

Instructions:

- You are encouraged to use no more than 3 hours to complete the exam, but you may use as long as you wish within the 48 hour window.
- You are allowed 1 page (front and back) cheat sheet. The cheat sheet must be scanned (or photographed with high resolution) and submitted along with the exam solutions. The time begins when you flip this page.
- You are also allowed a simple calculator, or Wolfram Mathematica to simplify equations. You may also write simple scripts in python or MATLAB to aid with your calculation. However, you do not need to for solving the problem, and you do not need to submit any code.
- You may print the exam and write your solutions or use lyx/latex. If you need extra sheets of paper, please label them carefully as to which question they are answering. Make your final answer clear.
- If you choose to handwrite your solutions, you must make sure that the digital scan / photograph is of high enough quality that we can see everything clearly. Anything we can't read, we will not grade.
- You may not discuss any problem with any other student while the exam submission portal is still open. You may not look for answers on the internet or in any notes outside of your cheatsheet.

Name: _____

Student ID: _____

| Scoring | |
|--------------|-------------|
| Q 1 | _____ / 10 |
| Q 2 | _____ / 10 |
| Q 3 | _____ / 20 |
| Q 4 | _____ / 20 |
| Q 5 | _____ / 20 |
| Q 6 | _____ / 20 |
| Total | _____ / 100 |

1. Classify the following as a supervised learning task (where training labels are needed) or unsupervised learning task (where training labels are not needed) **(2 point each)**

(a) Principal component analysis

Ans. Unsupervised

(b) logistic regression

Ans. Supervised

(c) Gaussian mixture model

Ans. Unsupervised

(d) boosted decision trees

Ans. Supervised

(e) clustering

Ans. Unsupervised

2. **True or False. (2 point each)**

(a) A graphical model with on average 2 directed edges pointing to each node requires less training samples for inference, than a similar graphical model with on average 5 directed edges pointing to each node.

Ans. True

(b) In a hidden Markov model, each hidden state is independent from all other hidden states

Ans. False. In a hidden Markov model, each state has a dependence on the previous and next state.

(c) In multiclass logistic regression, the event that a sample x belongs to class i is independent of the event that x belongs to class j whenever $i \neq j$.

Ans. False. The normalization term couples the probabilities of the event belonging to all the classes with each other.

(d) The difference between K-means and K-NN is that K-means is a supervised learning task and K-NN is an unsupervised learning task.

Ans. False. It is the other way around.

(e) Ensembling different models presents more advantages if each model acts as independently as possible.

Ans. True

3. I am given a bunch of animals and asked to classify them. I have some training data, given below, and I will use it to construct a decision tree, which I will use to classify future animals. In this problem, round all values to the nearest 0.001, and use log base 2.

| name | label | nose | ear shape | ear position |
|----------|----------|-------|-----------|--------------|
| Arthur | aardvark | small | round | high |
| D.W. | aardvark | small | round | high |
| Buster | rabbit | small | long | high |
| Bud | rabbit | small | long | high |
| Bitzi | rabbit | small | long | high |
| Francine | monkey | big | round | high |
| Muffy | monkey | big | round | low |
| Neal | moose | big | long | low |

- (a) What is the entropy of the labels over the entire dataset? **(3pts)**

Ans. Defining

$$\Pr(Y = \text{aardvark}) = \frac{2}{8}, \quad \Pr(Y = \text{rabbit}) = \frac{3}{8}, \quad \Pr(Y = \text{monkey}) = \frac{2}{8}, \quad \Pr(Y = \text{moose}) = \frac{1}{8}$$

the entropy of the collection is

$$\begin{aligned} H(X) &= -\Pr(Y = \text{aardvark}) \log_2(\Pr(Y = \text{aardvark})) - \Pr(Y = \text{rabbit}) \log_2(\Pr(Y = \text{rabbit})) \\ &\quad - \Pr(Y = \text{monkey}) \log_2(\Pr(Y = \text{monkey})) - \Pr(Y = \text{moose}) \log_2(\Pr(Y = \text{moose})) \\ &= -\frac{2}{8} \log_2\left(\frac{2}{8}\right) - \frac{3}{8} \log_2\left(\frac{3}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right) - \frac{1}{8} \log_2\left(\frac{1}{8}\right) \\ &\approx 1.906 \end{aligned}$$

- (b) For each feature, compute the information gain if that feature was to be used in the first split.

