, Factor2):
8      # your code
9      intersect = set(Factor1).intersection(set(Factor2))
10     intersect.remove("probs")
11     if len(intersect) == 0:
12         Factor = pd.merge(Factor1, Factor2, how = "cross")
13     else:
14         Factor = pd.merge(Factor1, Factor2, on = list(intersect), how = "left")
15
16     Factor["probs.x"] = Factor["probs.x"] * Factor["probs.y"]
17     Factor.rename(columns={"probs.x": "probs"}, inplace=True)
18     Factor.drop(columns=["probs.y"], inplace = True)
19
20     return Factor

```

B. `marginalizeFactor(factorTable, hiddenVar)`

In this section, I use `pandas.groupby` to eliminate `hiddenVar`.

- If `hiddenVar` in `factorTable`'s columns, we group all the columns except `hiddenVar` and "probs". However, if there no columns left, we only return `None`.
- If `hiddenVar` not in `factorTable`'s columns, we just return `factorTable`.

```

1  ## Marginalize a variable from a factor
2  ## table: a factor table in dataframe
3  ## hiddenVar: a string of the hidden variable name to be marginalized
4  ##
5  ## Should return a factor table that marginalizes margVar out of it.
6  ## Assume that hiddenVar is on the left side of the conditional.
7  ## Hint: you can look can pd.groupby
8  def marginalizeFactor(factorTable, hiddenVar):
9      # your code
10     labels = list(factorTable)
11     if hiddenVar in labels:
12         labels.remove(hiddenVar)
13         labels.remove("probs")
14         if labels:
15             return factorTable.groupby(labels)["probs"].sum().to_frame().reset_index()
16         else:
17             return None
18     else:
19         return factorTable

```

C. `marginalizeNetworkVariables(bayesNet, hiddenVar)`

Since the parameters is a list of factors and a list of `hiddenVar`, so we need to call `marginalizeFactor` for each hidden variable. And because we generate `None` in `marginalizeFactor` when there are no other columns except hidden variable and "probs", we need to filter the `None` element out.

```

1  ## Marginalize a list of variables
2  ## bayesnet: a list of factor tables and each table in dataframe type
3  ## hiddenVar: a string of the variable name to be marginalized
4  ## NOTE: hiddenVar should be a list of string of the variable name to be marginalized according to PDF.
5  ##
6  ## Should return a Bayesian network containing a list of factor tables that results
7  ## when the list of variables in hiddenVar is marginalized out of bayesnet.
8  def marginalizeNetworkVariables(bayesNet, hiddenVar):
9      # your code
10     for h in hiddenVar:
11         bayesNet = [marginalizeFactor(factorTable, h) for factorTable in bayesNet]
12         bayesNet = [factor for factor in bayesNet if factor is not None]
13
14     return bayesNet

```

D. evidenceUpdateFactor(factorTable, evidenceVars, evidenceVals)

In order to make the process of `evidenceUpdateNet`, I define a new function. In this function, I first check whether `evidenceVars` is a string, an empty list or a non-empty list.

- If `evidenceVars` is a string, we use a boolean series to filter evidence variable with their known value.
- If `evidenceVars` is a non-empty list, we can use `pandas.query` method to set multiple conditions.
- If `evidenceVars` is an empty list, we just return the initial `factorTable`.

```

1 def evidenceUpdateFactor(factorTable, evidenceVars, evidenceVals):
2     # use 'query' method.
3     if isinstance(evidenceVars, list):
4         if evidenceVars and set(evidenceVars).issubset(set(factorTable)):
5             query_condition = [f"({var}) == {val}" for var, val in zip(evidenceVars, evidenceVals)]
6             return factorTable.query("&".join(query_condition)).reset_index().drop(columns=["index"])
7         else:
8             return factorTable
9     else:
10        if evidenceVars and evidenceVars in set(factorTable):
11            return factorTable[factorTable[evidenceVars] == int(evidenceVals)].reset_index().drop(columns=["index"])
12        else:
13            return factorTable

```

E. evidenceUpdateNet(bayesnet, evidenceVars, evidenceVals)

For this function, we just need to call `evidenceUpdateFactor`.

```

1 ## Update BayesNet for a set of evidence variables
2 ## bayesnet: a list of factor and factor tables in dataframe format
3 ## evidenceVars: a vector of variable names in the evidence list
4 ## evidenceVals: a vector of values for corresponding variables (in the same order)
5 ##
6 ## Set the values of the evidence variables. Other values for the variables
7 ## should be removed from the tables. You do not need to normalize the factors
8 def evidenceUpdateNet(bayesnet, evidenceVars, evidenceVals):
9     # your code
10    return [evidenceUpdateFactor(factorTable, evidenceVars, evidenceVals) for factorTable in bayesnet]

```

F. inference(bayesnet, hiddenVar, evidenceVars, evidenceVals)

For the last inference function, we can use `functools.reduce` to loop through the `bayesnet` to join each factor. Then we call `marginalizeNetworkVariables` to marginalize hidden variable, and call `evidenceUpdateNet` to update evidence. Last, we normalize the probability.

