

**CS 4210 – Assignment #1 Completed
Maximum Points: 100 pts.**

Bronco ID: 013141414

Last Name: Mehrani

First Name: Arsham

Note 1: Your submission header must have the format as shown in the above-enclosed rounded rectangle.

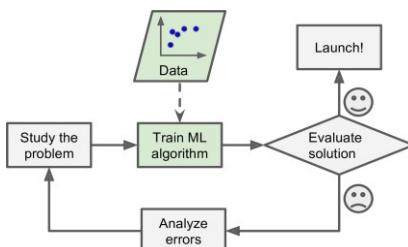
Note 2: Homework is to be done individually. You may discuss the homework problems with your fellow students, but you are NOT allowed to copy – either in part or in whole – anyone else’s answers.

Note 3: Your deliverable should be a .pdf file submitted through Gradescope until the deadline. Do not forget to assign a page to each of your answers when making a submission. In addition, source code (.py files) should be added to an online repository (e.g., github) to be downloaded and executed later.

Note 4: All submitted materials must be legible. Figures/diagrams must have good quality.

Note 5: Please use and check the Canvas discussion for further instructions, questions, answers, and hints.

1. [6 points] A computer program is said to learn from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E (Mitchell, 1997). Explain this definition of a machine learning system including in your answer details about E, T, P and how you **correlate** them with the **components** of the image below.



Machine Learning is defined as a programs ability to learn without being explicitly programmed. This is done through experience/training and is measured by a performance test designed specific to the algorithm.

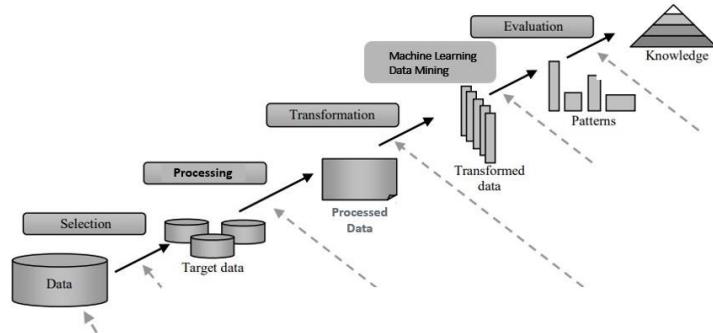
T = Task that needs to be done using Machine Learning

E = The training/experience data that the algorithm is going to use to improve

P = Performance testing using a portion of the same data set.

Study the problem: in this step the task T is analyzed to build solution. *Train ML algorithm* in this step the experience E is introduced to the model for training purposes. *Evaluate solution*, in this step the performance P is measured using the testing data set. *Analyze errors* This is part of the Task T stage, since in this step The algorithm is analyzed for errors to find the root of the problem. Once the problems are identified in last step the cycle continues to study these errors and training the model further. Once the performance P is satisfactory the algorithm is ready to launch.

2. [6 points] Some authors present a machine learning/data mining pipeline process with only 3 main phases instead of those 6 shown in the image below (see the dashed arrows). **Explain** the reasons why they probably decided to do that including in your answer **what** are those 3 main phases and their **corresponding relevance** to build knowledge.



The pipeline process is summarized into 3 main steps as described below. This is because steps 1-4 all have to do with the preparation of data to well suit the algorithm.

Steps:

1. *Preprocessing(1-4):*

This is probably the most time-consuming part of the ML process. The data selection is key to a good result. If the data is not optimized (i.e. features selection, dimensionality reduction, etc.) then the results that come out are prone to errors. Therefore, a clear data set that is optimized for a certain ML algorithm is essential to accurate results.

2. *Machine Learning(5):*

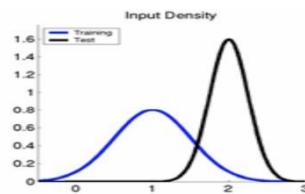
This step is where the ML algorithm is written to take in the training data. This algorithm needs to be able to use the processed data and make improvements to its output over time. The data produced by the end of this step will be used to build knowledge.