Nose (3pts)

$$\begin{aligned} H(X|\text{nose}) &= -\Pr(Y = \text{aardvark}|\text{nose} = \text{small}) \log_2(\Pr(Y = \text{aardvark}|\text{nose} = \text{small})) \\ &\quad - \Pr(Y = \text{rabbit}|\text{nose} = \text{small}) \log_2(\Pr(Y = \text{rabbit}|\text{nose} = \text{small})) \\ &\quad - \Pr(Y = \text{monkey}|\text{nose} = \text{big}) \log_2(\Pr(Y = \text{monkey}|\text{nose} = \text{big})) \\ &\quad - \Pr(Y = \text{moose}|\text{nose} = \text{big}) \log_2(\Pr(Y = \text{moose}|\text{nose} = \text{big})) \\ &= -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) \\ &\approx 0.951 \end{aligned}$$

$$IG = H(X) - H(X|\text{nose}) = 1.906 - 0.951 = 0.954$$

Ear shape (3pts)

$$\begin{aligned} H(X|\text{ear shape}) &= -\Pr(Y = \text{aardvark}|\text{ear shape} = \text{round}) \log_2(\Pr(Y = \text{aardvark}|\text{ear shape} = \text{round})) \\ &\quad - \Pr(Y = \text{rabbit}|\text{ear shape} = \text{long}) \log_2(\Pr(Y = \text{rabbit}|\text{ear shape} = \text{long})) \\ &\quad - \Pr(Y = \text{monkey}|\text{ear shape} = \text{round}) \log_2(\Pr(Y = \text{monkey}|\text{ear shape} = \text{round})) \\ &\quad - \Pr(Y = \text{moose}|\text{ear shape} = \text{long}) \log_2(\Pr(Y = \text{moose}|\text{ear shape} = \text{long})) \\ &= -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \\ &\approx 0.906 \end{aligned}$$

$$IG = H(X) - H(X|\text{ear shape}) = 1.906 - 0.906 = 1.0$$

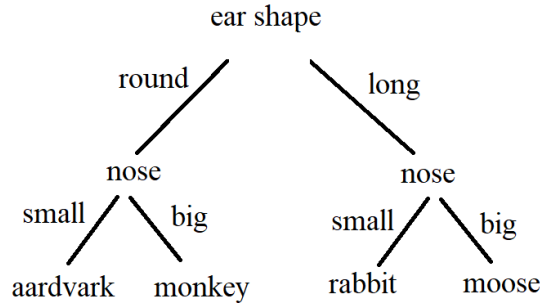
Ear position (3pts)

$$\begin{aligned}
 H(X|\text{ear position}) &= -\Pr(Y = \text{aardvark}) \log_2(\Pr(Y = \text{aardvark}|\text{ear position} = \text{high})) \\
 &\quad -\Pr(Y = \text{rabbit}) \log_2(\Pr(Y = \text{rabbit}|\text{ear position} = \text{high})) \\
 &\quad -\Pr(Y = \text{monkey}, \text{ear position} = \text{high}) \log_2(\Pr(Y = \text{monkey}|\text{ear position} = \text{high})) \\
 &\quad -\Pr(Y = \text{monkey}, \text{ear position} = \text{low}) \log_2(\Pr(Y = \text{moose}|\text{ear position} = \text{low})) \\
 &\quad -\Pr(Y = \text{moose}) \log_2(\Pr(Y = \text{moose}|\text{ear position} = \text{low})) \\
 &= -\frac{2}{8} \log_2\left(\frac{2}{6}\right) - \frac{3}{8} \log_2\left(\frac{3}{6}\right) - \frac{1}{8} \log_2\left(\frac{1}{6}\right) - \frac{1}{8} \log_2\left(\frac{1}{2}\right) - \frac{1}{8} \log_2\left(\frac{1}{2}\right) \\
 &\approx 1.344
 \end{aligned}$$

$$IG = H(X) - H(X|\text{ear position}) = 1.906 - 1.344 = 0.561$$

- (c) Construct the decision tree, using node purity to pick which node to split next, and information gain to pick which feature/split to make. Draw the final decision tree here. **(6pts)**

Ans.



