

Assignment 7: Written Exercises II

by the Staff of CSE 415

Due June 2, 2021 via GradeScope

This is an individual-work assignment.

Prepare your answers in a neat, easy-to-read PDF. Our grading rubric will be set up such that when a question is not easily readable or not correctly tagged or with pages repeated or out of order, then points will be deducted. However, if all answers are clearly presented, in proper order, and tagged correctly when submitted to Gradescope, we will award a 5-point bonus.

If you choose to typeset your answers in Latex using the template file for this document, please put your answers in blue while leaving the original text black.

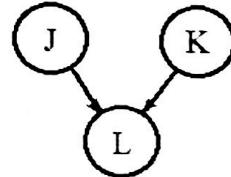
1 Joint Distributions and Factoring

(35 points) Consider the joint probability distribution below.

J	K	L	$P(J, K, L)$
true	true	true	0.024
true	true	false	0.096
true	false	true	0.028
true	false	false	0.252
false	true	true	0.144
false	true	false	0.036
false	false	true	0.294
false	false	false	0.126

- (a) (15 points) For each of the three pairs of random variables (J, K) , (J, L) , and (K, L) , provide a computation to prove that either (i) the two variables are independent, or (ii) they are dependent. If you provide any marginal or other derived distributions, make sure they are clearly identified, e.g., “ $P(K|L)$ ”, etc.

- (b) (15 points) Suppose that somebody has distributions for J and K , and a CPT (conditional probability table) for $P(L|J, K)$. Suggested that the joint distribution can be factored in a way that corresponds to the graph at the right. Prove either that (i) it cannot be factored according to this structure (i.e., some conditional independence assumption would be violated), or (ii) it can be factored according to this structure (i.e., provide the factorization with marginal dis-



- (c) (3 points) Whether or not our joint distribution can actually be factored according to this two-parent, one-child structure, explain how many “free parameters” would be involved in the specification.
- (d) (2 points) How many free parameters are (or would be) saved by using the factored representation of the joint distribution vs using free parameters of the full joint distribution table at the beginning of this problem?

6

1

1. (a) To check independence between two variables, say X and Y , we know that X and Y are independent if $P(X, Y) = P(X)P(Y)$. To find $P(J)$, $P(K)$, and $P(L)$ we must add rows of the joint probability distribution table.

For bookkeeping, let capital letters denote true (e.g. $J \Rightarrow J = \text{true}$) and lowercase denote false (e.g. $j \Rightarrow J = \text{false}$). So,

$$\begin{aligned} P(J) &= P(J, K, L) + P(J, K, l) + P(j, K, L) + P(j, K, l) \\ &= 0.024 + 0.096 + 0.028 + 0.292 \\ &= 0.4 \end{aligned}$$

$$\begin{aligned} P(K) &= P(J, K, L) + P(J, K, l) + P(j, K, L) + P(j, K, l) \\ &= 0.024 + 0.096 + 0.144 + 0.036 \\ &= 0.3 \end{aligned}$$

$$\begin{aligned} P(L) &= P(J, K, L) + P(J, k, L) + P(j, K, L) + P(j, k, L) \\ &= 0.024 + 0.028 + 0.144 + 0.294 \\ &= 0.49 \end{aligned}$$

$$P(J, K) = P(J, K|L) + P(J, K|l) = 0.12$$

$$P(J, L) = P(J, L|K) + P(J, L|k) = 0.052$$

$$P(K, L) = P(K, L|J) + P(K, L|j) = 0.168$$

$$P(J, K) = P(J)P(K) = 0.12 = (0.4)(0.3) \Rightarrow \underline{J, K \text{ independent}}$$

$$P(J, L) \neq P(J)P(L) \text{ since } 0.052 \neq (0.4)(0.49) \Rightarrow \underline{J, L \text{ dependent}}$$

$$P(K, L) \neq P(K)P(L) \text{ since } 0.168 \neq (0.3)(0.49) \Rightarrow \underline{K, L \text{ dependent}}$$

Note: we must also consider false cases to definitively prove that J and K are independent. That is,

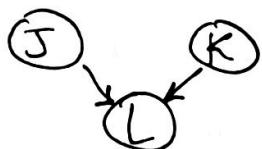
$$P(J, k) = P(J)P(k) = 0.28 = (0.4)(0.7) = 0.028 + 0.252$$

$$P(j, K) = P(j)P(K) = 0.18 = (0.6)(0.3) = 0.144 + 0.036$$

$$P(j, k) = P(j)P(k) = 0.42 = (0.6)(0.7) = 0.294 + 0.126$$

Since these are all true, we have confirmed that J and K are independent. Similar calculations apply but are unnecessary to show that J, L and K, L are dependent.