3. *Postprocessing(6):*

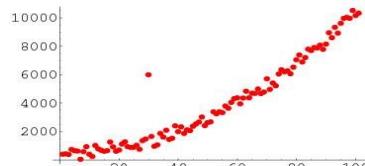
This step is where the data gathered from the ML algorithm is visualized and optimized for better understanding of the results. i.e. making graphs and models that easily describe the experiment. Other devices of learning such as logic (ex. If A then B, C is direct cause of D) are used as well to better paint the picture.

3. [10 points] Machine learning algorithms face multiple challenges while analyzing data, such as scalability, dimensionality, data distribution, sparsity, resolution, noise, outliers, missing values, and duplicated data. For each image below, **name** and **explain** what the corresponding challenge is from this list (you do not need to explain how to solve the challenge).

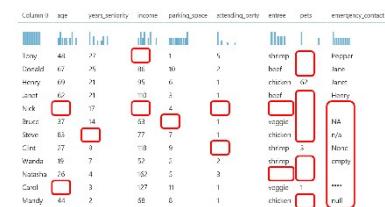
a.



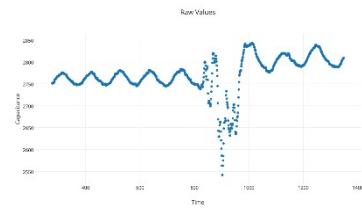
b.



c.



d.



e.

c1	c2	c3	c4	c5
0	0	0	5	0
2	0	0	0	0
0	0	1	0	0
0	5	0	0	1
3	0	0	3	0
0	4	0	0	0

- Data Distribution.* In this graph we see two distributions for test and training that have zero to no correlation.
- Outlier.* This graph demonstrates all the points aligned in a quadratic fashion while one point is way outside the acceptable boundaries. This is an example of an outlier or anomaly. Could happen because of an error or because of that certain test case.
- Missing values.* This table clearly is lacking some values in almost all its samples.
- Noise.* In this graph we see that during a certain period the experiment returned bad results that do not follow the logic of the rest of the experiment. This is likely because of a certain error that happened during that time that tampered with the experiment.
- Sparsity.* which is are many zero values that do not contribute to the experiment

4. [12 points – 2 points each] Analyze the dataset below and answer the proposed questions:

The Contact Lens Data

Age	Spectacle Prescription	Astigmatism	Tear Production Rate	Recommended Lenses
Young	Myope	No	Reduced	No
Presbyopic	Myope	No	Normal	No
Prepresbyopic	Myope	No	Reduced	No
Prepresbyopic	Myope	No	Normal	Yes
Presbyopic	Myope	Yes	Normal	Yes
Young	Myope	Yes	Normal	Yes
Young	Hypermetrope	No	Reduced	No
Prepresbyopic	Myope	Yes	Reduced	No
Presbyopic	Hypermetrope	No	Reduced	No
Young	Myope	Yes	Reduced	Yes

- a. *What is the most likely task that data scientists are trying to accomplish?*
The most likely task is to try and predict a nearsightedness in people with underlying conditions or based on prescription.

- b. *In general, what is a feature and how would you exemplify it with this data?*
Feature is a quality or a variable that are used to perform the task. In this case features include “Age”, “Spectacle Prescription”, “astigmatism” , etc. features are all the columns in a data set.

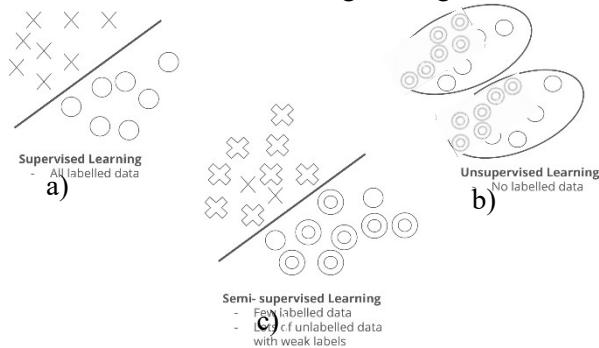
- c. *In general, what is a feature value and how would you exemplify it with this data?*
Each feature has a measurable quality. And in every sample these qualities are measured and noted in their designated area. For example here first row third column has a feature value of “No” which is essentially the measurement of that particular feature for this particular sample.

