

Name: Jonathan Farrell Kusuma  
Reg No: 499795092

## DSCI 552 MIDTERM

2 March 2023

For this exam one page of notes is allowed (both sides).

Calculators are allowed, but not smartphones, laptops or any device with internet connection.

The exam is 2 hours long and it is for 110 points. **You get a bonus of 10 points!**

**There are 6 problems and 20 pages total.**

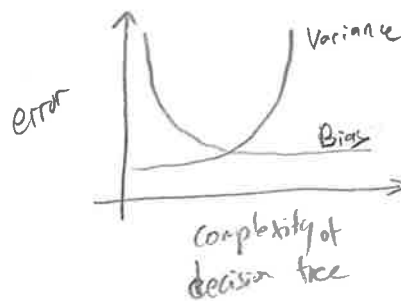
**Please remember to write your name**

Problem	Points
1	_____/42
2	_____/16
3	_____/14
4	_____/12
5	_____/12
6	_____/14
Total	_____/110

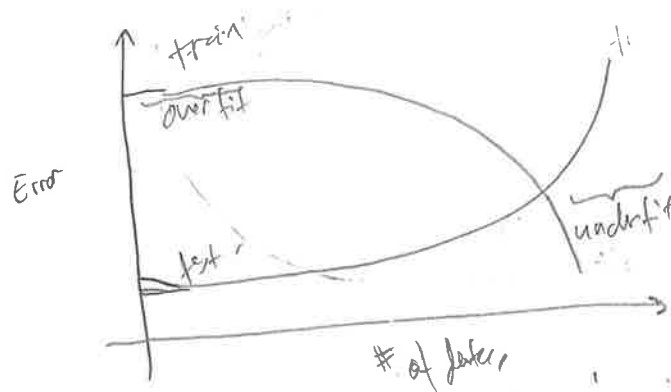
1. (42 points) Decision Trees and Bias/Variance Dilemma

- a. (6 points) Explain the bias/variance dilemma **specifically** in the context of **decision trees**. Draw a diagram of bias/variance to illustrate your explanation. Be sure to carefully label each part of your diagram.

The more complex the decision tree model is, the lower the bias, however the variance would be higher. Meanwhile, the less complex the model is, the lower the variance but higher bias.



- b. (6 points) Draw a diagram of train and test error curves that should be typical of decision trees. What is the relationship between train and test error curves to the curves in the bias/variance diagram?



the higher the train test, the higher the error.  
more training data increases ~~the~~ Variance  
and decreasing bias.

- c. (4 points) For the diagram in part b label the region where the decision tree is overfitting and where it is underfitting.

shown in (B)

- d. (6 points) What is the purpose of tree pruning? Describe the two types of tree pruning.

Prepruning: Stopping the tree earlier

Post pruning: Grow the whole tree then prune the subtrees that overfit on the pruning set

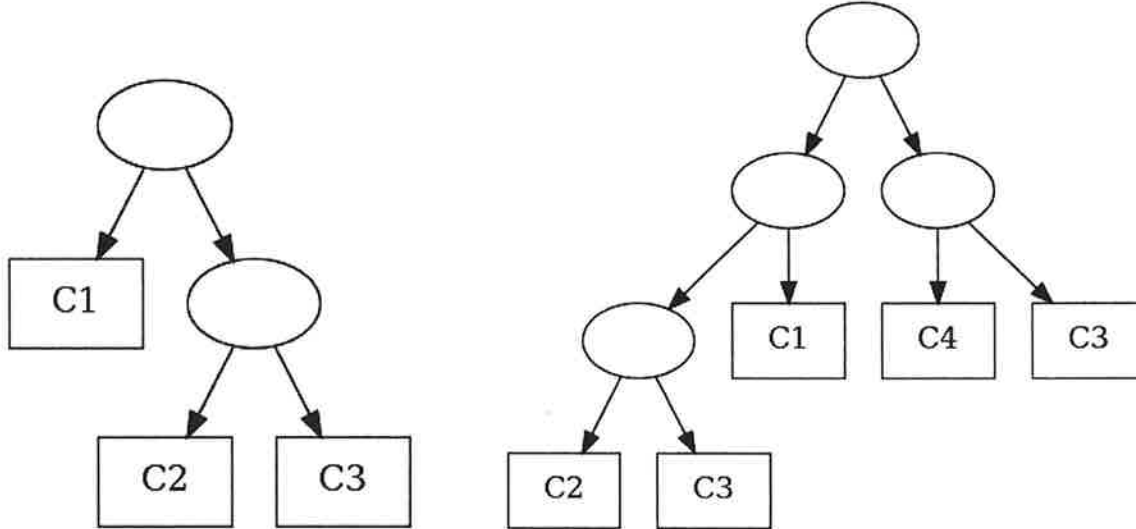
Purpose: Prevent overfitting for better generalization (decrease variance)

- e. (8 points) Minimum description length (MDL) principle.