```

1 ## Run inference on a Bayesian network
2 ## bayesnet: a list of factor tables and each table in dataframe type
3 ## hiddenVar: a string of the variable name to be marginalized
4 ## evidenceVars: a vector of variable names in the evidence list
5 ## evidenceVals: a vector of values for corresponding variables (in the same order)
6 ##
7 ## This function should run variable elimination algorithm by using
8 ## join and marginalization of the sets of variables.
9 ## The order of the elimination can follow hiddenVar ordering
10 ## It should return a single joint probability table. The
11 ## variables that are hidden should not appear in the table. The variables
12 ## that are evidence variable should appear in the table, but only with the single
13 ## evidence value. The variables that are not marginalized or evidence should
14 ## appear in the table with all of their possible values. The probabilities
15 ## should be normalized to sum to one.
16 def inference(bayesnet, hiddenVar, evidenceVars, evidenceVals):
17     # your code
18     bayesnet = reduce(joinFactors, bayesnet)
19     bayesnet = marginalizeNetworkVariables([bayesnet], hiddenVar)
20     bayesnet = evidenceUpdateNet(bayesnet, evidenceVars, evidenceVals)[0]
21     bayesnet["probs"] = bayesnet["probs"] / bayesnet["probs"].sum()
22     return bayesnet

```

III. WRITTEN PART

A. sub-problem 1

The total number of probabilities needed to store the full joint distribution is

$$2^8 * 4^3 * 8 = 2^{17}$$

B. sub-problem 2

First, I write down the code for both problem 2. I define an auxiliary function Problem2Or4 to do inference in problem 2 and 4.

```

## Problem 2
2 print("===== Problem 2 BEGIN =====")
# used in plt.pie
4 def makeautopct(values):
    def myautopct(pct):
6         return '{p:.2f}%'.format(p=pct)
    return myautopct
8 riskFactorNet = pd.read_csv('RiskFactorsData.csv')
income = readFactorTablefromData(riskFactorNet, ['income'])
10 smoke = readFactorTablefromData(riskFactorNet, ['smoke', 'income'])
exercise = readFactorTablefromData(riskFactorNet, ['exercise', 'income'])
12 long_sit = readFactorTablefromData(riskFactorNet, ['long_sit', 'income'])
stay_up = readFactorTablefromData(riskFactorNet, ['stay-up', 'income'])
14 bmi = readFactorTablefromData(riskFactorNet, ['bmi', 'income', 'exercise', 'long_sit'])
diabetes = readFactorTablefromData(riskFactorNet, ['diabetes', 'bmi'])
16 bp = readFactorTablefromData(riskFactorNet, ['bp', 'exercise', 'long_sit', 'income', 'stay-up', 'smoke'])
cholesterol = readFactorTablefromData(riskFactorNet, ['cholesterol', 'exercise', 'stay-up', 'income', 'smoke'])
18 stroke = readFactorTablefromData(riskFactorNet, ['stroke', 'bmi', 'bp', 'cholesterol'])
attack = readFactorTablefromData(riskFactorNet, ['attack', 'bmi', 'bp', 'cholesterol'])
20 angina = readFactorTablefromData(riskFactorNet, ['angina', 'bmi', 'bp', 'cholesterol'])
def Problem2Or4(flag, riskFactorNet, income, smoke, exercise, long_sit, stay-up, bmi, diabetes, bp, cholesterol,
    stroke, attack, angina):
22     def good_or_bad(flag, obsVals, riskFactorNet, income, smoke, exercise, long_sit, stay-up, bmi, diabetes, bp,
        cholesterol, stroke, attack, angina):
        # 1. Calculate p(diabetes|smoke, exercise, long_sit, stay-up)
24         risk_net = [income, smoke, long_sit, stay-up, exercise, bmi, diabetes]
        factors = riskFactorNet.columns

26         margVars = list(set(factors) - {'diabetes', 'smoke', 'exercise', 'long_sit', 'stay-up'})
        obsVars = ['smoke', 'exercise', 'long_sit', 'stay-up']

28         p = inference(risk_net, margVars, obsVars, obsVals)
        print(f"p(diabetes|smoke={obsVals[0]}, exercise={obsVals[1]}, long_sit={obsVals[2]}, stay-up={obsVals[3]}):")
30         print(p)
        plt.figure()
        plt.pie(p["probs"], labels = p["diabetes"], autopct=makeautopct(p["probs"]))
32         plt.title(f"p(diabetes|smoke={obsVals[0]}, exercise={obsVals[1]}, long_sit={obsVals[2]}, stay-up={obsVals[3]})")
        plt.savefig(f"./output/Problem{flag}-p-of-diabetes-GIVEN-smoke-{{obsVals[0]}}-exercise-{{obsVals[1]}}-long_sit-{{obsVals[2]}}-stay-up-{{obsVals[3]}}.pdf")

34         # 2. Calculate p(stroke|smoke, exercise, long_sit, stay-up)
        risk_net = [income, smoke, long_sit, stay-up, exercise, bmi, bp, cholesterol, stroke]
        factors = riskFactorNet.columns

40         margVars = list(set(factors) - {'stroke', 'smoke', 'exercise', 'long_sit', 'stay-up'})
        obsVars = ['smoke', 'exercise', 'long_sit', 'stay-up']

42         p = inference(risk_net, margVars, obsVars, obsVals)
        print(f"p(stroke|smoke={obsVals[0]}, exercise={obsVals[1]}, long_sit={obsVals[2]}, stay-up={obsVals[3]}):")
44         print(p)
        plt.figure()
        plt.pie(p["probs"], labels = p["stroke"], autopct=makeautopct(p["probs"]))
46         plt.title(f"p(stroke|smoke={obsVals[0]}, exercise={obsVals[1]}, long_sit={obsVals[2]}, stay-up={obsVals[3]})")
        plt.savefig(f"./output/Problem{flag}-p-of-stroke-GIVEN-smoke-{{obsVals[0]}}-exercise-{{obsVals[1]}}-long_sit-{{obsVals[2]}}-stay-up-{{obsVals[3]}}.pdf")

