

# STA130 Winter 2020

(Materials used in this presentation are provided by the U of T Statistical Sciences Department.

This presentation was prepared by Vivian Ngo.)

**[Github.com/vivianngo97/STA130-Winter-2020](https://github.com/vivianngo97/STA130-Winter-2020)**

**[viv.ngo@mail.utoronto.ca](mailto:viv.ngo@mail.utoronto.ca)**

# Agenda

- Reminders
- Vocabulary
- Group discussion
- Group work & presentations
- Let me know who your groups are \*\*\*

# Reminders

- We have FREE TA and office hours!

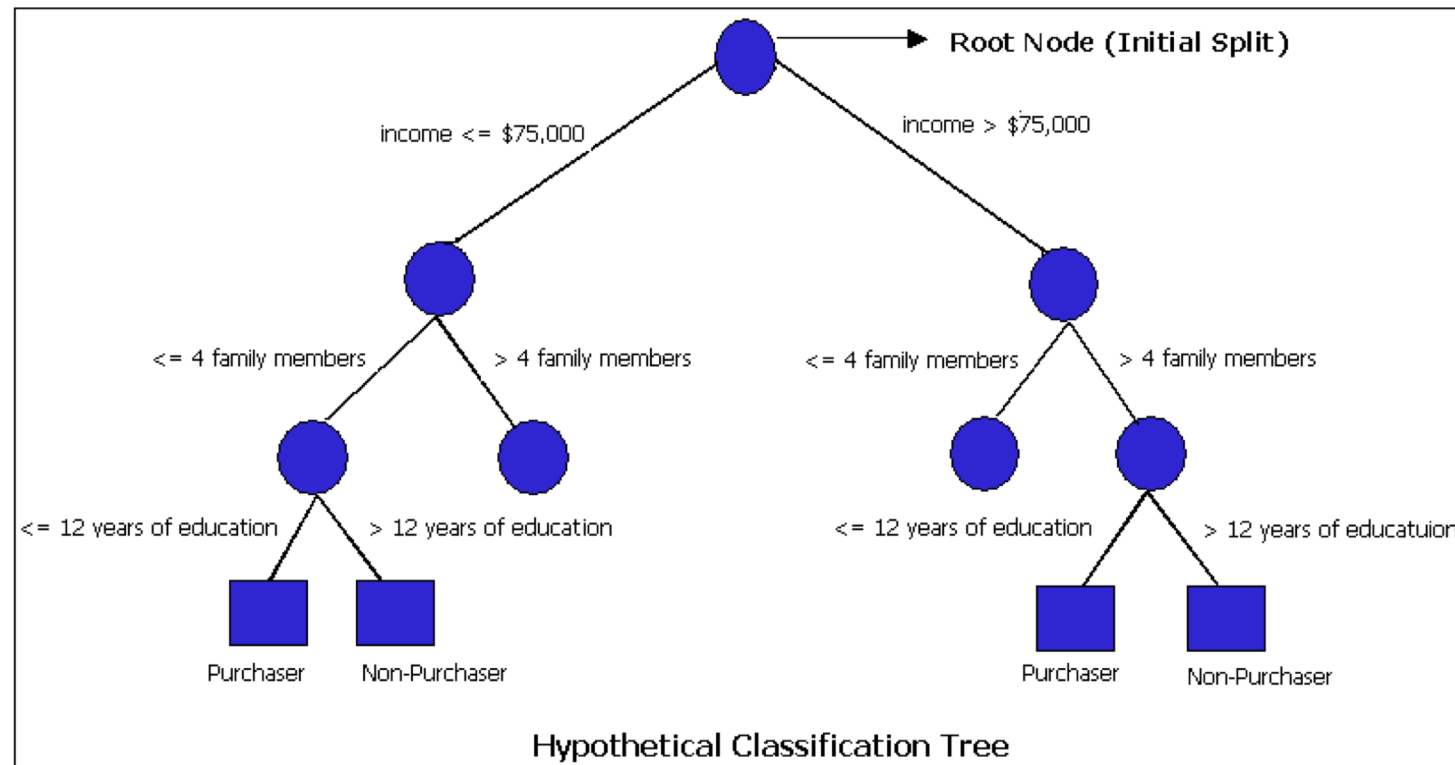


# Vocabulary

- Classification
- Receiver operating characteristic (ROC) curve
- Prediction
- Predictor(s)
- Covariate(s)
- Independent variable(s)
- Dependent variable(s)
- Input(s)
- Output(s)
- Training set/sample
- Validation
- Validation set/sample (or hold-out set or test set)
- Fitting a model
- Confusion matrix
- Category
- Tree
- Terminal node
- Stopping rule
- Threshold
- True positive (sensitivity)
- True negative (specificity)
- False positive
- False negative
- Accuracy
- Classifier
- Cutpoint
- Node(s)
- Terminal Node
- Binary
- Split(ting)

