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

Bronco ID:  015412449

Last Name: Langgeng

First Name: Leonardo

**Note 1:** Your submission header must have the format as shown in the above-enclosed rounded rectangle.
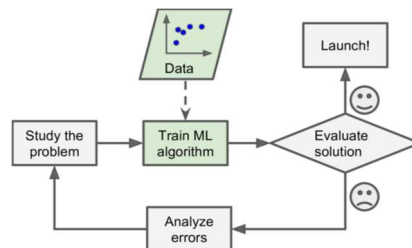**Note 2:** Homework is to be done individually.  You may discuss the homework problems with your fellow students, but you are NOT allowed to copy – either in part or in whole – anyone else's answers.
**Note 3:** Your deliverable should be a .pdf file submitted through Gradescope until the deadline. 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. The bold words/sentences provide information for a complete or accurate answer.

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 informing in your answer how **E, T, P correlate** with **each component** of the image below.



Experience are correlated with the data being fed to the ML algorithm. The more data being fed to the ML for training, the more experience the machine will learn from, and assuming there is no problem with the input data (bias, noise, etc), the machine will become more accurate.

Task are correlated with the objective of the ML. The machine will try to predict the output from a given input, based on the experience from training.
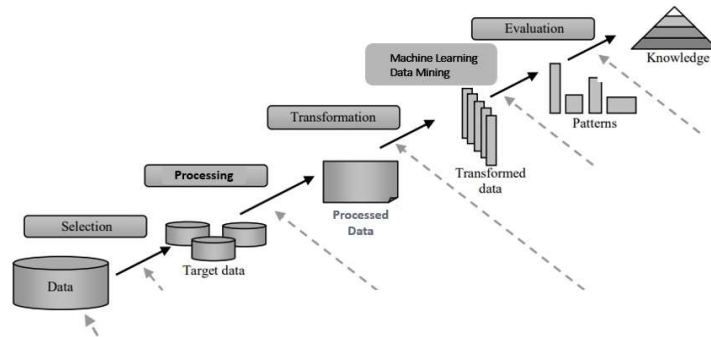
Performance are the accuracy of the machine's prediction. It either could be a % of error prediction.

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).  **Name** those 3 main phases and **explain** their corresponding relevance to build knowledge.

Preprocessing: preparing the dataset to be processed. This phase cleans up the dataset and making sure that the data are as optimal as possible thus making the ML process to be accurate with.
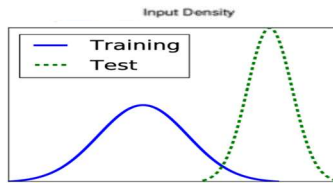
Machine learning: process where the machine will analyze the data using certain algorithm and outputs certain patterns to be analyzed later.

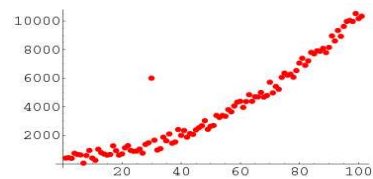Postprocessing: Analyzing and interpret the information produced by ML



3. [15 points – 3 points each] Machine learning algorithms face multiple challenges while analyzing data such as scalability, distribution, sparsity, resolution, class imbalance, 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).
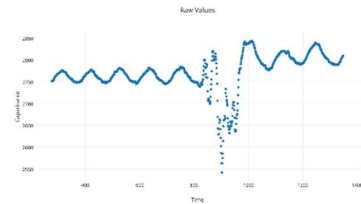
a.



b.



c.



d.



e.

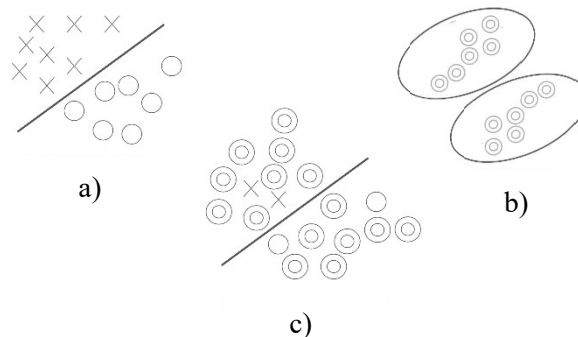| c1 | c2 | c3 | c4 | c5 |
|----|----|----|----|----|
| 0  | 0  | 0  | 5  | 0  |
| 2  | 0  | 0  | 0  | 0  |
| 0  | 0  | 1  | 0  | 0  |
| 0  | 5  | 0  | 0  | 1  |
| 3  | 0  | 0  | 3  | 0  |
| 0  | 4  | 0  | 0  | 0  |

A. Biased training/test sets. If the training and testing sets are not balanced/equal, the prediction from it will be biased toward the training set, not the expected result from test set.