48         # 3. Calculate p(attack|smoke, exercise, long_sit, stay-up)
        risk_net = [income, smoke, long_sit, stay-up, exercise, bmi, bp, cholesterol, attack]
        factors = riskFactorNet.columns

54         margVars = list(set(factors) - {'attack', 'smoke', 'exercise', 'long_sit', 'stay-up'})
        obsVars = ['smoke', 'exercise', 'long_sit', 'stay-up']

56         p = inference(risk_net, margVars, obsVars, obsVals)
        print(f"p(attack|smoke={obsVals[0]}, exercise={obsVals[1]}, long_sit={obsVals[2]}, stay-up={obsVals[3]}):")
60         print(p)
        plt.figure()
        plt.pie(p["probs"], labels = p["attack"], autopct=makeautopct(p["probs"]))
62
64

```

```

plt.title(f"p(attack|smoke={obsVals[0]}, exercise={obsVals[1]}, long.sit={obsVals[2]}, stay-up={obsVals
[3]})")
66 plt.savefig(f"./output/Problem{flag}-p-of-attack-GIVEN-smoke-{obsVals[0]}-exercise-{obsVals[1]}-long.sit-{
obsVals[2]}-stay-up-{obsVals[3]}.pdf")

68 ## 4. Calculate p(angina|smoke, exercise, long.sit, stay-up)
risk_net = [income, smoke, long.sit, stay-up, exercise, bmi, bp, cholesterol, angina]
70 factors = riskFactorNet.columns

72 margVars = list(set(factors) - {'angina', 'smoke', 'exercise', 'long.sit', 'stay-up'})
obsVars = ['smoke', 'exercise', 'long.sit', 'stay-up']

74 p = inference(risk_net, margVars, obsVars, obsVals)
76 print(f"p(angina|smoke={obsVals[0]}, exercise={obsVals[1]}, long.sit={obsVals[2]}, stay-up={obsVals[3]}):")
print(p)
78 plt.figure()
plt.pie(p["probs"], labels = p["angina"], autopct=make_autopct(p["probs"]))
80 plt.title(f"p(angina|smoke={obsVals[0]}, exercise={obsVals[1]}, long.sit={obsVals[2]}, stay-up={obsVals
[3]})")
plt.savefig(f"./output/Problem{flag}-p-of-angina-GIVEN-smoke-{obsVals[0]}-exercise-{obsVals[1]}-long.sit-{
obsVals[2]}-stay-up-{obsVals[3]}.pdf")

82 ## (a) Bad habits :
84 obsVals = [1, 2, 1, 1]
good_or_bad(flag, obsVals, riskFactorNet, income, smoke, exercise, long.sit, stay-up, bmi, diabetes, bp,
cholesterol, stroke, attack, angina)
86 ## (a) Good habits :
obsVals = [2, 1, 2, 2]
88 good_or_bad(flag, obsVals, riskFactorNet, income, smoke, exercise, long.sit, stay-up, bmi, diabetes, bp,
cholesterol, stroke, attack, angina)

90 def good_or_bad_b(flag, obsVals, riskFactorNet, income, smoke, exercise, long.sit, stay-up, bmi, diabetes, bp,
cholesterol, stroke, attack, angina):
92 ## 1. Calculate p(diabetes|bp, cholesterol, bmi)
risk_net = [income, smoke, long.sit, stay-up, exercise, bmi, bp, cholesterol, diabetes]
94 factors = riskFactorNet.columns

96 margVars = list(set(factors) - {'diabetes', 'bp', 'cholesterol', 'bmi'})
obsVars = ['bp', 'cholesterol', 'bmi']

98 p = inference(risk_net, margVars, obsVars, obsVals)
100 print(f"p(diabetes|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]}):")
print(p)
102 plt.figure()
plt.pie(p["probs"], labels = p["diabetes"], autopct=make_autopct(p["probs"]))
104 plt.title(f"p(diabetes|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]})")
plt.savefig(f"./output/Problem{flag}-p-of-diabetes-GIVEN-bp-{obsVals[0]}-cholesterol-{obsVals[1]}-bmi-{
obsVals[2]}.pdf")

106 ## 2. Calculate p(stroke|bp, cholesterol, bmi)
risk_net = [income, smoke, long.sit, stay-up, exercise, bmi, bp, cholesterol, stroke]
108 factors = riskFactorNet.columns

110 margVars = list(set(factors) - {'stroke', 'bp', 'cholesterol', 'bmi'})
obsVars = ['bp', 'cholesterol', 'bmi']

112 p = inference(risk_net, margVars, obsVars, obsVals)
114 print(f"p(stroke|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]}):")
print(p)
116 plt.figure()
plt.pie(p["probs"], labels = p["stroke"], autopct=make_autopct(p["probs"]))
118 plt.title(f"p(stroke|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]})")
plt.savefig(f"./output/Problem{flag}-p-of-stroke-GIVEN-bp-{obsVals[0]}-cholesterol-{obsVals[1]}-bmi-{
obsVals[2]}.pdf")

120 ## 3. Calculate p(attack|bp, cholesterol, bmi)
risk_net = [income, smoke, long.sit, stay-up, exercise, bmi, bp, cholesterol, attack]
122 factors = riskFactorNet.columns

124 margVars = list(set(factors) - {'attack', 'bp', 'cholesterol', 'bmi'})
obsVars = ['bp', 'cholesterol', 'bmi']

126 p = inference(risk_net, margVars, obsVars, obsVals)
128 print(f"p(attack|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]}):")
print(p)
130 plt.figure()
132 plt.pie(p["probs"], labels = p["attack"], autopct=make_autopct(p["probs"]))
plt.title(f"p(attack|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]})")
134 plt.savefig(f"./output/Problem{flag}-p-of-attack-GIVEN-bp-{obsVals[0]}-cholesterol-{obsVals[1]}-bmi-{
obsVals[2]}.pdf")

