Last Name:

First Name:

Matriculation Number:

Exam Artificial Intelligence 2

August 2, 2022

	To be used for grading, do not write here											
prob.	1.1	1.2	2.1	2.2	3.1	3.2	4.1	4.2	4.3	5.1	Sum	grade
total	7	10	8	14	15	10	10	6	6	9	95	
reached												

In the course Artificial Intelligence I/II we award bonus points for the first student who reports a factual error in an old exam. (Please report spelling/formatting errors as well.)

1 Probabilities

Problem 1.1 (Python)

Consider the Python program below.

1. Which mathematical function does the method foo compute?

Solution: The matrix product $a \cdot b$.

- 2. Assume random variables *X* with domain $\{0, ..., m 1\}$ and *Y* with domain $\{0, ..., n 1\}$. Assume the following Python objects
 - *C* holds the conditional probability distribution P(X | Y), i.e., C[i][j] = P(X = i | Y = j).
 - *D* holds the probability distribution P(Y), i.e., D[j] = P(Y = j).

Complete the definition of *E* in the program below in such a way that it holds the probability distribution P(X), i.e., E[i] = P(X = i). (Hint: This requires relatively little code.)

```
def foo(a,b):
    1 = len(a)
    m = len(a[0])
    n = len(b[0])
    res = []
    for i in range(1):
        row = []
        for j in range(n):
            s = 0
            for k in range(m):
                s += a[i][k] * b[k][j]
               row.append(s)
        res.append(row)
    return res
E =
```

Solution: E=foo(C, [[x] for x in D]). We also accepted E=foo(C,D).

Problem 1.2 (Calculations)

Assume three Boolean random variables X, Y, Z, whose joint distribution P(X, Y, Z) is given by

x	У	Z	P(X = x, Y = y, Z = z)
true	true	true	a
true	true	false	b
true	false	true	С
true	false	false	d
false	true	true	е
false	true	false	$\int f$
false	false	true	g
false	false	false	h

1. In terms of a, b, c, d, e, f, g, h, give P(X = true, Y = false).

Solution: P(X = true, Y = false) = P(X = true, Y = false, Z = true) + P(X = true, Y = false, Z = false) = c + d

2. In terms of a, b, c, d, e, f, g, h, give P(X = true | Y = false).

Solution: P(X = true | Y = false) = P(X = true, Y = false)/P(Y = false) = (c + d)/(c + d + g + h).

3. Which of the following are true if *Y* and *Z* are conditionally independent given *X*?

 $\Box a + b + e + f = a + c + e + g$ $\swarrow a = (a + b) \cdot (a + c)/(a + b + c + d)$ $\bowtie e = (f \cdot g)/h$ $\Box a = e$

2 Bayesian Reasoning

When working with an upper case Boolean random variable X, you may abbreviate the event X = true by the corresponding lower-case letter x. If you do that, make sure the distinction between upper and lower case letters is clear in your writing.

Problem 2.1 (Bayes' Rule)

Assume you are trying to predict whether a particular topic comes up in an exam. You have collected the following data:

8 pt

- 30% of all topics come up in the exam.
- 50% of all topics come up in the assignments.
- If a topic comes up in an exam, it was covered by an assignment 60% of the time.
- If a topic comes up in an exam, it also came up in a recent exam 80% of the time.

You model this situation using 3 Boolean random variables *E* (comes up in exam), *A* (covered by assignments), and *R* (came up in a recent exam).

- 1. By filling in the gaps below, state for each number in the text above, which probability it describes.
 - 1. P(**E=true**) = 0.32. P(**A=true**) = 0.53. P(**A=true** E = true |=)0.64. P(**R=true** E = true |=)0.8
- 2. When modeling this situation, is it reasonable to assume that *A* and *R* are stochastically independent? Why (not)?

Solution: No. Background knowledge indicates that *A* and *R* are often highly correlated (even if the details cannot be ascertained from the data given).

3. The topic you are interested in was covered by an assignment. Using Bayes' rule, calculate the probability that it will come up in the exam.

Solution: P(E = true | A = true) = P(A = true | E = true) * P(E = true)/P(A = true) = 0.6 * 0.3/0.5 = 0.36

Problem 2.2 (Bayesian Networks)

Consider the following situation:

- Covid can cause a sickness and/or fever.
- Fever itself is dangerous and can cause sickness.
- Tests can detect Covid. But a false-positive Covid test may cause sickness via a kind of Placebo effect.
- There are no other causal relationships.