Consider the two decision trees below. Assume they are generated from a dataset of 16 binary attributes and 4 classes,  $C_1$ ,  $C_2$ ,  $C_3$  and  $C_4$ . Assume

- Each internal node is coded using  $\log_2 d$  bits, where  $d$  is the number of attributes.
- Each leaf node is encoded using  $\log_2 K$  bits where  $K$  is the number of classes.
- For simplicity assume the cost of encode a tree is the total cost of encoding the internal nodes and leaf nodes.
- Each error is encoded using  $\log_2 N$  bits, where  $N$  is the number of training instances.

According to MDL principle which decision tree is better as a function of  $N$ ?



Tree 1 with 8 errors

Tree 2 with 4 error

Tree 1 because lower computing cost.

MDL: find a model that provides the most concise description of the data, quantifies model complexity in terms of the number of bits required to encode the data.

$$\log_2(6) + \log_2(2) + \log_2(3) = 5.585$$

$$\log_2(16) = 4$$

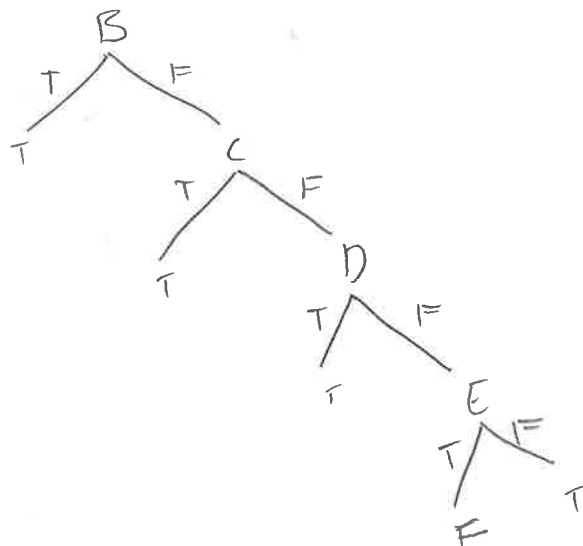
$$\log_2(4) = 2$$

$$\log_2(11) + \log_2(4) + \log_2(5) = 6.322$$

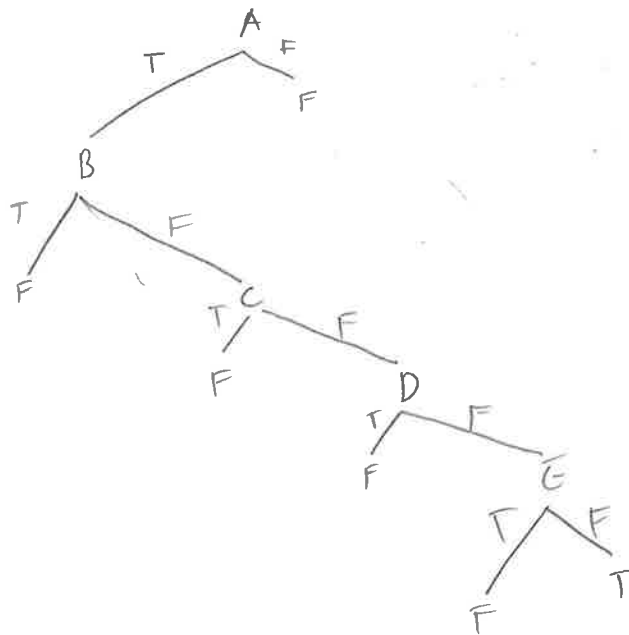
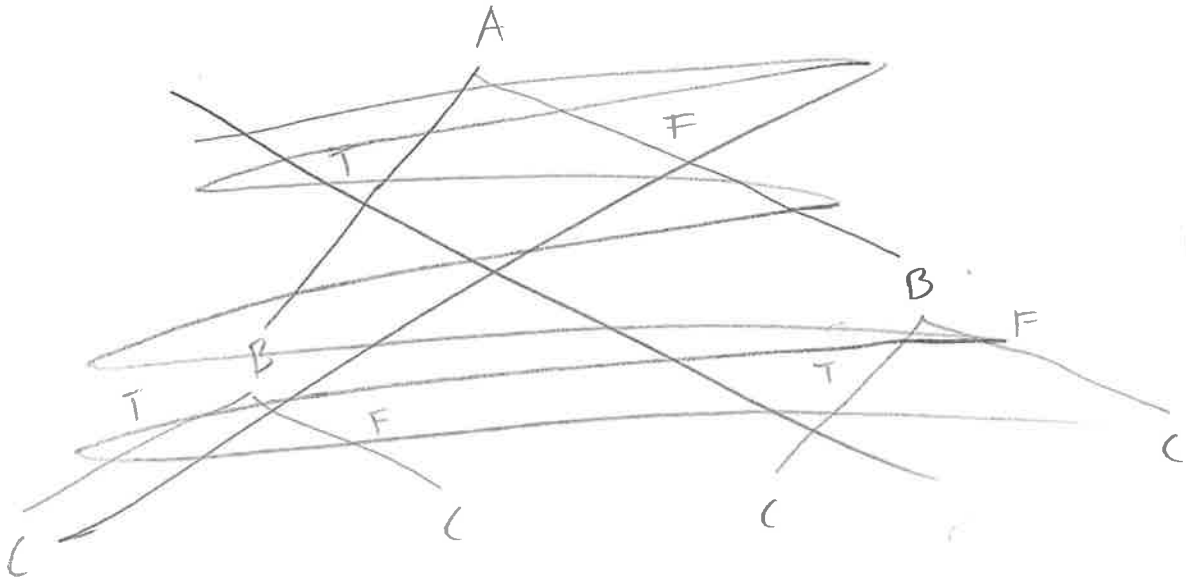
$$5.585 < 6.322$$