136 ## 4. Calculate p(angina|bp, cholesterol, bmi)
risk_net = [income, smoke, long.sit, stay-up, exercise, bmi, bp, cholesterol, angina]
138 factors = riskFactorNet.columns

140 margVars = list(set(factors) - {'angina', 'bp', 'cholesterol', 'bmi'})
obsVars = ['bp', 'cholesterol', 'bmi']

142 p = inference(risk_net, margVars, obsVars, obsVals)
144 print(f"p(angina|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]}):")
print(p)
146 plt.figure()
plt.pie(p["probs"], labels = p["angina"], autopct=make_autopct(p["probs"]))
148 plt.title(f"p(angina|bp={obsVals[0]}, cholesterol={obsVals[1]}, bmi={obsVals[2]})")
plt.savefig(f"./output/Problem{flag}-p-of-angina-GIVEN-bp-{obsVals[0]}-cholesterol-{obsVals[1]}-bmi-{
obsVals[2]}.pdf")

```

```

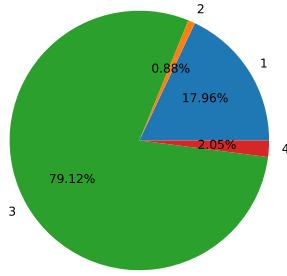
150  ## (b) Poor health:
152  obsVals = [1, 1, 3]
152  good_or_bad_b(flag, obsVals, riskFactorNet, income, smoke, exercise, long_sit, stay_up, bmi, diabetes, bp,
154  cholesterol, stroke, attack, angina)
154  ## (b) Good health:
156  obsVals = [3, 2, 2]
156  good_or_bad_b(flag, obsVals, riskFactorNet, income, smoke, exercise, long_sit, stay_up, bmi, diabetes, bp,
156  cholesterol, stroke, attack, angina)
158 Problem2Or4(2, riskFactorNet, income, smoke, exercise, long_sit, stay_up, bmi, diabetes, bp, cholesterol, stroke,
158 attack, angina)
160 print("===== Problem 2 END =====")
160 print()

```

(a)

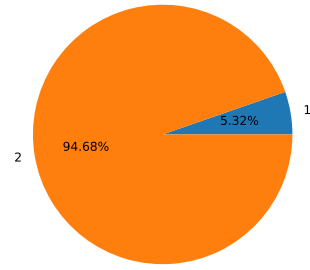
The probability of the outcome if I have bad habits (smoke, don't exercise, long sitting and stay up), can be calculated by inference, the pie charts of the probability are listed below.

p(diabetes|smoke=1, exercise=2, long_sit=1, stay_up=1)



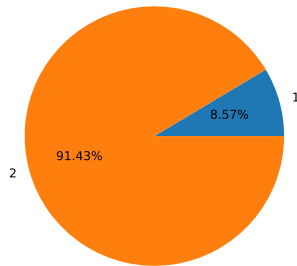
(a)

p(stroke|smoke=1, exercise=2, long_sit=1, stay_up=1)



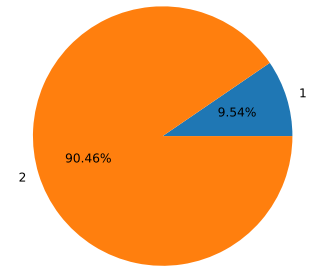
(b)

p(attack|smoke=1, exercise=2, long_sit=1, stay_up=1)



(c)

p(angina|smoke=1, exercise=2, long_sit=1, stay_up=1)

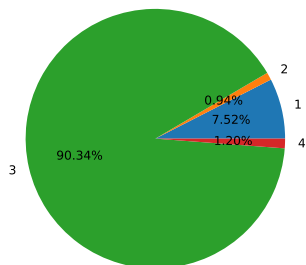


(d)

FIG. 1: The probability of the outcome if I have bad habits (smoke, don't exercise, long sitting and stay up)

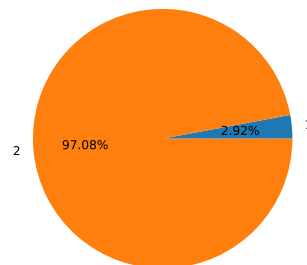
The probability of the outcome if I have good habits (don't smoke, exercise, no long sitting and don't stay up), can be calculated by inference, the pie charts of the probability are listed below.

$p(\text{diabetes}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



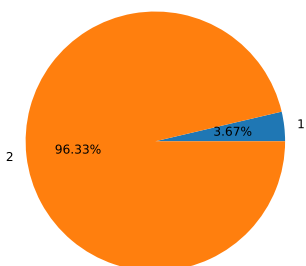
(a)

$p(\text{stroke}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



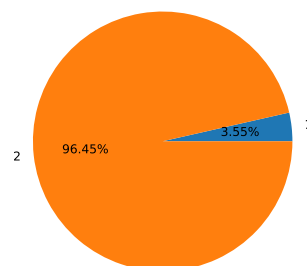
(b)

$p(\text{attack}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



(c)

$p(\text{angina}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



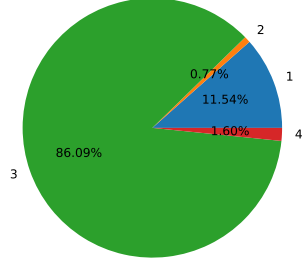
(d)

FIG. 2: The probability of the outcome if I have good habits (don't smoke, exercise, no long sitting and don't stay up)

(b)

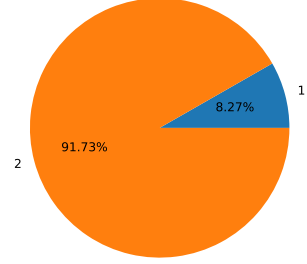
The probability of the outcome if I have poor health (high blood pressure, high cholesterol, and overweight), and probability of the outcome if I have good health, can be calculated by inference. The pie charts are listed below.

p(diabetes|bp=1, cholesterol=1, bmi=3)



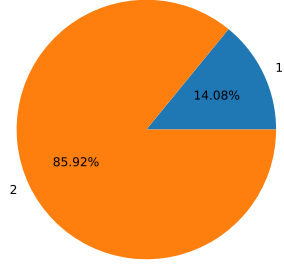
(a)

p(stroke|bp=1, cholesterol=1, bmi=3)



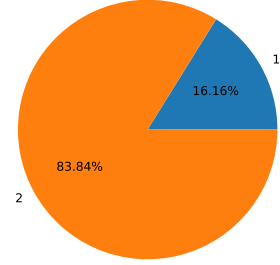
(b)

p(attack|bp=1, cholesterol=1, bmi=3)



(c)

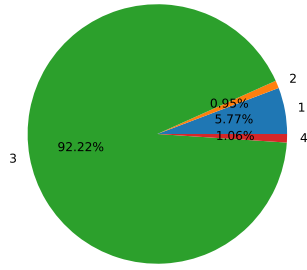
p(angina|bp=1, cholesterol=1, bmi=3)



(d)

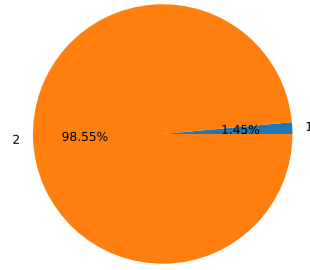
FIG. 3: The probability of the outcome if I have poor health (high blood pressure, high cholesterol, and overweight)

$p(\text{diabetes}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



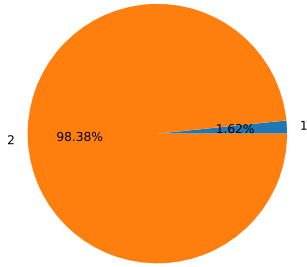
(a)

$p(\text{stroke}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



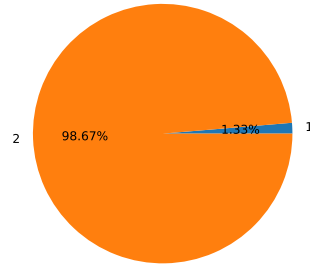
(b)

$p(\text{attack}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



(c)

$p(\text{angina}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



(d)

FIG. 4: The probability of the outcome if I have good health (low blood pressure, low cholesterol, and normal weight)

C. sub-problem 3

Now we are going to evaluate the effect a person's income has on their probability of having one of the four health outcomes (diabetes, stroke, heart attack, angina). I first write down the code.