You want to model this situation using Boolean random variables *C* (Covid infection), *F* (fever), *S* (sickness), and *T* (positive Covid test).

1. Give a good variable ordering for forming a Bayesian network for this situation.

Solution: Causes should occur before effects, so e.g., CFTS or CTFS.

2. Give the resulting network.

Solution: $C \to F \to S$ and $C \to T \to S$ and $C \to S$.

3. You have a fever, feel sick, and have tested positive for Covid. Now you want to determine if you have Covid. What are the query, evidence, and hidden variables?

Solution: Query: C, evidence: F, S, T, hidden: none.

4. Assume your network is *C* → *F* → *T* ← *S* (which *may or may not* be a correct solution to the above question). Which probabilities are stored in the conditional probability table of node *T*?

Solution: The probability distribution P(T | F, S), i.e., P(T = x | F = y, S = z) as a function of Booleans x, y, z.

5. Again using the network $C \rightarrow F \rightarrow T \leftarrow S$, give the formula for

$$P(C = true, T = true, S = true)$$

in terms of the entries of the conditional probability table of that network.

Solution:

 $\begin{aligned} P(c,t,s) &= P(c,t,s,f) + P(c,t,s,\neg f) \\ &= P(c) \cdot P(f \mid c) \cdot P(s \mid f, c) \cdot P(t \mid c, f, s) + P(c) \cdot P(\neg f \mid c) \cdot P(s \mid \neg f, c) \cdot P(t \mid c, \neg f, s) \\ &= P(c) \cdot P(f \mid c) \cdot P(s) \cdot P(t \mid f, s) + P(c) \cdot P(\neg f \mid c) \cdot P(s) \cdot P(t \mid \neg f, s) \\ &= P(c) \cdot P(s) \cdot \left(P(f \mid c) \cdot P(t \mid f, s) + P(\neg f \mid c) \cdot P(t \mid \neg f, s) \right) \end{aligned}$

Markovian Reasoning 3

Problem 3.1 (Hidden Markov Models)

Consider the following situation:

- You make daily observations about your business B. Each day business is either good (b_1) or bad (b_2) .
- You know this is caused by the weather W, which can be rainy (w_1) , cloudy (w_2) , or sunny (w_3) .
- You have previously obtained the following information:
 - when the weather is rainy, your business is good 36% of the time,
 - when the weather is cloudy, your business is good 84% of the time,
 - when the weather is sunny, your business is good 90% of the time,
 - half the time, the weather is the same as on the previous day,
 - when the weather changes from one day to the next, each change is equally likely.

You want to model this situation as a hidden Markov model with two families of random variables indexed by day number d.

1. Give the state and evidence variables and their domains.

Solution: State variables $W_d \in \{w_1, w_2, w_3\}$, evidence variables $B_d \in \{b_1, b_2\}$

2. How can you tell that the sensor model is stationary here?

Solution: The business-weather relation is the same for each day.

- 3. What order does the model have?
- 4. Complete the following sentences:
 - 5. The transition model *T* is given by the matrix

$T = \left(\begin{array}{c} \end{array} \right)$		where	$T_{ij} = P(W_{d+1} = w_j \mid W_d = w_i$
Solution: $T = \begin{pmatrix} 0.5 \\ 0.25 \\ 0.25 \end{pmatrix}$	$\begin{array}{ccc} 0.25 & 0.25 \\ 0.5 & 0.25 \\ 0.25 & 0.5 \end{array}$		

6. The sensor model *S* is given by the matrix

、

$$S = \begin{pmatrix} \\ \end{pmatrix}$$
 where $S_{ij} = P(B_d = b_j | W_d = w_i).$

		(0.36	0.64
Solution:	S =	0.84	0.16
		0.9	0.1

7. Let *T* be as above and let \mathbf{v} be a 3-dimensional vector whose coefficients sum to 1 What is the intuitive meaning of the property $T \cdot \mathbf{v} = \mathbf{v}$?

Solution: \mathbf{v} is a probability distribution of the weather that is a fixed point of the transition model, i.e., the distribution will stay the same when predicting the future.