- d. *In general, what is dimensionality and how would you exemplify it with this data?*
Dimensionality is referred to the number of features each sample in a data set contains. Basically this is the number of columns a data set has. In this case the dimensionality is 5, because it has 5 attributes or features.

- e. *In general, what is an instance and how would you exemplify it with this data?*
A sample or an instance is each row of the data set. Each row signifies a particular case that was examined. In this case each of these samples happen to be a person. For each sample all the features are measured, and the results are recorded.

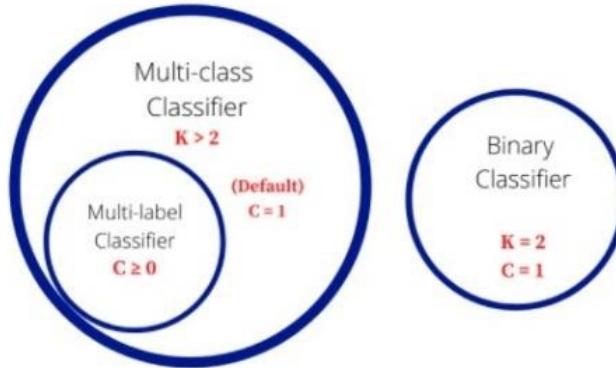
- f. *In general, what is a class and how would you exemplify it with this data?*
Class/labels are used for classification of each sample. In this example we don't see the class vector.

5. [6 points] Identify and explain what **kind of machine learning** (supervised, unsupervised, semi supervised, reinforcement) **system** should be used for each scenario below including in your answer **data labels** information. Hint: check the images to figure out which data sample is labelled.



- supervised learning*, because some of the elements are known and others are undefined and could be used as test cases
- unsupervised*, elements are not labeled, and algorithm is working to identify clusters
- semi-supervised* because some of the figures are defined and others are not (X marks)

6. [6 points] Explain the **tasks** addressed by each classifier below.



K = Total number of classes in the problem statement

C = Number of classes an item maybe assigned to

Binary Classifier:

This classifier has 2 possible outcomes/classes, and each item can only be assigned to one of these not both nor 0 classes.

Multi-class Classifier:

This classifier has **strictly more** than 2 possible outcomes/classes and by default each item is assigned **one**.

Multi-Label Classifier:

This classifier has **strictly more** than two outcomes/classes, but the items can be assigned 0 or more of these classes.

7. [54 points] Regarding the training data shown in question 4:
- [12 points] Find a **maximally specific hypothesis** that fits the data **following** the strategy of **Find-S algorithm**. Hint: start from the hypothesis $h \leftarrow (\emptyset, \emptyset, \emptyset, \emptyset)$.

My solution is on next page.

- [12 points] Complete the given python program (find_s.py) that will read the file contact_lens.csv and output the hypothesis of **Find-S algorithm** (the hypothesis you got in part a). The output should be in this format: ['Sunny', '?', 'Strong', '?']. Add the link to the online repository as the answer to this question.

The link: https://github.com/Arsham1024/Introduction_ML

- [15 points] Derive the decision tree produced by the standard ID3 algorithm. Show your calculations for entropy and information gain for all splits. Plot your final tree at the end.

My solution is on next page.

- [15 points] Complete the given python program (decision_tree.py) that will read the file contact_lens.csv and output the decision tree of **ID3** (the tree you got in part c). PS: if this tree is different from yours shown in letter c), try to explain why. Add the link to the online repository as the answer to this question.

The link : https://github.com/Arsham1024/Introduction_ML

Important Note: Answers to all questions should be written clearly, concisely, and unmistakably delineated. You may resubmit multiple times until the deadline (the last submission will be considered).

NO LATE ASSIGNMENTS WILL BE ACCEPTED. ALWAYS SUBMIT WHATEVER YOU HAVE COMPLETED FOR PARTIAL CREDIT BEFORE THE DEADLINE!