- (d) Using this tree, infer the labels of the animals in the following test set and report the test error. **(2pts)**

Ans.

| name | true label | inferred label (fill in) | nose | ear shape | ear location |
|---------|------------|--------------------------|-------|-----------|--------------|
| Bambi | moose | aardvark | small | round | high |
| Thumper | rabbit | moose | big | long | low |
| George | monkey | monkey | big | round | low |

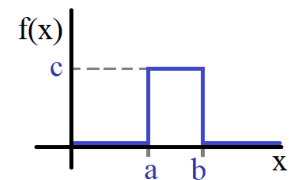
Test error = 2/3

4. **Gradient boosting using pulse functions** I am now faced with the task of estimating people's height based on their age. Again, I have some training data, which I will use to construct a boosted regression model. In this problem, round all values to the nearest 0.001.

| training sample | name | age (x_i) | height in ft (y_i) |
|-----------------|--------|---------------|------------------------|
| 1 | Alice | 25 | 6 |
| 2 | Brian | 60 | 5.5 |
| 3 | Carlos | 15 | 5 |
| 4 | Dianne | 1 | 2 |
| 5 | Esther | 5 | 3 |
| 6 | Freddy | 35 | 6 |

The goal is to return a function $f(\text{age}) \approx \text{height}$. We will use gradient boosting, with the set of weak regressors as simple pulse functions, parametrized by scalars a , b , and c .

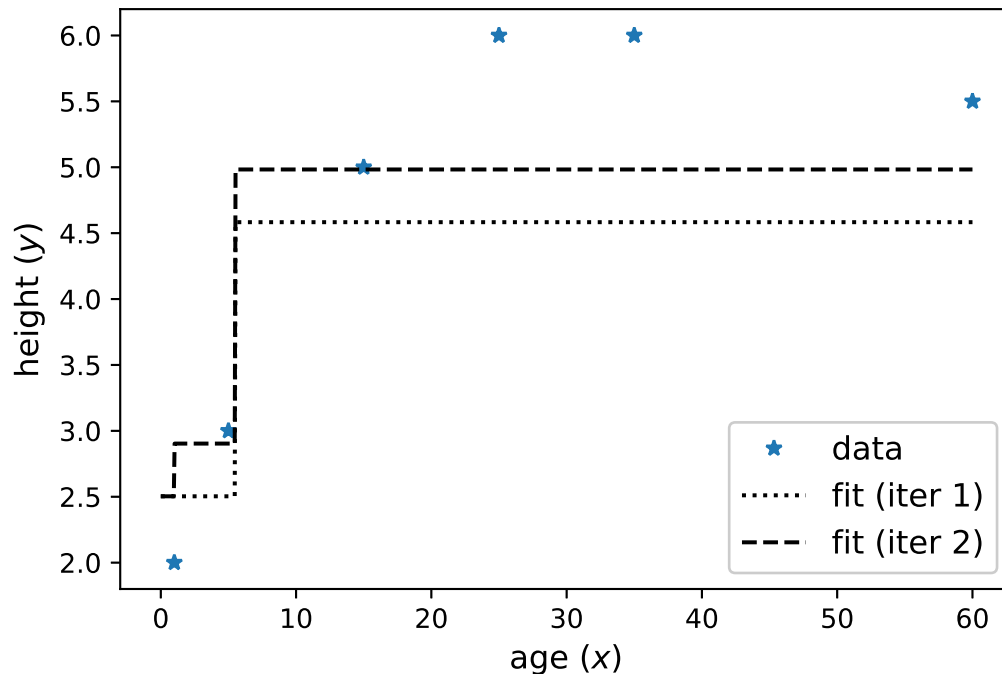
$$\mathcal{H} = \left\{ f : f(x) = \begin{cases} c & \text{if } a < x < b \\ 0 & \text{else.} \end{cases} \right\}$$