- 8. It was sunny yesterday, and your business is good today. You want to use filtering to obtain the probability distribution of today's weather. You proceed as follows:
 - 9. Give the recursive filtering equation for $f_{1:d+1}$.

Solution: We compute $f_{1:1}$ by applying the filtering equation once. $f_{1:d+1} = \alpha(O_{d+1} \cdot T^t f_{1:d})$

10. Give the initial value $f_{1:0}$ to use in this case.

Solution: $f_{1:0} = \langle 0, 0, 1 \rangle$

11. Give the diagonal sensor matrix O_1 to use in this case.

		(0.36	0	0)	
Solution:	$O_1 =$	0	0.84	0	
		0	0	0.9)	

12. Compute the resulting distribution.

Fully compute all values including the normalization. (This does not require approximations or a calculator.)

Solution: $f_{1:1} = \alpha \cdot O_1 \cdot T^t \cdot f_{1:0} = \alpha \langle 0.36 \cdot 0.25, 0.84 \cdot 0.25, 0.9 \cdot 0.5 \rangle$ = 4/3(0.09, 0.21, 0.45) = (0.12, 0.28, 0.6).

Problem 3.2 (Utility and Decision Processes)

Consider an agent moving on an 8×8 grid as indicated in the picture below. The agent can move up, down, right, and left except where restricted by the edges of the grid. Every action results in moving one step in that direction with probability 75% and no move otherwise. The agent's goal is to get to the location (7, 7).



1. Choose an appropriate reward function and model this situation as a Markov Decision Process.

Solution: One possible model is

- $S = \{0, \dots, 7\}^2$
- $A((i, j)) = \{u, d, l, r\} \setminus E_i \setminus F_j$ where $E_0 = \{l\}, E_7 = \{r\}, F_0 = \{d\}, F_7 = \{u\}$ and $E_i = F_j = \emptyset$ otherwise
- $P(s' \mid a, s)$ is 0.75 if s' is the result of moving a from s, 0.25 if s' = s, 0 otherwise
- A typical choice is any function *R* that is high for (7, 7) and slightly negative for other states. E.g., R(s) = 1 for s = (7, 7) and R(s) = -0.1 otherwise.
- 2. Give an optimal policy π^* .

Solution: Any policy that maps state (7,7) to *d* or *l* and every other state to any legal action that is *u* or *r*. E.g., $\pi^*(s) = u$ if $u \in A(s)$, otherwise $\pi^*(s) = r$ if $r \in A(s)$, otherwise $\pi^*((7,7)) = d$.

3. Now ignore the rewards, and assume we use a fixed utility i + j for the field (i, j). Compute the expected utility of moving up once in state (1, 1).

Solution:
$$EU(u) = 0.75 \cdot U((1, 2)) + 0.25 \cdot U((1, 1)) = 2.25 + 0.5 = 2.75$$

4. Now assume the agent is unable to tell whether an action resulted in a move or not. Explain informally how that would change the modeling.

Solution: We would need a POMDP. A state in the POMDP is a so-called belief state, a probability distribution for the MDP-state that the agent is in.

4 Learning

Problem 4.1 (Decision Trees and Lists)

10 pt

You observe the set of values below for 6 games of a sports team. You want to predict the result based on weather, location, and opponent.

#	Weather	Location	Opponent	Result
1	Rainy	Home	Weak	Win
2	Sunny	Away	Weak	Win
3	Sunny	Home	Strong	Loss
4	Sunny	Away	Weak	Win
5	Cloudy	Away	Strong	Loss
6	Sunny	Home	Strong	Loss

1. Draw the decision tree that arises if attributes are chosen according to the priority *Location*, *Weather*, *Opponent*.



2. Give all minimal sets *A* of attributes such that A > Result.

Solution: {Weather, Location} and {Opponent}

3. How can such a minimal set *A* be exploited to find a decision tree?

Solution: Only the attributes in *A* are needed in the tree. So smaller sets yield smaller trees.

4. You want to build a decision list for the result using tests with literals of the form *attribute* = *number*. Give the shortest possible decision list.

Solution: If *Opponent* = *Weak* then *Result* = *Win* else *Result* = *Loss*.

5. Now stop using the above observations. Instead, assume some set of observations for which a decision list exists, and assume that the determination *Weather*, *Location* ≻ *Result* holds. In the worst case, what is the number of tests in the shortest decision list?