Question 7

	1	2	3	4	class
	Age	Spectacle Prescription	Astigmatism	Tear Production Rate	Recommended Lenses
1	Young	Myope	No	Reduced	No
2	Presbyopic	Myope	No	Normal	No
3	Presbyopic	Myope	No	Reduced	No
4	Presbyopic	Myope	No	Normal	Yes
5	Presbyopic	Myope	Yes	Normal	Yes
6	Young	Myope	Yes	Normal	Yes
7	Young	Hypermetrope	No	Reduced	No
8	Presbyopic	Myope	Yes	Reduced	No
9	Presbyopic	Hypermetrope	No	Reduced	No
10	Young	Myope	Yes	Reduced	Yes

Part A) ① hypothesis $(\emptyset, \emptyset, \emptyset, \emptyset)$

② Loop through all (+) samples: Samples: 4, 5, 6, 10

Sample ④ $h(\text{presbyopic}, \text{Myope}) \text{ NO, Normal}$

Sample ⑤ $h(\text{?}, \text{Myope}) \text{ ? , Normal}$

Sample ⑥ $h(\text{?}, \text{Myope}) \text{ ? , Normal}$

Sample ⑩ $h(\text{?}, \text{Myope}) \text{ ? , ? }$

③ Return $h(\text{?}, \text{Myope}) \text{ ? , ? }$

Part C)

$$\text{Entropy}(s) = -p_e \lg_2 P_e - p_o \lg_2 P_o$$

find $E(s)$:

$$S = 10 \Rightarrow E(s) = -\frac{4}{10} \frac{\lg 0.4}{\lg 2} - \frac{6}{10} \frac{\lg 0.6}{\lg 2}$$

$$E(s) = 0.971$$

Age:

$$\text{Gain}(S, A) = \text{Entropy}(s) - \sum \frac{|S_i|}{|s|} \text{Entropy}(S_i)$$

$$\begin{aligned} \text{Gain}(S, \text{Age}) &= 0.971 + \left[-S_{\text{young}} - S_{\text{pres}} - S_{\text{prepre}} \right] \\ &\downarrow \\ &= 0.971 + \left[-0.4(1) - 0.3(0.971) - 0.3(0.971) \right] \end{aligned}$$

$$\boxed{\text{Gain}(S, \text{Age}) = 0.0202}$$

Question 7:

a) strategy of finding S algorithm $\rightarrow (\emptyset, \emptyset, \emptyset, \emptyset)$

b) Complete python program final S.py
with respect contact lenses

- Output [young, ?, strong, ?, ?]

→ will output same as @

c) Derive the Decision tree

- ID3 algorithm → for all splits?

- Show calculation for Entropy & info Gain

- plot the final tree!

d) Complete Decision-tree.py

→ Read save Contact-Lens.csv

- Output the tree

- if Different explain why?

from c)

Spectacle Prescription (sp)

$$\begin{aligned} \text{Gain}(S, sp) &= 0.971 + \left[-S_{\text{myope}} - S_{\text{hyper}} \right] \\ &= 0.971 + \left[-\frac{8}{10}(1) - 0 \right] \end{aligned}$$

$$\boxed{\text{Gain}(S, sp) = 0.171}$$

$$S_{\text{young}} [2+, 2-]$$

$$S_{\text{presbyopic}} [1+, 2-] \rightarrow \frac{3}{10} - \frac{1}{3} \left(\frac{\lg \frac{1}{3}}{\lg 2} \right) - \frac{2}{3} \left(\frac{\lg \frac{2}{3}}{\lg 2} \right)$$

$$S_{\text{presbyopic}} [1+, 2-] \rightarrow \frac{3}{10}$$

$$S_{\text{myope}} [4+, 4-] \rightarrow E = 1$$

$$S_{\text{Hyper}} [0+, 2-] \rightarrow E = 0$$

Astigmatism:

$$\text{Gain}(S, As) = 0.971 + \left[-S_{yes} - S_{no} \right]$$

$$= 0.971 + \left[-(0.4)(0.811) - (0.6)(0.650) \right]$$

$$\boxed{\text{Gain}(S, As) = 0.2566}$$

$$S_{yes} [3+, 1-] \rightarrow 0.4$$

$$S_{no} [1+, 5-] \rightarrow 0.6$$

$$- \frac{1}{6} \frac{\log \frac{1}{6}}{\log 2} - \frac{5}{6} \frac{\log \frac{5}{6}}{\log 2}$$