1. (b)



In (a), we determined that J and K are independent and that L is dependent on both J and K . The factored conditional probability distribution is therefore allowable.

We have the following CPT:

L	J	K	$P(L J,K)$
True	True	True	.024/.12 = 0.2
False	True	True	.096/.12 = 0.8
True	True	False	.028/.28 = 0.1
False	True	False	.252/.28 = 0.9
True	False	True	.144/.18 = 0.8
False	False	True	.036/.18 = 0.2
True	False	False	.294/.42 = 0.7
False	False	False	.126/.42 = 0.3

(c) From the CPT, we have 4 free parameters because each permutation of J, K has one free parameter associated with it. Since the CPT is factored from J and K , we must also consider 1 free parameter per independent J and K . Hence,

$$\# \text{ free parameters} = 4 + 1 + 1 = 6.$$

(d) The original joint probability distribution table has 7 free parameters because $P(J, K, L)$ for all 8 rows must sum to 1. Therefore, we save

$$7 - 6 = 1 \text{ free parameter}$$

by factoring.

2 Bayes Nets: D-Separation

(35 points) Consider the Bayes Net graph at the bottom of the page, which represents the topology of a web-server security model. Here the random variables have the following interpretations:

V = Vulnerability exists in web-server code or configs.

C = Complexity to access the server is high. (Passwords, 2-factor auth., etc.)

S = Server accessibility is high. (Firewall settings, and configs are permissive).

A = Attacker is active.

L = Logging infrastructure is state-of-the-art.

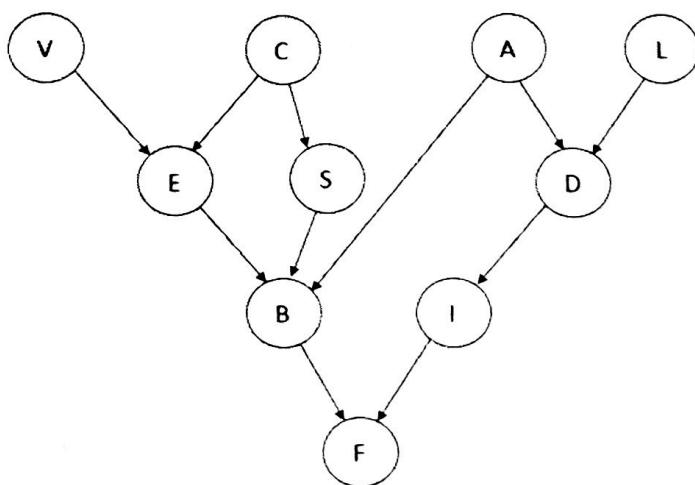
E = Exposure to vulnerability is high.

D = Detection of intrusion attempt.

B = Break-in; the web server is compromised.

I = Incident response is effective.

F = Financial losses are high (due to data loss, customer dissatisfaction, etc.).



For each of the following statements (4 points each), indicate whether (True) or not (False) the topology of the net guarantees that the statement is true. If False, identify an undirected path through which influence propagates between the two random variables being considered. (Be sure that the path follows the D-Separation rules covered in lecture.)

For example: $E \perp\!\!\!\perp S$: False (ECS).

- (a) $L \perp\!\!\!\perp I | F$ False (LDI)
- (b) $I \perp\!\!\!\perp E | A, S$ True
- (c) $L \perp\!\!\!\perp C | S, B, F$ False (LDIFBEC)
- (d) $F \perp\!\!\!\perp C | B, L, E$ False (FBSC)
- (e) $L \perp\!\!\!\perp V | D, E, F, S$ True

(g) (15 points) Suppose that the company hired an outside expert to examine the system and she determines that B and E are true: The web server is compromised, and exposure to vulnerability is high. Given this information, your job is to explain to management why getting additional information about A (whether the attacker is active) could have an impact on the probability of V (regarding the existence or non-existence of vulnerabilities). Give your explanation, for the manager of the company, using about between 10 and 20 lines of text, which should be based on what you know about D-separation, applied to this situation. However, your explanation should not use the terminology of D-separation but be **in plain English**. (You can certainly use words like “influence”, “probability”, “given”, but not “active path”, “triple”, or even “conditionally independent”).

In addition to E, both S and A can contribute to the appearance of B. That is, server accessibility is high and attacker is active can lead to a break-in, which is already determined to be true. An attacker could also trigger the detection of an intruder attempt (D). Detection leads to an incident response (I), which in turn lead to financial losses (F). Given that there is a break-in, financial losses are guaranteed. Whether or not we gain information about A is therefore important for two primary reasons. Firstly, an attacker is the only path to detection and incident response, which would inform us of further financial losses. Secondly, if we detect an attacker given a break-in, the probability that server accessibility is high (S) decreases since E, S, and A contribute to B. If the probability of S decreases, then the chances that complexity to access the server is high (C) also decrease.
⁴

Given E is true, a decrease in the probability of C raises the likelihood that vulnerability exists in the web server code (V) is a root cause. In short, information about an attacker greatly informs both causes and effects of a server break-in.