Formally, our goal is to minimize the squared error, e.g. solving

$$\begin{aligned} & \underset{f_t \in \mathcal{H}}{\text{minimize}} && \sum_{i=1}^m (F(x_i) - y_i)^2 \\ & \text{subject to} && F = f_1 + \dots + f_T. \end{aligned}$$

Ans.



(cont.)

- (a) **(2pts)** For the training set, graph height vs age, using the graph on the previous page.
(b) **(1pts)** For the “zeroth” weak learner, we just compute the mean label

$$f^{(0)}(x_i) = y_0 = \frac{1}{m} \sum_{i=1}^m y_i.$$

What is the mean squared error of picking the average height y_0 as the prediction for all the datapoints? (Round to nearest 0.01.) **Ans.**

$$f^{(0)}(x_i) = y_0 = \frac{1}{6}(6 + 5.5 + 5 + 2 + 3 + 6) = 4.58$$

$$\mathbf{m.s.e.} = \frac{1}{6}((6 - y_0)^2 + (5.5 - y_0)^2 + (5 - y_0)^2 + (2 - y_0)^2 + (3 - y_0)^2 + (6 - y_0)^2) = 2.37$$

- (c) **Iteration 1** To find the first weak learner, we will solve the least squares problem, where $z_i^{(0)} = y_i - F^{(0)}(x_i)$ is the current residual.

$$\underset{f \in \mathcal{H}}{\text{minimize}} \quad \sum_{i=1}^m (f(x_i) - z_i^{(0)})^2. \quad (1)$$

- i. **(2pts)** Frame problem (1) as an optimization problem over a , b , and c .

$$\underset{a, b, c}{\text{minimize}} \quad \sum_{a < x_i < b} (c - z_i^{(0)})^2 + \sum_{x_i > b, x_i < a} (z_i^{(0)})^2.$$

- ii. **(2pts)** I have optimized this problem and found that the optimum values of a and b are $a = 0$ and $b = 5.5$. Find the best value of c , rounding to the nearest 0.1.

Ans. There are a number of ways to getting to this answer. Using some python code and just sweeping values is one choice. Another is to notice that, with fixed a and b as given, the problem reduces to

$$\underset{c}{\text{minimize}} \quad \sum_{a < x_i < b} (c - z_i^{(0)})^2 = (c - z_4^{(0)})^2 + (c - z_5^{(0)})^2$$

which is achieved by picking the mean

$$c = \frac{z_4^{(0)} + z_5^{(0)}}{2} = \frac{2.0 + 3.0}{2} - 4.58 = -2.08 \approx -2.1.$$

- iii. **(2pts)** For the aggregate learner $F^{(1)}(x) = f^{(1)}(x) + f^{(0)}(x)$, draw $(x_i, F^{(1)}(x_i))$ on the plot in the previous page.
iv. **(2pts)** Compute the current loss and new residual vector $z^{(1)}$. Round to the nearest 0.01.

$$\mathbf{loss} = 5.53, \quad \mathbf{residual} \ z^{(1)} = (1.42, 0.92, 0.42, -0.58, 0.42, 1.42)$$

- (d) **Iteration 2** To find the second weak learner, we will again solve (1), replacing $z^{(0)}$ with the newly computed residual vector $z^{(1)}$. Again, I have optimized this problem and found that the optimum values of a and c are $a = 1$ and $c = 0.4$.

- i. **(2pts)** Find the best value of b , rounding to the nearest 0.1.