```

1  ## Problem 3
2  print("===== Problem 3 BEGIN =====")
3  ## 1. Calculate  $p(\text{diabetes} = 1 \mid \text{incomes} = i)$ 
4  risk_net = [income, smoke, long_sit, stay_up, exercise, bmi, diabetes]
5  factors = riskFactorNet.columns
6
7  margVars = list(set(factors) - {'diabetes', 'income'})
8  obsVars = ['income']
9
10 p_diabetes = []
11 for obsVals in range(1,9):
12     p = inference(risk_net, margVars, obsVars, [obsVals])
13     p_diabetes.append(p["probs"][0])
14     print(p_diabetes)
15
16 ## 2. Calculate  $p(\text{stroke} = 1 \mid \text{incomes} = i)$ 
17 risk_net = [income, smoke, long_sit, stay_up, exercise, bmi, bp, stroke]
18 factors = riskFactorNet.columns
19
20 margVars = list(set(factors) - {'stroke', 'income'})
21 obsVars = ['income']
22
23 p_stroke = []
24 for obsVals in range(1,9):
25     p = inference(risk_net, margVars, obsVars, [obsVals])
26     p_stroke.append(p["probs"][0])
27     print(p_stroke)
28
29 ## 3. Calculate  $p(\text{attack} = 1 \mid \text{incomes} = i)$ 
30 risk_net = [income, smoke, long_sit, stay_up, exercise, bmi, bp, attack]
31 factors = riskFactorNet.columns
32
33 margVars = list(set(factors) - {'attack', 'income'})
34 obsVars = ['income']
35
36 p_attack = []
37 for obsVals in range(1,9):
38     p = inference(risk_net, margVars, obsVars, [obsVals])
39     p_attack.append(p["probs"][0])
40     print(p_attack)
41
42 ## 4. Calculate  $p(\text{angina} = 1 \mid \text{incomes} = i)$ 
43 risk_net = [income, smoke, long_sit, stay_up, exercise, bmi, bp, angina]
44 factors = riskFactorNet.columns
45
46 margVars = list(set(factors) - {'angina', 'income'})
47 obsVars = ['income']
48
49 p_angina = []
50 for obsVals in range(1,9):
51     p = inference(risk_net, margVars, obsVars, [obsVals])
52     p_angina.append(p["probs"][0])
53     print(p_angina)
54
55 ## plot
56 plt.figure()
57 plt.plot(list(range(1,9)), p_diabetes, label="diabetes")
58 plt.plot(list(range(1,9)), p_stroke, label="stroke")
59 plt.plot(list(range(1,9)), p_attack, label="attack")
60 plt.plot(list(range(1,9)), p_angina, label="angina")
61 plt.legend()
62 plt.xlabel(r"$i$")
63 plt.ylabel(r"$P(y=1 \mid \text{income}=i)$")
64 plt.title("The probability of outcome given income status")
65 plt.savefig("./output/Problem3.pdf")
66
67 print("===== Problem 3 END =====")
68 print()

```

And the probability given income status are showed below, the horizontal axis is $i = 1, 2, \dots, 8$, and the vertical axis is $P(y = 1 \mid \text{income} = i)$, where y is the outcome.

Then start to analysis. We can see from Figure 5 that the **more** of someone's outcome, the **less** probability they will have one of the four health outcomes. This can be a fact that the low-income person often work more busily, thus have less exercise, more long sitting, more staying up, more smoke, from problem 2 we know that these bad habits will more likely to cause high blood pressure, high cholesterol, and overweight. Ultimately, they will cause the four health outcomes.

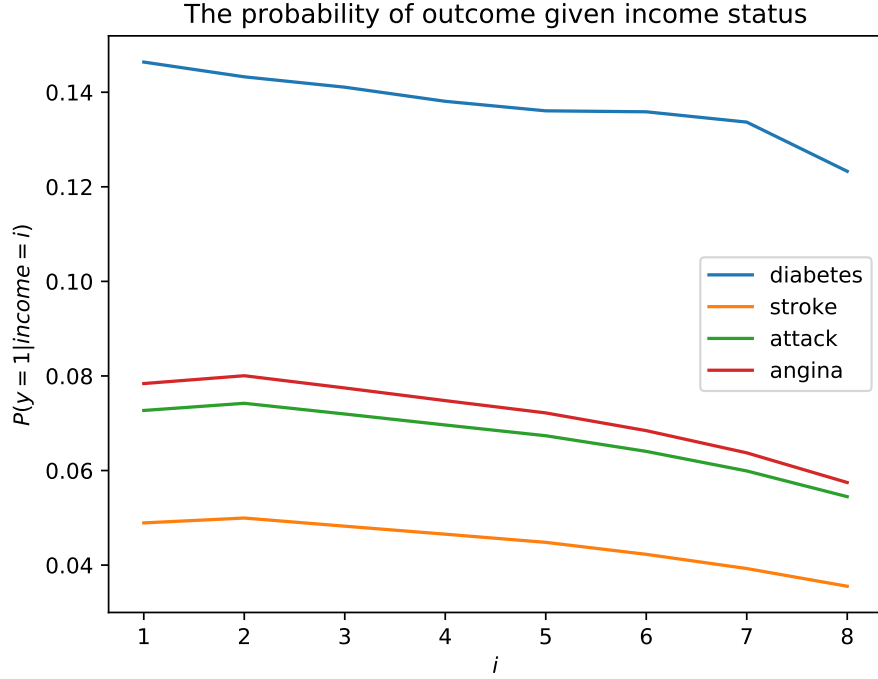


FIG. 5: The probability of outcome given income status, i.e. $P(y = 1|income = i)$

D. sub-problem 4

The **assumption** made in the first network is that *there are no direct effects between habits and health outcomes*. Now we are going to check whether the assumption is valid. In this problem, we can reuse the code in problem 2. I just need to reset the network.