B. Outliers. Depending on algorithm used, it may throw off the model produced by the machine as outliers do not line up with the rest of the data points. Outliers are real data point, thus it should still be considered/included in the dataset.

C. Missing values. It makes the machine harder to plot the data points with missing values, and in some/most cases we need to replace the missing values.

D. Noise. It puts a lot of inaccurate data points and thus will throw off machines prediction.

E. Sparse data. If the data is sparse, it is very hard for the machine to find all of the possible combinations of the features to build the model.

4. [18 points – 3 points each] Analyze the dataset below and answer the proposed questions:
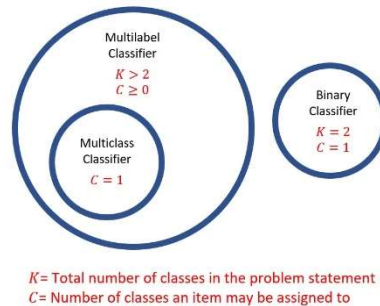
The Contact Lens Data

| 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?
To predict whether or not we recommend lenses for certain people/patient based on certain information

b. **In general**, what is a feature and how would you **exemplify** it with **this data**?
Feature are the known trait from each data point. In this case, it is age, spectacle prescription, astigmatism, and tear rate

c. **In general**, what is a feature value and how would you **exemplify** it with **this data**?
The value of each feature of each data point. Age could be young, presbyopic, and prepresbyopic

d. **In general**, what is dimensionality and how would you **exemplify** it with **this data**?
The number of features from a dataset. The table have 4 datasets as we using recommended lenses as class/result

e. **In general**, what is an instance and how would you **exemplify** it with **this data**?
Instance are the data sample. The table have 10 instance (10 people/patient)

f. **In general**, what is a class and how would you **exemplify** it with **this data**?
Class is the expected result from each instance. In this case, it is whether or not we recommend lenses to the instance/patient

5. [9 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 information about **data labels**. Hint: check the images to figure out which data sample is labelled.



a)

b)

c)

A. Supervised learning. It is a classification tasks, and we try to predict the cutoff line / the divider between 2 or more class.

B. Unsupervised learning. The image trying to find interesting patterns or grouping in the data points.
C. Semi supervised learning. It is trying to map the data of a lot of unlabeled data points with a bit of labeled points.

6. [9 points] Explain the **tasks** addressed by each classifier below.



$K$ = Total number of classes in the problem statement
$C$ = Number of classes an item may be assigned to

Binary classifier guessing single output from 2 possible class (binary output; yes or no)