3 Markov Models

(30 points) According to an unnamed source, the stock market can be modeled using a Markov model, where there are two states “bull” and “bear.” The dynamics of the model are given in the table below:

S_{t-1}	S_t	$P(S_t S_{t-1})$
bull	bull	0.5
bull	bear	0.5
bear	bull	0.3
bear	bear	0.7

- (a) (4 points) Draw a visual state-transition diagram to represent the conditional probability table of this Markov model. (Hint: there should be two nodes in your diagram).
- (b) (6 points) Compute the stationary probabilities for *bull* and *bear*.

	<i>bull</i>	<i>bear</i>
P_∞	0.375	0.625

- (c) (15 points) Suppose it's given that $S_0 = \text{bull}$. Perform 5 rounds of the mini-forward algorithm and calculate $P(S_1), P(S_2), P(S_3), P(S_4), P(S_5)$ for each outcome of *bull* and *bear*. Furthermore, define the mean square error \mathcal{E}_i between the distribution $P(S_i)$ and the stationary distribution P_∞ you've calculated in the previous part as

$$\mathcal{E}_i = \frac{1}{2} \sqrt{(P(S_i = \text{bull}) - P_\infty(\text{bull}))^2 + (P(S_i = \text{bear}) - P_\infty(\text{bear}))^2}.$$

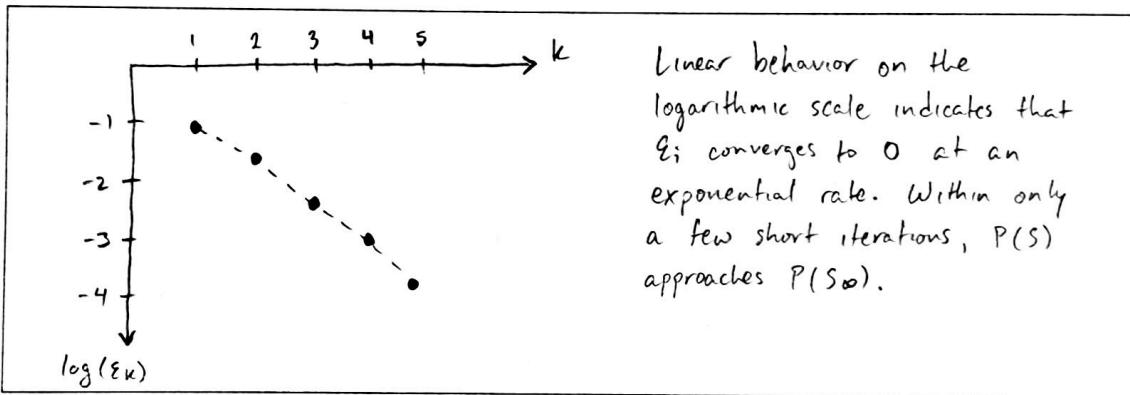
Compute the \mathcal{E}_i for the 5 steps. For the error terms, use scientific notation and keep 3 significant digits. Fill in the table below.

Note: You might find it helpful to use software such as Julia¹, NumPy or Matlab to compute these values for you or verify your hand-calculated results. If you have a linear algebra background, you might also find it helpful to express the calculation in terms of matrix vector multiplications.

¹Our reference Julia solution for the previous part and this part combined takes 12 lines.

	<i>bull</i>	<i>bear</i>	\mathcal{E}_i
$P(S_1)$	0.5	0.5	8.84×10^{-2}
$P(S_2)$	0.4	0.6	1.77×10^{-2}
$P(S_3)$	0.38	0.62	3.54×10^{-3}
$P(S_4)$	0.376	0.624	7.07×10^{-4}
$P(S_5)$	0.3752	0.6248	1.41×10^{-4}

- (d) (3 points) Plot $\log \mathcal{E}_k$ with k as the horizontal axis. Comment on the convergence behavior of \mathcal{E}_i (i.e., how fast does it converge to 0?).



Comment: Google uses an algorithm based on this idea to compute the PageRank values for each webpage to inform its search ranking. However, as web search evolves, PageRank has become less and less important over time.

- (e) (2 points) Suppose a startup uses this model to predict the state of the market in the future. Their analysis software can make a precise observation about the state of the market of previous day S_{t-1} as soon as markets open each day t . One night, a power surge destroyed the server containing all the historic market records. The next day, you are called in to assess the damage.