Tear Production Rate

$$\text{Gain}(S, TPR) = 0.971 + \left[-S_{Reduced} - S_{Normal} \right]$$

$$= 0.971 + \left[-0.6(0.650) - 0.4(0.811) \right]$$

$$S_{Reduced} [1+, 5-] \xrightarrow{\text{Done this before}}$$

$$S_{Normal} [3+, 1-]$$

$$\boxed{\text{Gain}(S, TPR) = 0.2566}$$

true so far:

$$\boxed{\text{Gain}(S, Age) = 0.0202} \quad (4)$$

$$\boxed{\text{Gain}(S, SP) = 0.171} \quad (3)$$

$$\boxed{\text{Gain}(S, As) = 0.2566} \quad (2)$$

$$\boxed{\text{Gain}(S, TPR) = 0.2566} \quad (1)$$

TPR

Normal

$$\{2, 4, 5, 6\}$$

$$[3+, 1-]$$

Reduced

$$\{1, 3, 7, 8, 9, 10\}$$

$$[1+, 5-]$$

same but I choose

Tear Production Rate!

Moving to second level:

Normal:

$$S_{=4} \Rightarrow E(S) = -\frac{3}{4} \frac{\log \frac{3}{4}}{\log 2} - \frac{1}{4} \frac{\log \frac{1}{4}}{\log 2}$$

$$\boxed{E(S) = 0.811}$$

Age:

$$\text{Gain}(S, Age) = 0.811 + \left[-S_{young}^0 - S_{pre}^0 - S_{prepre}^0 \right]$$

$$= 0.811 - (0.5)$$

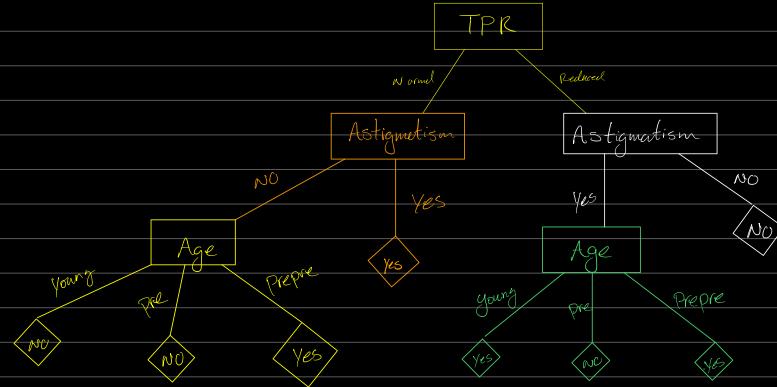
$$\boxed{\text{Gain}(S, Age) = 0.311}$$

$$S_{young} [1+, 0-] \rightarrow E=0$$

$$S_{pre} [1+, 1-] \rightarrow E=1$$

$$S_{prepre} [1+, 0-] \rightarrow 0$$

Final Tree



Age	Spectacle Prescription	Astigmatism	Tear Production Rate	Recommended Lenses
Young	Myope	No	Reduced	No
Presbyopic	Myope	No	Normal	No
Prepresbyopic	Myope	No	Reduced	No
Prepresbyopic	Myope	No	Normal	Yes
Presbyopic	Myope	Yes	Normal	Yes
Young	Myope	Yes	Normal	Yes
Young	Hypermetrope	No	Reduced	No
Prepresbyopic	Myope	Yes	Reduced	No
Presbyopic	Hypermetrope	No	Reduced	No
Young	Myope	Yes	Reduced	Yes