```

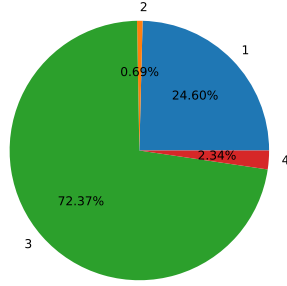
## Problem 4
2 print("===== Problem 4 BEGIN =====")
income = readFactorTablefromData(riskFactorNet, ['income'])
4 smoke = readFactorTablefromData(riskFactorNet, ['smoke', 'income'])
exercise = readFactorTablefromData(riskFactorNet, ['exercise', 'income'])
6 long_sit = readFactorTablefromData(riskFactorNet, ['long_sit', 'income'])
stay_up = readFactorTablefromData(riskFactorNet, ['stay-up', 'income'])
8 bmi = readFactorTablefromData(riskFactorNet, ['bmi', 'income', 'exercise', 'long_sit'])
diabetes = readFactorTablefromData(riskFactorNet, ['diabetes', 'bmi', 'smoke', 'exercise']) # add 'smoke' and '
exercise'
10 bp = readFactorTablefromData(riskFactorNet, ['bp', 'exercise', 'long_sit', 'income', 'stay-up', 'smoke'])
cholesterol = readFactorTablefromData(riskFactorNet, ['cholesterol', 'exercise', 'stay-up', 'income', 'smoke'])
12 stroke = readFactorTablefromData(riskFactorNet, ['stroke', 'bmi', 'bp', 'cholesterol', 'smoke', 'exercise']) #
add 'smoke' and 'exercise'
attack = readFactorTablefromData(riskFactorNet, ['attack', 'bmi', 'bp', 'cholesterol', 'smoke', 'exercise']) #
add 'smoke' and 'exercise'
14 angina = readFactorTablefromData(riskFactorNet, ['angina', 'bmi', 'bp', 'cholesterol', 'smoke', 'exercise']) #
add 'smoke' and 'exercise'
16 Problem2Or4(4, riskFactorNet, income, smoke, exercise, long_sit, stay-up, bmi, diabetes, bp, cholesterol, stroke,
attack, angina)
18 print("===== Problem 4 END =====")
print()

```

(a)

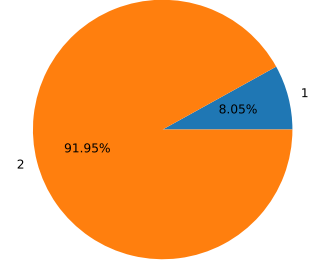
The probability of the outcome if I have bad habits (smoke, don't exercise, long sitting and stay up), can be calculated by inference, the pie charts of the probability are listed below.

$p(\text{diabetes}|\text{smoke}=1, \text{exercise}=2, \text{long_sit}=1, \text{stay_up}=1)$



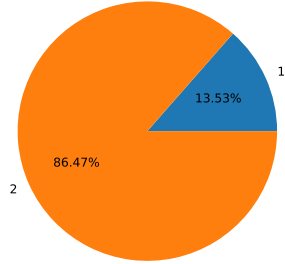
(a)

$p(\text{stroke}|\text{smoke}=1, \text{exercise}=2, \text{long_sit}=1, \text{stay_up}=1)$



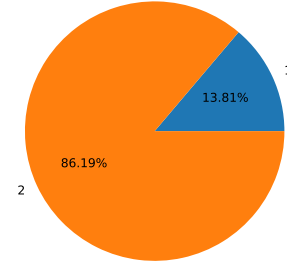
(b)

$p(\text{attack}|\text{smoke}=1, \text{exercise}=2, \text{long_sit}=1, \text{stay_up}=1)$



(c)

$p(\text{angina}|\text{smoke}=1, \text{exercise}=2, \text{long_sit}=1, \text{stay_up}=1)$

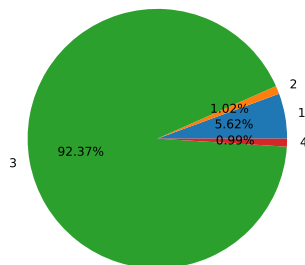


(d)

FIG. 6: The probability of the outcome if I have bad habits (smoke, don't exercise, long sitting and stay up)

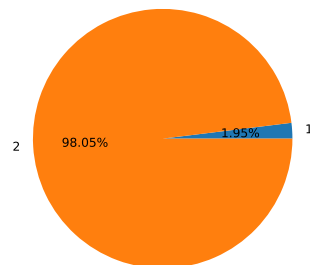