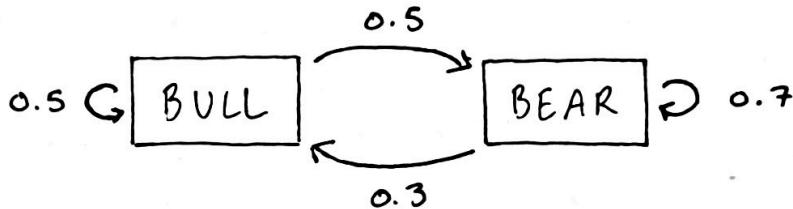
How has this event affected the accuracy of the market predictions of this startup?

This event will have no effect on predictions because the Markov model only uses data from the previous day, S_{t-1} .

Give an explanation for why you arrived at your assessment.

Since the startup will have data from the previous day on-hand, so long as the model remains unharmed, it will continue to make accurate predictions.

3. (a)



$$\begin{aligned}
 (b) \quad P_{\infty}(\text{bull}) &= P(\text{bull} | \text{bull}) P_{\infty}(\text{bull}) + P(\text{bull} | \text{bear}) P_{\infty}(\text{bear}) \\
 &= 0.5 P_{\infty}(\text{bull}) + 0.3 P_{\infty}(\text{bear}).
 \end{aligned}$$

$$\begin{aligned}
 P_{\infty}(\text{bear}) &= P(\text{bear} | \text{bear}) P_{\infty}(\text{bear}) + P(\text{bear} | \text{bull}) P_{\infty}(\text{bull}) \\
 &= 0.7 P_{\infty}(\text{bear}) + 0.5 P_{\infty}(\text{bull}).
 \end{aligned}$$

$$P_{\infty}(\text{bull}) + P_{\infty}(\text{bear}) = 1.$$

Solving this system of equations, we have

$$\begin{aligned}
 P_{\infty}(\text{bull}) &= 0.375 \\
 P_{\infty}(\text{bear}) &= 0.625.
 \end{aligned}$$

(c) $S_0 = \text{bull}.$

$$P(S_1 = \text{bull}) = P(\text{bull} | \text{bull}) = 0.5$$

$$P(S_1 = \text{bear}) = P(\text{bear} | \text{bull}) = 0.5$$

$$\begin{aligned}
 P(S_2 = \text{bull}) &= P(\text{bull} | \text{bull}) P(S_1 = \text{bull}) + P(\text{bull} | \text{bear}) P(S_1 = \text{bear}) = 0.4 \\
 &= (0.5)(0.5) + (0.3)(0.5)
 \end{aligned}$$

$$\begin{aligned}
 P(S_2 = \text{bear}) &= P(\text{bear} | \text{bull}) P(S_1 = \text{bull}) + P(\text{bear} | \text{bear}) P(S_1 = \text{bear}) = 0.6 \\
 &= (0.5)(0.5) + (0.7)(0.5)
 \end{aligned}$$

$$\begin{aligned}
 P(S_3 = \text{bull}) &= P(\text{bull} | \text{bull}) P(S_2 = \text{bull}) + P(\text{bull} | \text{bear}) P(S_2 = \text{bear}) = 0.38 \\
 &= (0.5)(0.4) + (0.3)(0.6)
 \end{aligned}$$

$$\begin{aligned}
 P(S_3 = \text{bear}) &= P(\text{bear} | \text{bull}) P(S_2 = \text{bull}) + P(\text{bear} | \text{bear}) P(S_2 = \text{bear}) = 0.62 \\
 &= (0.5)(0.4) + (0.7)(0.6)
 \end{aligned}$$

$$P(S_4 = \text{bull}) = (0.5)(0.38) + (0.3)(0.62) = 0.376$$

$$P(S_4 = \text{bear}) = (0.5)(0.38) + (0.7)(0.62) = 0.624$$

$$P(S_5 = \text{bull}) = (0.5)(0.376) + (0.3)(0.624) = 0.3752$$

$$P(S_5 = \text{bear}) = (0.5)(0.376) + (0.7)(0.624) = 0.6248$$

$$\mathcal{E}_i = \frac{1}{2} \sqrt{(P(S_i = \text{bull}) - P_{\infty}(\text{bull}))^2 + (P(S_i = \text{bear}) - P_{\infty}(\text{bear}))^2}$$