Solution: 5 (not longer because we can use the determination and use one test for every one of the 6 combinations of values (no test is needed for the last one because we can use the final else-case); not shorter because we might indeed need those 6 cases)

Problem 4.2 (Statistical Learning)

6 pt

You observe the values below for 20 games of a sports team. You want to predict the result based on weather and opponent.

		Number of		
Weather	Opponent	wins	losses	
Rainy	Weak	3	1	
Cloudy	Weak	0	1	
Sunny	Weak	4	2	
Rainy	Strong	0	2	
Cloudy	Strong	2	3	
Sunny	Strong	0	2	

1. What is the hypothesis space for this situation, seen as an *inductive learning problem*?

Solution: The set of functions {*Rainy*, *Cloudy*, *Sunny*} \times {*Weak*, *Strong*} \rightarrow {*Win*, *Loss*}.

2. Explain whether we can learn the function by building a decision tree.

Solution: It does not. Even all attributes together, i.e., *Weather* and *Opponent*, do not determine the result. So no decision tree exists.

- 3. To apply Bayesian learning, we model this situation as a Bayesian network $W \rightarrow R \leftarrow O$ using random variables *W* (weather), *O* (opponent), and *R* (game result). What are the resulting entries of the conditional probability table for the cases
 - 1. P(W = rainy) = 3/10
 - 2. P(R = win | O = weak) = 7/11

Solution: P(W = rainy) = 3/10 and P(R = win | O = weak) = 7/11

Problem 4.3 (Support Vector Machines)

Consider a set of points $\mathbf{x} = \langle \mathbf{x}_1, \mathbf{x}_2 \rangle$ in \mathbb{R}^2 that are classified as either y = +1 or y = -1.

1. Give the hypothesis space for finding a linear separator.

Solution: The set of functions $\mathbf{w} \cdot \mathbf{x} + b$ for real numbers \mathbf{w}_1 , \mathbf{w}_2 , b. Alternatively, one can use \mathbb{R}^3 with some explanation that it holds the tuples (\mathbf{w}_1 , \mathbf{w}_2 , b).

6 pt

2. Given a linear separator *h*(**x**), which formula computes the classification *y* of a vector **x**?

Solution: $y = sgn(h(\mathbf{x}))$

3. What is the point of transforming a dataset into a higher-dimensional space?

Solution: It's possible that no linear separator exists for the original dataset, but a linear separator exists for the transformed dataset in the bigger space.

4. In this context, briefly discuss the usefulness of the transformation $F(\langle \mathbf{x}_1, \mathbf{x}_2 \rangle) = \langle \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_1 + \mathbf{x}_2 \rangle$.

Solution: It is not useful because it is linear. If a linear separator exists afterwards, one existed before as well.

5 Natural Language Processing

Problem 5.1 (Information Retrieval)

Consider the following three texts

- d_1 : "The man is tall."
- *d*₂: "The tall man sees the woman."
- *d*₃: "The woman shouts."

Let $D = \{d_1, d_2, d_3\}.$

Below we use alphabetical order for the vector components:

is, man, sees, shouts, tall, the, woman

Simplify all results as much as possible without introducing approximate values.

1. What is the idea of cosine similarity for comparing a query against the documents in *D*?

Solution: The query and each document are represented as a vector representing word frequencies. Vectors pointing in the same directions are considered similar. So the documents can be ranked by the angle between them and the query.

2. Give the vector $tf(_, d_2)$

Solution: $tf(_, d_2) = \langle 0, 1/6, 1/6, 0, 1/6, 1/3, 1/6 \rangle$.

3. Give the vector $idf(_, D)$.

Solution: $idf(_,D) = \log_{10}(3/\langle 1,2,1,1,2,3,2\rangle) = \langle k,l,k,k,l,0,l\rangle$ with $k = \log_{10} 3$ and $l = \log_{10} 1.5$.

4. For $d \in D$ and a word *t*, give the definition of tfidf(t, d, D).

Solution: $tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$

5. What is the benefit of using tfidf instead of tf for using cosine similarity?

Solution: tfidf gives more weight to words that occur in fewer documents. Otherwise, many documents would falsely appear similar just because the most common words appear in most of them.