The probability of the outcome if I have good habits (don't smoke, exercise, no long sitting and don't stay up), can be calculated by inference, the pie charts of the probability are listed below.

$p(\text{diabetes}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



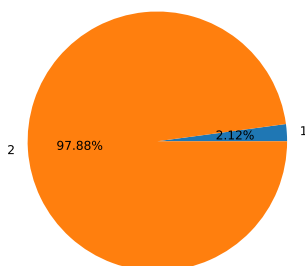
(a)

$p(\text{stroke}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



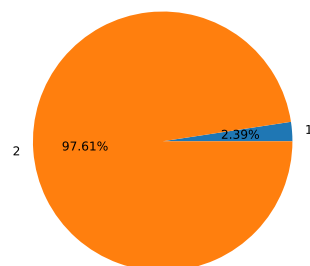
(b)

$p(\text{attack}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



(c)

$p(\text{angina}|\text{smoke}=2, \text{exercise}=1, \text{long_sit}=2, \text{stay_up}=2)$



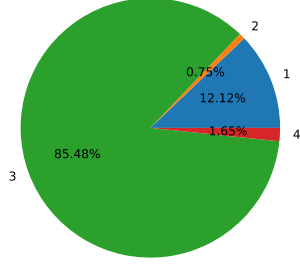
(d)

FIG. 7: The probability of the outcome if I have good habits (don't smoke, exercise, no long sitting and don't stay up)

(b)

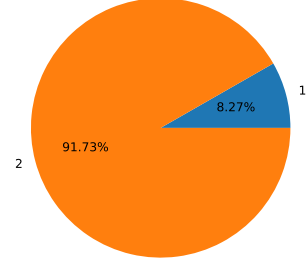
The probability of the outcome if I have poor health (high blood pressure, high cholesterol, and overweight), and probability of the outcome if I have good health, can be calculated by inference. The pie charts are listed below.

p(diabetes|bp=1, cholesterol=1, bmi=3)



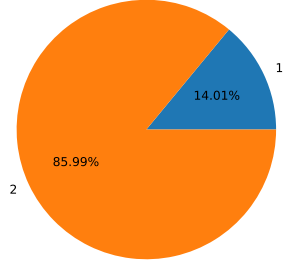
(a)

p(stroke|bp=1, cholesterol=1, bmi=3)



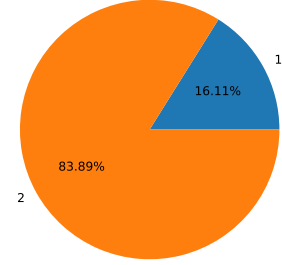
(b)

p(attack|bp=1, cholesterol=1, bmi=3)



(c)

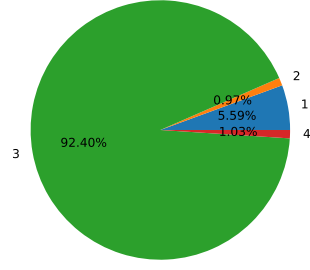
p(angina|bp=1, cholesterol=1, bmi=3)



(d)

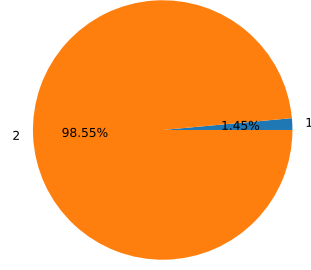
FIG. 8: The probability of the outcome if I have poor health (high blood pressure, high cholesterol, and overweight)

$p(\text{diabetes}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



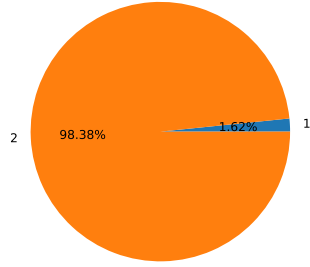
(a)

$p(\text{stroke}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



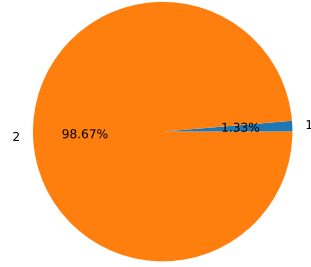
(b)

$p(\text{attack}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



(c)

$p(\text{angina}|\text{bp}=3, \text{cholesterol}=2, \text{bmi}=2)$



(d)

FIG. 9: The probability of the outcome if I have good health (low blood pressure, low cholesterol, and normal weight)

conclusion of sub-problem 4

Then, we can compare the result with problem 2. Compare Figure 1 with Figure 6, Figure 2 with Figure 7, Figure 4 with Figure 9, Figure 3 with Figure 8, we can find that they are all different. So the assumption is invalid.

E. sub-problem 5

When there are no edges between the four outcomes, the **assumption** made is that *there are no direct effects between the four outcomes*. Now we make a third network to test the assumption.