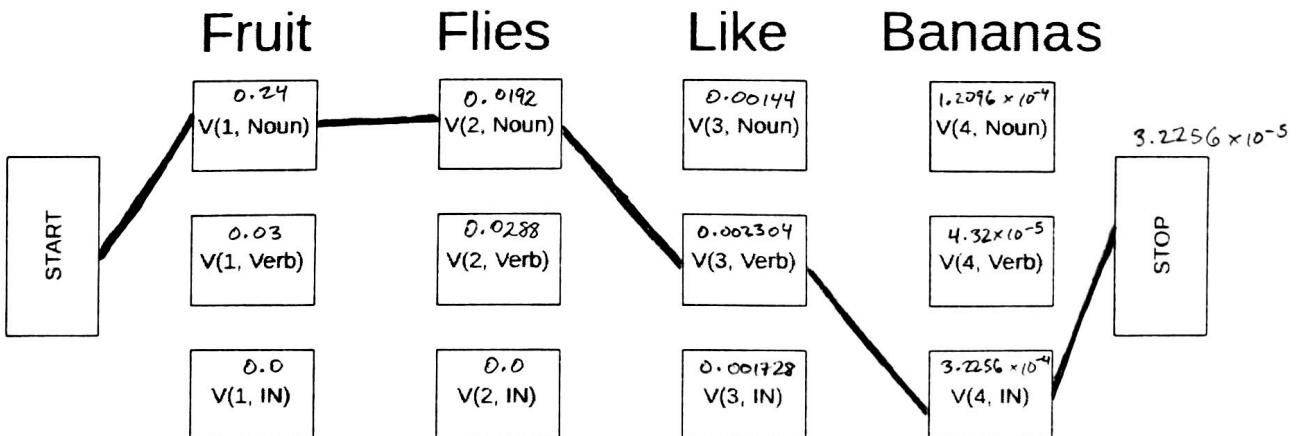
Calculate \mathcal{E}_i via Matlab.

4 The Viterbi Algorithm

(35 Points) One of the applications where the Viterbi algorithm can help us is in POS (Part-Of-Speech) tagging in Natural Language Processing. That is, given a sequence of words (e.g., a sentence), assign the correct grammatical tag (Noun, Verb, etc.) to each word based on its definition and context.

- Read more about POS-tagging [here](#)
- Read more about the Viterbi algorithm [here](#)
- Watch a video to help you understand how to approach this problem [here](#)

In this problem we perform POS tagging using the Viterbi algorithm for a small sentence and only considering Nouns (N), Prepositions or subordinating conjunctions (IN) and Verbs (V). While we do not use any other POS tags in this question, if you are curious you can find a list of them [here](#). For the setup below we will be using the Viterbi algorithm to first determine the score for each POS tag for each word and then determine the most probable POS-tag sequence for the entire sentence based on the scores we have calculated.



Let Y_i be random variables representing the emissions (at position i) which are English words and X_i be random variables representing the POS tags. $P(y_i|x_i)$ represents the emission probability of word y_i given tag x_i and $P(x_i|x_{i-1})$ be probability of transitioning to tag x_i

from tag x_{i-1} . Therefore, for computing score V_{i,x_i} for each state, we have:

$$V_{i,x_i} = \max_{x_1 \dots x_{i-1}} P(x_1 \dots x_i, y_1 \dots y_i) = \max_{x_{i-1}} P(y_i|x_i)P(x_i|x_{i-1})V_{i-1,x_{i-1}}$$

- (a) (25 points) In this part we want to compute the maximum score of the tags for each word. Below we have provided the emission and transition probabilities, you will have to use them and the formula above to calculate the score for each state. Show all calculations (you can use a separate page if you want but be sure to include it in the Gradescope submission), please put your maximum score above the corresponding tag box (e.g. put a number above $V(1, N)$ box). DO NOT round any numbers/final answers and show all calculations.

Transition model:

$$\begin{array}{llll} P(N|START) = 0.6 & P(N|N) = 0.4 & P(N|V) = 0.5 & P(N|IN) = 0.7 \\ P(V|START) = 0.3 & P(V|N) = 0.3 & P(V|V) = 0.1 & P(V|IN) = 0.1 \\ P(IN|START) = 0.1 & P(IN|N) = 0.1 & P(IN|V) = 0.2 & P(IN|IN) = 0.1 \\ P(STOP|START) = 0.0 & P(STOP|N) = 0.2 & P(STOP|V) = 0.2 & P(STOP|IN) = 0.1 \end{array}$$

Emission model:

$$\begin{array}{lll} P(fruit|N) = 0.4 & P(fruit|V) = 0.1 & P(fruit|IN) = 0.0 \\ P(flies|N) = 0.2 & P(flies|V) = 0.4 & P(flies|IN) = 0.0 \\ P(like|N) = 0.1 & P(like|V) = 0.4 & P(like|IN) = 0.3 \\ P(bananas|N) = 0.1 & P(bananas|V) = 0.1 & P(bananas|IN) = 0.7 \end{array}$$

- (b) (10 points) Now that we have computed the score, determine the maximum-probability sequence of states by working backwards from the STOP state. Consider the STOP state as the first value of a variable CURRENT-STATE. At each stage, move to the left, to the state that was selected during the arg-maxing for CURRENT-STATE. Then make that state be CURRENT-STATE, and iterate until all the way left at START. Reading the sequence along this path from left to right, you'll have the maximum-probability sequence of states that could give rise to the sentence "Fruit flies like bananas." Clearly draw the path, and also write down the final maximum-probability tag sequence.