- f. (8 points) Domingos (2012) points out that overfitting can be caused by noise, but bad learning algorithms can also cause overfitting. For the Boolean training dataset below, draw a decision tree that will **only** classify correctly the positive instances in the training dataset and **no other positive instances** (it will ignore all negative instances).

A	B	C	D	E	Class
T	F	F	F	F	T
T	T	F	F	F	T
T	F	T	F	F	T
T	F	F	T	F	T
T	F	F	F	T	T
F	T	T	T	T	F
F	F	F	F	F	F
F	T	F	F	F	F
F	F	T	F	F	F
F	F	F	T	F	F
F	F	F	F	T	F



- g. (4 points) Using the dataset in the previous part, draw the smallest decision tree that will classify the entire dataset correctly with zero training error, i.e. without considering the **no other positive instances** restriction.



2. (16 points) Density estimation

- a. (4 points) An entomologist is studying the behaviors of dung beetles by collecting a dataset of the number of attempts individual dung beetles need to successfully push a ball of dung uphill. The dataset collected is a dataset of  $N$  beetles

$X = \{x^t\}$ , where beetle  $t$  failed on the first  $x^t - 1$  attempts, and succeeded on the last attempt. The entomologist assumes the beetles are not intelligent enough to learn across attempts, so he uses a geometric distribution

$p(x) = (1 - p_g)^{x-1} p_g$ , where  $p_g$  is the probability of success. Write down the likelihood equation for parameter  $p_g$ .

$$L = \prod_{t=1}^N p_g \cdot (1 - p_g)^{x^t - 1}$$

$$\ln L = \ln \prod_{t=1}^N p_g \cdot (1 - p_g)^{x^t - 1}$$

$$\ln L = \sum_{t=1}^N \ln p_g + \sum_{t=1}^N (x^t - 1) \ln(1 - p_g)$$

- b. (8 points) Derive maximum likelihood estimate of  $p_g$ .

$$\frac{dL}{dp_g} = \frac{d}{dp_g} \left( \prod_{t=1}^N p_g \cdot (1 - p_g)^{x^t - 1} \right) =$$

$$\frac{dL}{dp_g} = \sum_{t=1}^N 1 - \sum_{t=1}^N \frac{1}{p_g} + \sum_{t=1}^N \frac{1}{1 - p_g} = 0$$

finds this derivative equals to zero.  
finds the maximum of the likelihood estimate.

$$\hat{p}_g = \frac{1}{\sum_{t=1}^N x^t}$$

- c. (4 points) To the surprise of the entomologist the beetles in this dataset only needed about half the number of attempts as reported in entomology literature. Suppose the entomologist was able to obtain the prior density from literature. Write down the equation the entomologist needs to solve to incorporate the **prior density**.