Ans. This one is a little harder to do analytically. Since $b > a$, we can consider this case by case:

- $1 < b < 5$, which gives an mse of 0.92
- $5 < b < 15$, which gives an mse of 0.88
- $15 < b < 25$, which gives an mse of 0.85
- $25 < b < 35$, which gives an mse of 0.69
- $35 < b < 60$, which gives an mse of 0.53

- $b > 60$, which gives an mse of 0.43

So, any value of b larger than 60 will minimize the loss.

¹

- ii. **(3pts)** For the aggregate learner $F^{(2)}(x) = f^{(2)}(x) + f^{(1)}(x) + f^{(0)}(x)$, draw $(x_i, F^{(1)}(x_i))$ on the plot from before.

Ans. there can be some ambiguity as to where the two vertical boundaries lie.

- iii. **(2pts)** Compute the current loss and new residual vector $z^{(2)}$. Round to the nearest 0.01.

$$\text{loss} = 2.60, \quad \text{mse} = 0.43, \quad \text{residual } z^{(2)} = (1.02, 0.52, 0.017, -0.50, 0.097, 1.02)$$

- (e) **(2pts)** Using this final aggregated learner, predict and give the test set heights based on age, by filling in the table below. Round to the nearest 0.01. Report the test mean squared error.

Ans.

| person | age | true height | fitted height |
|--------|-----|-------------|---------------|
| Gary | 12 | 5 | 4.98 |
| Harris | 57 | 5 | 4.98 |
| Ivy | 3 | 2 | 2.9 |

Test mean squared error: 0.27

¹There is a bug in the problem, in that the minimum values of a, b, c should have been $a = 15.0, b > 60, c = 1.20$. So, anyone who wrote anything here got full credit, and we used that value of b to calculate the rest of the problem.

5. **Kmeans and Gaussian mixture models** I am on an alien planet, and there are three types of fruit. They are of an assortment of sizes, which I have measured and given below, in centimeters:

$$x_1 = 6, \quad x_2 = 1, \quad x_3 = 12, \quad x_4 = 14, \quad x_5 = 3, \quad x_6 = 5, \quad x_7 = 1, \quad x_8 = 1, \quad x_9 = 2$$

- (a) **(3pts)** Using the 3 initial cluster center values of $c_1 = 1, c_2 = 2, c_3 = 3$, list the points that are assigned to each cluster in this first iteration. (Write down x_i and the values of i .)

| points in cluster 1 | points in cluster 2 | points in cluster 3 |
|---------------------|---------------------|---------------------------|
| x_7, x_8, x_2 | x_9 | x_1, x_3, x_4, x_5, x_6 |

- (b) **(3pts)** Perform one iteration of Kmeans (using 2 norm distance) give the new cluster centers.

$$c_1 = 1, \quad c_2 = 2, \quad c_3 = \frac{6 + 12 + 14 + 3 + 5}{5} = 8$$

- (c) An alien approaches me and tells me that the fruits have names, called "Ajax", "Basic", and "C++".

- Ajax is the smallest type of fruit, with a mean size of 30cm and a variance of 25cm.
- Basic is a medium sized fruit, with mean size 50cm, and a variance of 5cm.
- C++ is a large sized fruit, with mean size 60cm, and a variance of 100cm.

Additionally, Ajax is a super common fruit, occurring 3x as often as the other two, combined. Basic occurs 4x as often as C++.

(4pts) Fit a Gaussian mixture model to this story. That is, write out the PDF to model the size of a berry picked at random from this planet.

$$\begin{aligned} f_{\text{Ajax}}(x) &= \frac{1}{5\sqrt{2\pi}} \exp\left(-\frac{(x-30)^2}{50}\right) \\ f_{\text{Basic}}(x) &= \frac{1}{\sqrt{5} \cdot 2\pi} \exp\left(-\frac{(x-50)^2}{10}\right) \\ f_{\text{C++}}(x) &= \frac{1}{10\sqrt{2\pi}} \exp\left(-\frac{(x-60)^2}{200}\right) \\ f_{\text{fruit size}}(x) &= \frac{3}{4}f_{\text{Ajax}}(x) + \frac{1}{5}f_{\text{Basic}}(x) + \frac{1}{20}f_{\text{C++}}(x) \end{aligned}$$