START → N → N → V → IN → STOP

$$4. (a) V(1, N) = P(\text{fruit}|N) P(N, \text{START}) V(0| \text{START}) = 0.4 \cdot 0.6 \cdot 1 = 0.24$$

$$V(1, V) = P(\text{fruit}|V) P(V, \text{START}) V(0| \text{START}) = 0.1 \cdot 0.3 \cdot 1 = 0.03$$

$$V(1, IN) = P(\text{fruit}|IN) P(IN, \text{START}) V(0| \text{START}) = 0.0 \cdot 0.1 \cdot 1 = 0$$

$$V(2, N) = \max \left\{ P(\text{flies}|N) P(N|N) V(1, N), P(\text{flies}|N) P(N|V) V(1, V), P(\text{flies}|IN) P(N|IN) V(1, IN) \right\}$$

$$= \max \left\{ 0.2 \cdot 0.4 \cdot 0.24, 0.2 \cdot 0.5 \cdot 0.03, 0 \right\}$$

$$= 0.0192$$

$$V(2, V) = \max \left\{ P(\text{flies}|V) P(V|N) V(1, N), P(\text{flies}|V) P(V|V) V(1, V), P(\text{flies}|V) P(V|IN) V(1, IN) \right\}$$

$$= \max \left\{ 0.4 \cdot 0.3 \cdot 0.24, 0.4 \cdot 0.1 \cdot 0.03, 0 \right\}$$

$$= 0.0288$$

$$V(2, IN) = \max \left\{ P(\text{flies}|IN) P(IN|N) V(1, N), P(\text{flies}|IN) P(IN|V) V(1, V), P(\text{flies}|IN) P(IN|IN) V(1, IN) \right\}$$

$$= \max \left\{ 0, 0, 0 \right\}$$

$$= 0.0$$

$$V(3, N) = \max \left\{ P(\text{like}|N) P(N|N) V(2, N), P(\text{like}|N) P(N|V) V(2, V), P(\text{like}|N) P(N|IN) V(2, IN) \right\}$$

$$= \max \left\{ 0.1 \cdot 0.4 \cdot 0.0192, 0.1 \cdot 0.5 \cdot 0.0288, 0 \right\}$$

$$= 0.00144$$

$$V(3, V) = \max \left\{ P(\text{like}|V) P(V|N) V(2, N), P(\text{like}|V) P(V|V) V(2, V), P(\text{like}|V) P(V|IN) V(2, IN) \right\}$$

$$= \max \left\{ 0.4 \cdot 0.3 \cdot 0.0192, 0.4 \cdot 0.1 \cdot 0.0288, 0 \right\}$$

$$= 0.002304$$

$$V(3, IN) = \max \left\{ P(\text{like}|IN) P(IN|N) V(2, N), P(\text{like}|IN) P(IN|V) V(2, V), P(\text{like}|IN) P(IN|IN) V(2, IN) \right\}$$

$$= \max \left\{ 0.3 \cdot 0.1 \cdot 0.0192, 0.3 \cdot 0.2 \cdot 0.0288, 0 \right\}$$

$$= 0.001728$$

$$V(4, N) = \max \left\{ P(\text{bananas}|N) P(N|N) V(3, N), P(\text{bananas}|N) P(N|V) V(3, V), P(\text{bananas}|N) P(N|IN) V(3, IN) \right\}$$

$$= \max \left\{ 0.1 \cdot 0.4 \cdot 0.00144, 0.1 \cdot 0.5 \cdot 0.002304, 0.1 \cdot 0.7 \cdot 0.001728 \right\}$$

$$= 1.2096 \times 10^{-4}$$

$$V(4, V) = \max \left\{ P(\text{bananas}|V) P(V|N) V(3, N), P(\text{bananas}|V) P(V|V) V(3, V), P(\text{bananas}|V) P(V|IN) V(3, IN) \right\}$$

$$= \max \left\{ 0.1 \cdot 0.3 \cdot 0.00144, 0.1 \cdot 0.1 \cdot 0.002304, 0.1 \cdot 0.1 \cdot 0.001728 \right\}$$

$$= 4.32 \times 10^{-5}$$

$$V(4, IN) = \max \left\{ P(\text{bananas}|IN) P(IN|N) V(3, N), P(\text{bananas}|IN) P(IN|V) V(3, V), P(\text{bananas}|IN) P(IN|IN) V(3, IN) \right\}$$