$$\cancel{P(X)} = \cancel{P(X) / \theta_{prior}}$$

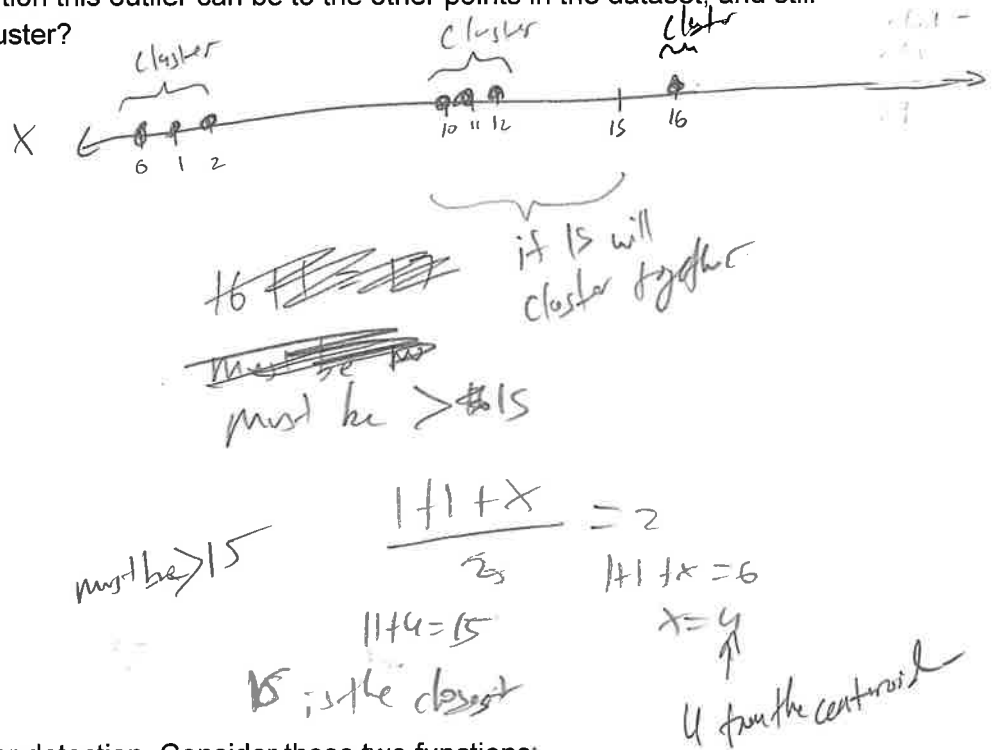
$$\hat{P}_0 = \frac{\sum_{i=1}^N \cancel{X_i}}{N}$$

$$\hat{P}_i = \frac{\sum_{t=1}^T X_t^+}{N}$$



3. (14 points) Clustering

- a. (8 points) Show K-mean clustering is not robust to outliers. Consider this one-dimensional dataset of 6 instances  $X = \{0, 1, 2, 10, 11, 12\}$ . For  $K=2$  clusters add one outlier to the dataset that will cause the K-mean clustering to place the outlier in its own cluster, and the rest of the dataset in the other cluster. What is the closest location this outlier can be to the other points in the dataset, and still be in its own cluster?



- b. (6 points) Outlier detection. Consider these two functions:

- $d_k(x)$ : the distance to the  $k$ -th nearest neighbor to instance  $x$
- $ave_k(x)$ : the average  $d_k(n)$  over  $n$ , where instance  $n$  is in the set of the  $k$  nearest neighbor of instance  $x$

Describe how to combine these two functions to use it for outlier detection, where  $k$  is a hyperparameter that we can change. Use the dataset in part a. to describe your solution.

Use the  $d_k(x)$  to calculate all the mean distance and use the  $ave_k(x)$  to compare if the newly calculated mean distance with the new point is higher or lower. If the mean distance is higher, then the point is likely an outlier.

4. (12 points) Dimension Reduction

- a. (4 points) In Principal Component Analysis (PCA) what does the eigenvalue  $\lambda_i$  of the  $i$ th component represent?

The eigenvalue  $\lambda$  is the distance from two points  
it is the projection directions in a covariance matrix  
of PCA.

- b. (4 points) What are the similarities and differences between Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA)?

PCA is an unsupervised dimension reduction that learns  
the projection directions that maximize data variance  
while LDA is a supervised learning method  
used for classification & dimensionality reduction.

- c. (4 points) Describe the distance metrics used by Isomap and Laplacian Eigenmaps. What is similar about these two metrics?

distance metric Isomap = Geodesic distance  
distance metric Laplacian Eigenmaps = uses a dot matrix  
of neighborhood preserving  
embeddings.