- (d) **(5pts)** Using the above model, return the probability that a berry of size 40cm is Ajax. Hint: Remember Bayes' rule. Round to the nearest 0.01.

$$\begin{aligned} \Pr(\text{Ajax}|40) &= \frac{\Pr(40|\text{Ajax})\Pr(\text{Ajax})}{\Pr(40|\text{Ajax})\Pr(\text{Ajax}) + \Pr(40|\text{Basic})\Pr(\text{Basic}) + \Pr(40|\text{C++})\Pr(\text{C++})} \\ &\approx 0.97 \end{aligned}$$

- (e) **(5pts)** Ajax and C++ are poisonous, and Basic is not. I find a fruit of size 50cm and eat it. What is the probability that I have been poisoned? Round to the nearest 0.01.

$$\Pr(\text{Poisoned}|50) = \Pr(\text{Ajax}|50) + \Pr(\text{C++}|50) \approx 0.03$$

6. **Teamwork.** Alice, Bob, and Carlos are all working on their problem set together. Each question is a True/False question. The group works by each guessing the answer, and writing down the majority vote answer.

- (a) **(3pts)** Given that each friend is independently 80% sure that the answer is False, what is the probability that they will write down the answer True? Round to the nearest 0.001. **Ans.**

$$\Pr(2T's, 1F) = \underbrace{\binom{3}{2}}_{\text{"3 choose 2" = 3}} \cdot 0.2^2 \cdot 0.8 + 0.2^3 \approx 0.104$$

- (b) **(3pts)** Suppose that each friend independently has p chance of being correct, where p is a probability value between 0 and 1. Given the guesses of each three friends, we define the maximum likelihood guess to be

$$\hat{s}_{MLE} = \operatorname{argmax}_{\hat{s} \in \{T, F\}} \Pr(\text{answer is } \hat{s} | A, B, C).$$

Show that as long as $p > 0.5$, then \hat{s}_{MLE} is in fact the guess given by majority vote.

Hint 1: show that the probability that the majority vote is correct is also > 0.5 .

Hint 2: In general, $f(x) > c$ for all $x > b$ if f is strictly monotonically increasing and $f(b) \geq c$.

Ans.

$$\Pr(\text{majority vote guess is right}) = \binom{3}{2} \cdot p^2 \cdot (1-p) + p^3 = 3p^2 - 2p^3.$$

To show that this value is always > 0.5 if $p > 0.5$, it suffices to show that for $f(p) = 3p^2 - 2p^3$,

$$f'(p) = 6p - 6p^2 > 0 \forall p < 1, \quad f(0.5) = 0.5.$$

To deal with the last case of $p = 1$, note that $f(1) = 1$. Thus the probability of majority vote being correct, for $p > 0.5$, is always > 0.5 .

- (c) **(3pts)** Give the mean and variance of \hat{s}_{MLE} , given that Alice, Bob, and Carlos are all simultaneously and independently p percent sure the answer is True. Use the numerical encoding $\hat{s} = 1$ if the answer is True and $\hat{s} = 0$ if the answer is False. (No need to simplify the equations.)

Hint: recall the mathematical definition of expectations.