```

1  ## Problem 5
2  print("===== Problem 5 BEGIN =====")
3  ## 1. Calculate  $p(\text{stroke} = 1 \mid \text{diabetes} = 1)$  and  $p(\text{stroke} = 1 \mid \text{diabetes} = 3)$  with the network in problem 4.
4  risk_net = [income, smoke, long_sit, stay_up, exercise, bmi, diabetes, bp, cholesterol, stroke]
5  factors = riskFactorNet.columns
6
7  margVars = list(set(factors) - {'stroke', 'diabetes'})
8  obsVars = ['diabetes']
9
10 obsVals = [1]
11 p = inference(risk_net, margVars, obsVars, obsVals)
12 print("p(stroke = 1 | diabetes = 1) with the network in problem 4:")
13 print(p["probs"][0])
14
15 obsVals = [3]
16 p = inference(risk_net, margVars, obsVars, obsVals)
17 print("p(stroke = 1 | diabetes = 3) with the network in problem 4:")
18 print(p["probs"][0])
19
20 ## 2. Calculate  $p(\text{stroke} = 1 \mid \text{diabetes} = 1)$  and  $p(\text{stroke} = 1 \mid \text{diabetes} = 3)$  with third network.
21 stroke = readFactorTablefromData(riskFactorNet, ['stroke', 'bmi', 'bp', 'cholesterol', 'smoke', 'exercise', 'diabetes']) # add 'diabetes'
22
23 risk_net = [income, smoke, long_sit, stay_up, exercise, bmi, diabetes, bp, cholesterol, stroke]
24 factors = riskFactorNet.columns
25
26 margVars = list(set(factors) - {'stroke', 'diabetes'})
27 obsVars = ['diabetes']
28
29 obsVals = [1]
30 p = inference(risk_net, margVars, obsVars, obsVals)
31 print("p(stroke = 1 | diabetes = 1) with third network:")
32 print(p["probs"][0])
33
34 obsVals = [3]
35 p = inference(risk_net, margVars, obsVars, obsVals)
36 print("p(stroke = 1 | diabetes = 3) with third network:")
37 print(p["probs"][0])
38
39 print("===== Problem 5 END =====")

```

After run the code, we get the output as:

```

===== Problem 5 BEGIN =====
p(stroke = 1 | diabetes = 1) with the network in problem 4:
0.04441718130847403
p(stroke = 1 | diabetes = 3) with the network in problem 4:
0.03995473229116679
p(stroke = 1 | diabetes = 1) with third network:
0.07654239679086357
p(stroke = 1 | diabetes = 3) with third network:
0.034456134859939015
===== Problem 5 END =====

```

FIG. 10: $P(\text{stroke} = 1 | \text{diabetes} = 1)$ and $P(\text{stroke} = 1 | \text{diabetes} = 3)$ in network 2 and network 3

From the output, we can see that the probability of $P(\text{stroke} = 1 | \text{diabetes} = 1)$ and $P(\text{stroke} = 1 | \text{diabetes} = 3)$ are both different for network 2 and network 3, so the assumption is invalid.

F. sub-problem 6

After run the initial **BayesNetworkTestScript.py**, I get the output as below, which is the same as that in the PDF document [1].

```
inference starts
  gauge probs
0      0 0.315
1      1 0.685
  fuel gauge probs
0      0      0 0.81
1      0      1 0.19
  fuel gauge probs
0      0      0 0.257143
1      1      0 0.742857
      probs battery fuel gauge
0 0.888889      0      1      0
1 0.111111      0      0      0
inference ends
income dataframe is
      probs income
0 0.050848      1
1 0.059429      2
2 0.074042      3
3 0.094414      4
4 0.116356      5
5 0.150725      6
6 0.164430      7
7 0.289755      8
      smoke long_sit exercise diabetes probs
0      1      1      2      1 0.136815
1      1      1      2      2 0.008916
2      1      1      2      3 0.837218
3      1      1      2      4 0.017052
```

FIG. 11: My code runs correctly on all of the examples in **BayesNetworkTestScript.py**.

IV. CONCLUSION

In this homework, I combine the knowledge in the class with a real-world example, and use Bayes Network to analyze risk factors for certain health problems. During the procedure, I get familiar with `pandas` package, and enhance my ability of data handling.

-
- [1] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.