Both metrics are low-dimensional representation of  
data.

5. (12 points) Naive Bayes Classification

- a. (8 points) Use the Naive Bayes assumption and the dataset table below to classify the words: **Credit Card Deal**. Show your work, not just the final answer.

Words	SPAM
Interest Free Card	No
Cash Credit Gift	Yes
Mortgage Interest Deal	No
Cash Back Credit Card	No
Debt Free Deal	No
Credit Card Interest	No
Exclusive Free Deal	Yes
Card Interest Mortgage	Yes

CCD

$$S \rightarrow \frac{3}{8}$$

$$NS \rightarrow \frac{5}{8}$$

$$P(\text{Credit} | \text{Spam}) = \frac{1}{3}$$

$$P(\text{card} | \text{Spam}) = \frac{1}{3}$$

$$P(\text{deal} | \text{Spam}) = \frac{1}{3}$$

$$P(\text{Credit} | \text{NS}) = \frac{2}{5}$$

$$P(\text{card} | \text{NS}) = \frac{3}{5}$$

$$P(\text{deal} | \text{NS}) = \frac{2}{5}$$

$$\frac{1}{3} \times \frac{3}{8} + \frac{1}{3} \times \frac{3}{8} + \frac{1}{3} \times \frac{3}{8} = 0.375$$

$$\frac{2}{5} \times \frac{3}{5} + \frac{2}{5} \times \frac{2}{5} = 0.875$$

$$0.875 > 0.375$$

So Credit Card Deal will be classified as Not Spam.

$$P(S | \text{CCD}) =$$

$$P(N | \text{CCD}) =$$

$$\frac{P(\text{CCD} | S) P(S)}{P(\text{CCD})} =$$

$$\frac{P(\text{CCD} | \text{NS}) P(\text{NS})}{P(\text{CCD})} =$$

$$\frac{1}{3} \times \frac{1}{3} \times \frac{1}{3} = 0.037$$

$$\frac{2}{5} \times \frac{3}{5} \times \frac{2}{5} = 0.096$$

$$0.096 > 0.037$$

b. (4 points) Describe how you would classify the words: **Credit Card Promotion**.

for Not spam

to compare find the  $P(\text{credit} | \text{spam}) \times P(\text{spam})$  +  $P(\text{card} | \text{spam}) \times P(\text{spam})$  +  $P(\text{promotion} | \text{spam}) \times P(\text{spam})$  and  
 $\Rightarrow P(\text{credit} | \text{Nspam}) \times P(\text{Nspam}) + P(\text{card} | \text{Nspam}) \times P(\text{Nspam}) + P(\text{promotion} | \text{Nspam}) \times P(\text{Nspam})$  for spam

the higher probability will be classified with the class associated with the equation this uses the Bayes theorem.

\* We can't use Naive Bayes because  $P(\text{promotion} | \text{spam})$  is zero.

6. (14 points) Association rules

Use the dataset in Question 5. Given the association rule: **Interest**  $\rightarrow$  **Card**

a. (4 points) What is the support of this rule?

$$\begin{aligned} \text{Support}(\text{Interest} \rightarrow \text{Card}) &= P(\text{Interest}, \text{Card}) = \frac{\# \text{ of "interest" and "card"}}{\# \text{ data}} \\ &= \frac{3}{8} = 0.375 \end{aligned}$$

b. (4 points) What is the confidence of this rule?

$$\begin{aligned} \text{Confidence}(\text{Interest} \rightarrow \text{Card}) &= P(\text{Card} | \text{Interest}) \\ &= \frac{3}{4} = 0.75 \end{aligned}$$

c. (6 points) Show why if this rule has low confidence:

**Credit Interest**  $\rightarrow$  **Card**

Then this rule can be pruned:

**Interest**  $\rightarrow$  **Card Credit**

$$\begin{aligned} \text{Confidence}(\text{Interest} \rightarrow \text{Card}) &= P(\text{Card} | \text{Interest}) = \frac{3}{4} \\ \text{Confidence}(\text{Credit Interest} \rightarrow \text{Card}) &= \frac{0}{4} \end{aligned}$$

$$P(\text{Interest} | \text{Credit Card}) = \frac{1}{2}$$

because confidence of "credit interest" is zero

<extra sheet>

<extra sheet>

<extra sheet>



<extra sheet>

<extra sheet>

<extra sheet>

<extra sheet>