# CPSC 100 2016W2: Practice Midterm

## Data Representation

Convert the binary number 011111111111111 to hexadecimal and decimal

Convert the hexadecimal number 0x12CD3 to binary.

Convert the hexadecimal number 0x45 to decimal.

Convert the decimal number 1034 to binary.

Convert decimal number 132 to hexadecimal.

DNA has four nucleotides (A, C, T, and G) (i.e., a base-4 system). Codons are "words" of three nucleotides that code for amino acids. For example AAA codes for an amino acid called lysine. What is the maximum number of amino acids that could be represented if we use four nucleotides and each codon can only be made up of three nucleotides?

There are just 20 amino acids. In a base-4 system (i.e., a system where each "digit" can be A, C, T, or G), what is the smallest number of digits that I need per position, in order to be able to represent all 20 amino acids using three-nucleotide codons?

What is rasterization?

Henry's friend has given him the dimensions (the number of rows and columns) of an image, and the number of bytes used to represent each pixel. How can Henry use that information to calculate the approximate size of the file?

Note: This calculation would not get you the exact file size since files also include metadata (data about the file (e.g., header information in a bitmap image file).

How does the blur filter work?

## CPSC 100 2016W2: Practice Midterm
## Apriori

Show all the steps to finding the frequent itemsets with >50% support.

| Transaction | Items |
|:---:|:---|
| T1 | Coffee, Tea, Juice, Water |
| T2 | Tea, Juice, |
| T3 | Coffee, Juice |

Frequent itemset of size 1
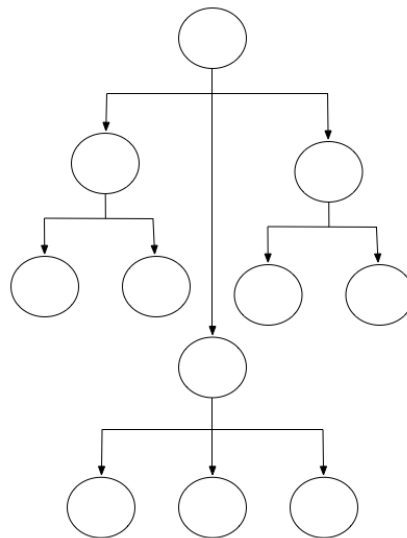
Frequent itemset of size 2

Frequent itemset of size 3

**Frequent itemsets with > 50% support**

# CPSC 100 2016W2: Practice Midterm

Explain how the itemsets with ≥50% support would change if we added T4: Lemonade in the original table.

## Decision Trees

Based on the following diagram, determine how many of the following there are:



- Nodes:
- Edges:
- Leaves:
- Depth:
- Parent:

## CPSC 100 2016W2: Practice Midterm

You are procrastinating from studying for your CPSC 100 midterm and you are trying to decide whether you should study, so you make a decision tree to help you determine if you should study for each of the five chapters of the textbook.

| Did I read the chapter? | How well do I understand the text? | How long will it take me to review? | What impact will the chapter have on the exam? | How difficult are the questions? | Should I study? |
|---|---|---|---|---|---|
| Yes | Confident | Long | Significant | Hard | No |
| No | Fairly well | Short | Trivial | Medium | Yes |
| Yes | Confused | Long | Significant | Medium | No |
| Yes | Confused | Medium | Significant | Hard | No |
| No | Confident | Medium | Trivial | Easy | Yes |

For each attribute (i.e., Read Chapter or Not, Undersanding of Text, …, Question Difficulty), what is the overall entropy if we split on that attribute?

| Read Chapter or Not | Understanding of Text | Review Time | Chapter Impact | Question Difficulty |
|---|---|---|---|---|
|  |  |  |  |  |