Multilabel classifier allows multiple class output from multiple possible class (result could be a percentage of each possible class)
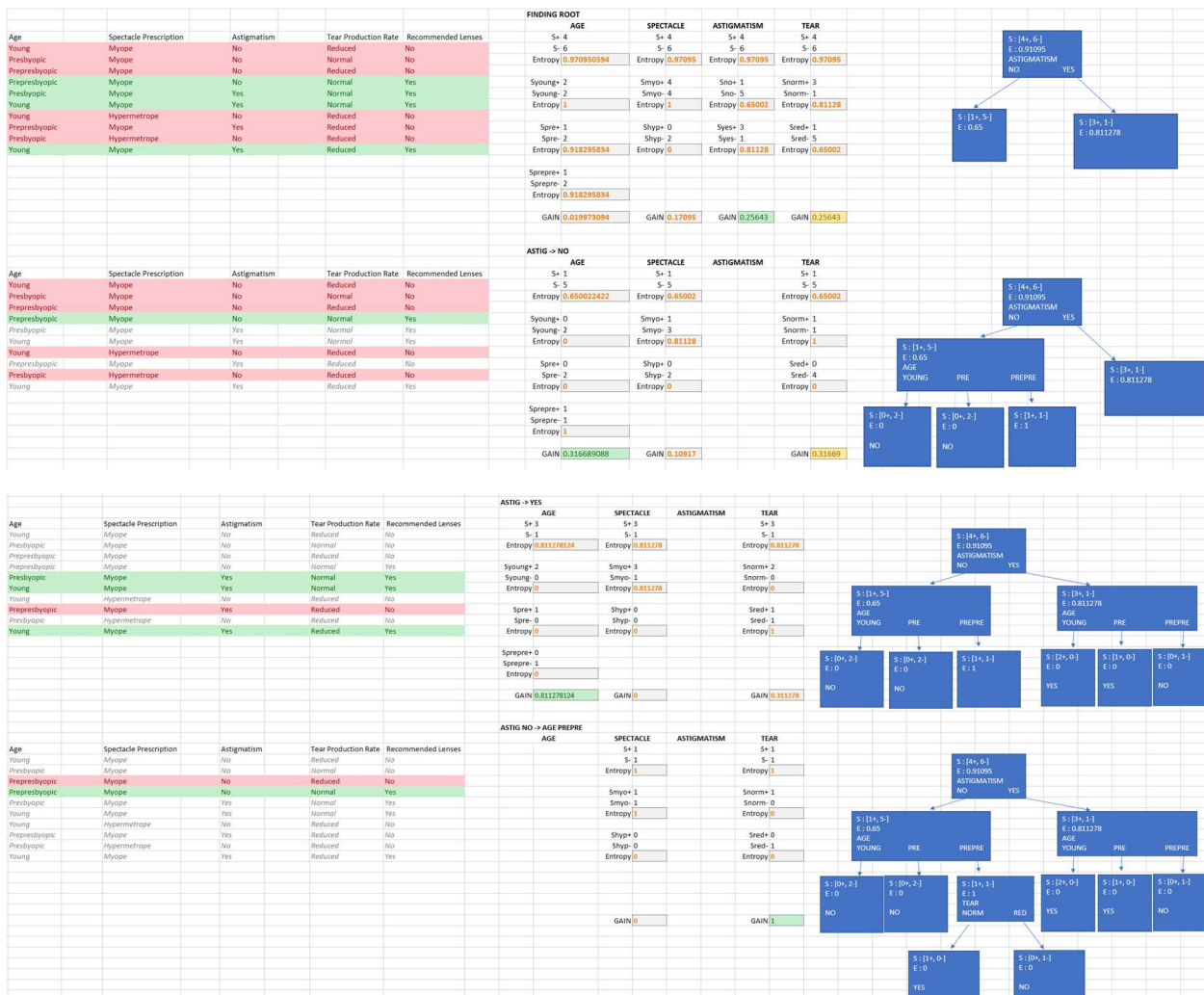
Multiclass classifier produce 1 output class from multiple possible class (there is n# options, chose 1)

7. [37 points] Regarding the training data shown in question 4:

    a. [20 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. Excel attachment

    b. [15 points] Complete the given python program (decision_tree.py) that will read the file contact_lens.csv and output a decision tree. Add the link to the online repository as the answer to this question.
https://github.com/leolanggeng/contact_lens_ML

    c. [2 points] The tree you got in part b) should be the same one you got in part a), but there are probably some differences. Try to explain why.
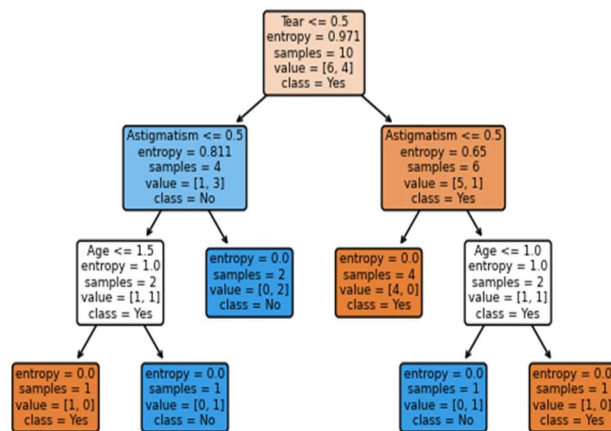
While finding the root, astigmatism and tear have the same gain score, and also while finding the leaf of astigmatism -> age and tear also have the same gain. From this alone, there are at least 3 other possible trees, 1 the tree that I made and 2 alternate path to create the tree. This shows as the first time I ran the code I produced a different tree, but by repeating the code I do get the same tree that I made. There are also the case of the binary tree, the code cannot handle multi branch options so the 3 possibilities of age are combined.

**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!**