Ans.

$$\begin{aligned} \mathbb{E}[\hat{s}_{MLE}] &= \Pr(\hat{s}_{MLE} = 1) = 3p^2 - 2p^3 \\ \text{Var}(\hat{s}_{MLE}) &= \mathbb{E}[\hat{s}_{MLE}^2] - \mathbb{E}[\hat{s}_{MLE}]^2 \\ &= \Pr(\hat{s}_{MLE} = 1) - \Pr(\hat{s}_{MLE} = 1)^2 \\ &= 3p^2 - 2p^3 - (3p^2 - 2p^3)^2 \end{aligned}$$

- (d) **(4pts)** Again, Alice, Bob, and Carlos all simultaneously and independently believe the answer is True with probability p . The professor walks in. The three ask her for the answer. She mumbles "I think it's True". The professor really isn't that bright, though, so rather than trusting her completely, the students decide on a modified majority vote scheme: if all of the students believe the answer is False, then they write False. But, otherwise, the students will write down True, following the advice of the professor. We call what the students write down \hat{s}_{MAP} , where $\hat{s}_{MAP} = 1$ if they write down True, and $\hat{s}_{MAP} = 0$ if they write down False. What is the mean and variance of \hat{s}_{MAP} ?

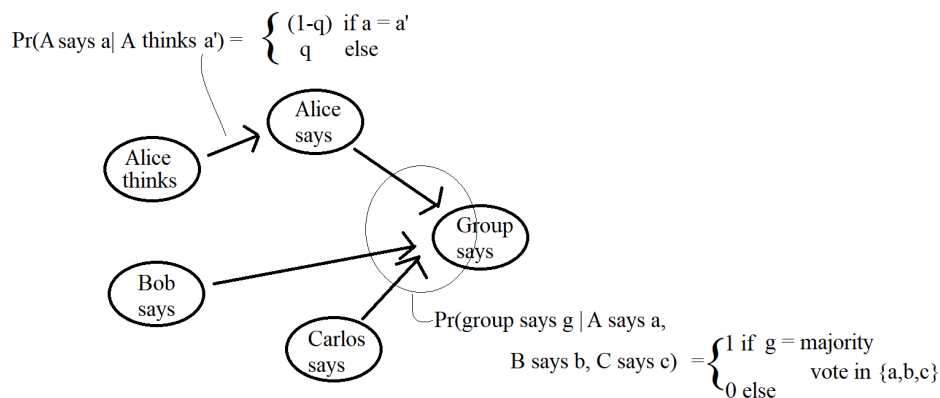
Ans.

$$\begin{aligned} \mathbb{E}[\hat{s}_{MAP}] &= 1 - \Pr(\hat{s}_{MAP} = 0) = 1 - (1-p)^3 \\ \text{Var}(\hat{s}_{MAP}) &= \mathbb{E}[\hat{s}_{MAP}^2] - \mathbb{E}[\hat{s}_{MAP}]^2 \\ &= \Pr(\hat{s}_{MAP} = 1) - \Pr(\hat{s}_{MAP} = 1)^2 \\ &= 1 - (1-p)^3 - (1 - (1-p)^3)^2 \end{aligned}$$

- (e) Now it is the next day, and the professor has gone on vacation; we must rely only on the guesses of the three students. While they are all good friends, the score overall is curved, so there is incentive for friends to sabotage each other, in order to get better scores. After doing some problems, Bob and Carlos start to suspect Alice of sabotaging.

If Alice is sabotaging, then her strategy is simply to flip her answers with probability q , where $0 \leq q \leq 1$.

- i. **(3pts)** Draw directed edges on the graphical model corresponding to this scenario, and write the conditional probability on each edge. **Ans.**



- ii. **(4pts)** Given that Alice, Bob, and Carlos are independently 70% correct on any given problem, and their majority voting scheme is also 70% correct, find q .

Ans.

$$\underbrace{\Pr(\text{correct})}_{0.7} = \underbrace{\Pr(\text{Bob, Carlos both correct})}_{0.7^2} + 2 \cdot \underbrace{\Pr(\text{Bob or Carlos correct, not both})}_{0.7 \cdot 0.3} \Pr(\text{Alice said correct})$$

Solved, this gives $\Pr(\text{Alice said correct}) = 0.5$. Breaking this down further,

$$\Pr(\text{Alice said correct}) = \Pr(\text{Alice correct, didn't flip}) + \Pr(\text{Alice incorrect, flipped}) = 0.7(1-q) + 0.3q = 0.7 - 0.4q$$

which gives $q = 0.5$.