Draw the tree(s) that split on attributes with the greatest reduction in overall entropy.
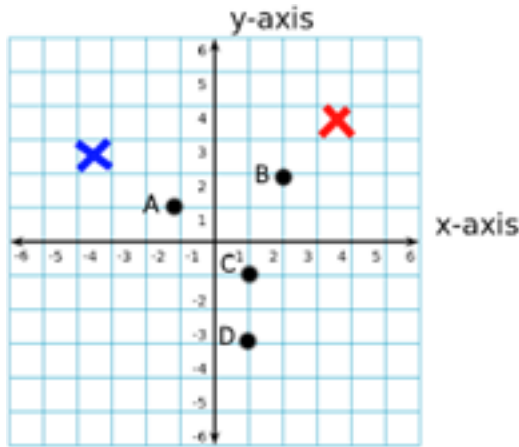
## Clustering

Given the following 9 items, how would you group these items (what measure of quality would you use)? How many data points and what are the data dimensions based on what how you clustered these items?



What are the benefits of clustering and how does it help in data mining?

## CPSC 100 2016W2: Practice Midterm

Below is a graph with data points A, B, C, and D and two centroids (blue and red X's).



What will happen in the following 2 steps (cluster assignment and move centroid) in the k-means clustering algorithm, given the following table?

| Data point | Distance to Red Centroid (4,4) | Distance to Blue Centroid (-4,3) |
|---|---|---|
| A (-1,1) | 5.8 | 3.6 |
| B (2,2) | 2.8 | 6.1 |
| C (1,-1) | 5.8 | 6.4 |
| D (1,-3) | 7.6 | 7.8 |

1. **Cluster assignment:** Which clusters will these data points be assigned to?

2. **Move centroid**: Given the following data, what coordinates will the red and blue centroid move to? State which of the following calculations should be performed (average or median and for which points. Also, state which of the following will determine where the red centroid, and blue centroid will move to.

| Calculation | New Centroid |
|---|---|
| 1) Average of points A & D | **(0,-1)**<br>x-coord = (-1 +1)/2 = **0**<br>y-coord = (1+-3)/2 = **-1** |
| 2) Median of points A & D | **(0,-1)**<br>X-coord = -1, 1, even number so take the average (-1+1)/2 = **0**<br>y-coord = -3,1. (-3+1)/2 = **-1** |
| 3) Average of points B & C & D | **(1.3,-0.7)**<br>x-coord = (2+2+1)/3 = **1.3**<br>y-coord = (2+-1+-3)/3 =**-0.67** |
| 4) Median of points B & C & D | **(1,-1)**<br>x-coord = 1,1,2 median is **1**<br>y-coord = -3,-1,2, median is **-1** |
| 5) Average of points B & C | **(1.5,0.5)**<br>x-coord = (2+1)/2 = **1.5**<br>y-coord = (2+-1)/2 = **0.5** |
| 6) Median of points B & C | **(1.5,0.5)**<br>X-coord = 1,2 even number so take the average (1+2)/2 = **1.5**<br>y-coord = -1,2. (-1+2)/2 = **0.5** |
| 7) Average of point A | **(-1,1)** |
| 8) Median of point A | **(-1,1)** |

The two calculations used should be numbers _____ and _____.
Red centroid will move to _____.
Blue centroid will move to _____.

# Decimal, binary, and hex conversion table
(leading 0's, shown in gray, are useful for some conversions)

| dec | bin | hex |
|-----|------|-----|
| 00  | 0000 | 0   |
| 01  | 0001 | 1   |
| 02  | 0010 | 2   |
| 03  | 0011 | 3   |
| 04  | 0100 | 4   |
| 05  | 0101 | 5   |
| 06  | 0110 | 6   |
| 07  | 0111 | 7   |

| dec | bin  | hex |
|-----|------|-----|
| 08  | 1000 | 8   |
| 09  | 1001 | 9   |
| 10  | 1010 | A   |
| 11  | 1011 | B   |
| 12  | 1100 | C   |
| 13  | 1101 | D   |
| 14  | 1110 | E   |
| 15  | 1111 | F   |