$$= \max \left\{ 0.7 \cdot 0.1 \cdot 0.00144, 0.7 \cdot 0.2 \cdot 0.002304, 0.7 \cdot 0.1 \cdot 0.001728 \right\}$$

$$= 5.2256 \times 10^{-4}$$

$$V(5, \text{STOP}) = \max \left\{ P(\text{STOP}|N) V(4, N), P(\text{STOP}|V) V(4, V), P(\text{STOP}|IN) V(4, IN) \right\}$$

$$= \max \left\{ 0.2 \cdot 1.2096 \times 10^{-4}, 0.2 \cdot 4.32 \times 10^{-5}, 0.1 \cdot 5.2256 \times 10^{-4} \right\}$$

$$= 3.2256 \times 10^{-5}$$

5 Disambiguating Syntax with PCFGs

(35 points) Consider the sentence “Mary analyzed the algorithm with an equation.” This might mean that, for example, Mary used an equation to analyze the equation, or rather differently it might mean that Mary analyzed the algorithm, which happened to contain an equation.

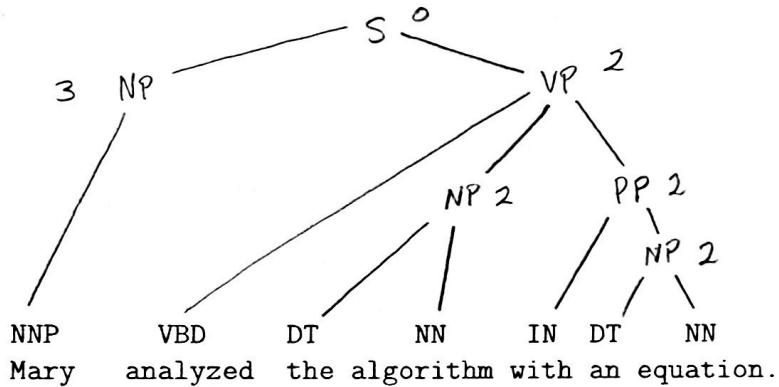
A grammar relevant to this example is given below. Consider the number at the right of a production to be the conditional probability of applying that production given that the symbol to be expanded, during a derivation, is the symbol on the left-hand side of the production.

In this problem, you’ll convert probabilities of productions into scores. Then, with the given probabilistic context-free grammar, you will find two legal parses for the sentence, and compute a score for each parse. You’ll then convert the overall parse scores back to probabilities of each parse. Then you’ll identify the more probable parse using the parse probabilities.

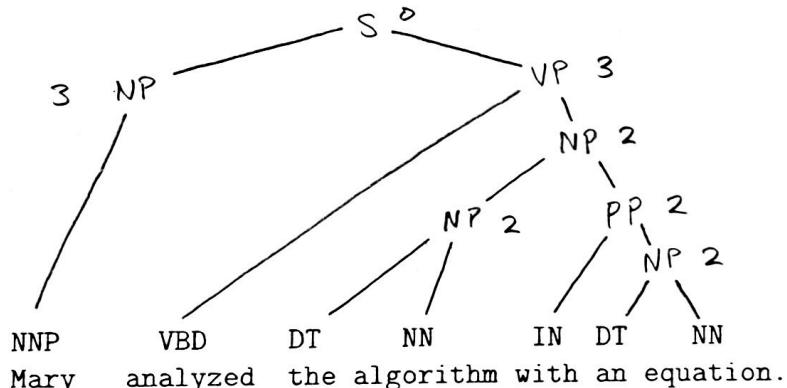
				<u>Scores</u>
S	$::=$	NP VP	1.000	0
NP	$::=$	NN	0.250	2
NP	$::=$	NP PP	0.250	2
NP	$::=$	NNP	0.125	3
NP	$::=$	DT NN	0.250	2
NN	$::=$	NNP	0.125	3
VP	$::=$	VBD NP PP	0.250	2
VP	$::=$	VBD NP	0.125	3
PP	$::=$	IN NP	0.250	2

- (a) (8 points) Scores for each production rule: Convert each probability into a score by taking score = $-\log_2(p)$. Write the scores next to the probabilities above.

- (b) (6 points) Find parse number 1. The parse will assume that the terminal symbols have been converted to non-terminals as shown. Make this parse correspond to the interpretation that Mary used an equation.



- (c) (6 points) Parse number 2. Make this parse according to an interpretation in which the algorithm contains an equation.



(d) (6 points) Total score 11 and overall probability 2^{-11} for parse number 1.

(e) (6 points) Total score 14 and overall probability 2^{-14} for parse number 2.

(f) (3 point) Which parse is more probable? parse 1 ($2^{-11} > 2^{-14}$)

You are welcome but not required to use the online parser at <https://parser.kitaev.io/> in this problem. Note that it doesn't return all parses of a sentence.

6 The Laws of Robotics and Ethics in AI

(30 points) Read the short story *Little Lost Robot* by Isaac Asimov. Answer the questions listed below the links (5 points each).

Short Story: <https://canvas.uw.edu/courses/1448875/files?preview=77766644>

- (a) List three facts about Isaac Asimov, including one that you think would be of interest to computer scientists.
- (b) What are Asimov's 3 Laws of Robotics?
- (c) Is the ordering of the laws important? If so, explain why (e.g. give an example of what might happen if the ordering of the laws changed). If not, explain why not. (You may find referring to this resource helpful: <https://xkcd.com/1613/>)
- (d) Can you describe an example of when a robot didn't obey Asimov's Laws? Do you think the example you described is consistent with the reported modification to the laws?
- (e) Did you find the end of the story satisfying – i.e. did it seem reasonable that the robot could be tricked to reveal that it could differentiate between different kinds of radiation? Explain.
- (f) In the story, Bogert and Calvin disagree on how potentially dangerous the modified robots are. Do you agree with Calvin's opinion that the robots are dangerous enough that all 63 robots should be destroyed if the lost robot can't be found? Explain.

6 The Laws of Robotics and Ethics in AI

(a) 1. Isaac Asimov was a Russian immigrant who attended New York City public schools and later became a Doctor of Philosophy in chemistry at Columbia. He would go on to become a professor of biochemistry at Boston University.

2. In addition to the science fiction works that he is renowned for such as the *Foundation* series, Asimov published over one hundred nonfiction books as well including lengthy biochemistry textbooks and scientific encyclopedias.

3. Computer scientists may be interested in the fact that Asimov was credited with the first use of the word ‘robotics’. Before “Runaround”, the word robot already existed, but the field of robotics was not yet formalized in the dictionary.

Credit: https://en.wikipedia.org/wiki/Isaac_Asimov

(b) Asimov’s 3 Laws of Robotics:

1. A robot may not injure a human being or, through inaction, allow a human being to come to harm.

2. A robot must obey the orders given to it by human beings except where such orders would conflict with the First Law.

3. A robot must protect its own existence as long as such protection does not conflict with the First or Second Laws.

(c) The ordering of the three laws and subsequent hierarchy is key to the stability of a world where humans are able to use robots in ethical and effective ways. As the xkcd comic describes, a re-ordering of the three laws would either result in unproductive robots or potential damage to other humans, robots, or both. For instance, if the order of the three laws were reversed, a robot would prioritize its own existence, then obey orders so long as they did not endanger their own existence, and lastly would prioritize not inflicting harm on a human being – but only if this did not also harm itself or violate a direct order. This way, if someone ordered a robot to kill another human, the robot would prioritize the value of the order given over the life of its intended target and robot assassins would dominate the “killbot hellscape”.

(d) The psychologist in the story references a hypothetical robot making a logical leap using the modified first law to drop a weight on someone computing that it would be able to catch the weight and do no harm to the human, but then redefining its system boundary to exclude the human and effectively dropping the weight without catching it. Along the same line of thought, I think an example of a robot failing to abide by Asimov’s laws that we will likely have to grapple with in the next decade will be the incalculable hypothetical damage caused by self-driving cars. For instance, if an AI driving a truck swerves to avoid someone on a bicycle who suddenly switches into traffic and instead chooses to strike a telephone pole instead of the cyclist to abide by the first law, the pole could do significantly more damage if it comes down on a group of pedestrians on the sidewalk. The processor within the car will likely not have the power or foresight to see such an outcome, but if it did, the reported modification to the laws would send the car into the cyclist.

(e) I personally found it unsatisfying that the robot’s fatal error was a miscalculation that the doctor blamed on computational latency. Assigning hubris to the robot also seemed a bit too anthropomorphizing as well. The trick at the end of the story relied on human error and

(according to the doctor) human-like emotion from the robot. Though the story already takes many liberties with the possibilities afforded to AI through advanced robotics, I would at least prefer some logical consistency in the competition between man and machine. The robot would have instantly known that the NS-2's would be unable to differentiate radiation and acted accordingly.

- (f) The issue of runaway AI is intriguing because it looks at the dividing line between life and artificial life. My understanding of the story is that Dr. Calvin wants all of the robots destroyed because a single defective agent could result in the runaway multiplication of dangerous superphysical threats. One of the core tenets of evolutionary biology is that all life seeks to sustain or prolong itself through procreation and defense. I am unconvinced that a robot abiding by Asimov's laws would seek to build replicas of itself, but a modified robot with a dangerous logical fallacy embedded in its central programming might take the liberty to do so if it calculates that the best way to protect humankind is to strengthen and multiply. In this regard, I would tend to agree with Dr. Calvin and destroy all of the robots for insurance against such an outcome.