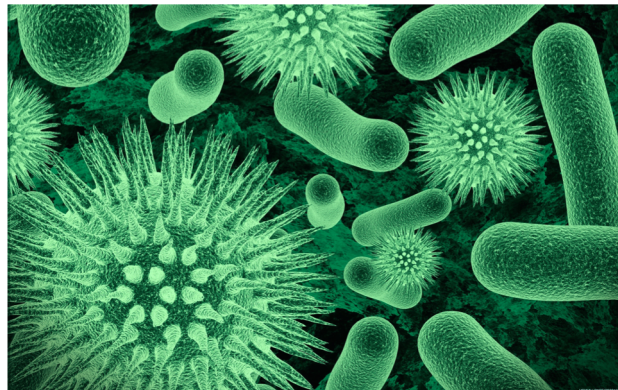
# Group discussion

- 1. Can you think of any real-life examples where you may want to develop a classification tree?



# Group discussion

- 1. Can you think of any real-life examples where you may want to develop a classification tree?
- 2. [From Practice Problem Question 2] Suppose you developed a classification tree to diagnosis whether or not somebody has Disease X, which is a very serious and life-threatening illness if left untreated. The overall accuracy of your tree was 77%; false-positive rate was 32%; and false-negative rate: 7.9%.



# Group discussion

- 1. Can you think of any real-life examples where you may want to develop a classification tree?
- 2. [From Practice Problem Question 2] Suppose you developed a classification tree to diagnosis whether or not somebody has Disease X, which is a very serious and life-threatening illness if left untreated. The overall accuracy of your tree was 77%; false-positive rate was 32%; and false-negative rate: 7.9%.
- Suppose that your colleague also created a classifier for the same purpose. Its overall accuracy is 81%; false-positive rate is 6.4%; and false-negative rate is 39%. Explain which of these two classifiers you would prefer to use to diagnosis Disease X.

# Group discussion

- 3. Consider the same 2 classifiers for Disease X, but now suppose the treatment is very expensive and has many bad side effects; e.g. people taking the treatment tend to get very sick, similar to chemotherapy. In this case which classifier would you prefer?



# Group discussion

- 3. Consider the same 2 classifiers for Disease X, but now suppose the treatment is very expensive and has many bad side effects; e.g. people taking the treatment tend to get very sick, similar to chemotherapy. In this case which classifier would you prefer?
- **Note:** There are many ways to measure the performance of your classifier and the context for which you are using the classifier matters: accuracy vs false negatives vs false positives

# Group discussion

- 4. Suppose you developed a classification tree only to later discover that the values for one of your covariates is missing for a number of observations. Can you use the classification tree you built to make a prediction for these individuals? Explain. [See Practice Problem 1e for example]

# Group discussion

- 5. Imagine you were interested in making a classifier to predict what movie somebody would be most interested in. To do this, you first gathered data from a sample of your closest friends. You validated and tested your classifier using different subsets of this data. Now you wish you use your classifier to predict which movie Dr. Moon/ Bolton, your TA, your parents, etc. would like. How well do you think your classifier will perform in each of these cases?

# Extra questions to consider

- What is accuracy? Is accuracy always a good metric?
- What is the confusion table? What are the cells in this table?
- What is the ROC curve? When to use it?
- What is a decision/ classification tree? How do we train them?
- How do we handle categorical variables in decision trees? Continuous variables?

# Oral Presentations

- **In poster project groups**
- **Remember:**
- **THE 4 C'S:** Calm; Confident; Clear; Concise
- ***Tips for giving a great oral presentation: Content***
  - What is the main message you want to get across?
  - Create an (organized) outline of your presentation
  - Define terms early
  - Make clear transitions between parts of your presentation
- Make your data/ figures meaningful
- Summarize
- ***Tips for giving a great oral presentation: Delivery***
  - Be confident, make eye contact and avoid reading
  - Avoid filler words – “ummm”, “like”, “you know”
  - Speak slowly and it's ok to pause (and breathe!)
  - Remember to enunciate all the parts of each word
  - Practice! Practice! Practice!

# Oral presentations

- **Topic one:**

- Explain how to make a ROC curve and the type of information it provides.
- Based on the ROC curves you created for Practice Problem 4c, describe the accuracy of each of the two trees.
- Does this fit your expectations based on the description of each classifier?
- Which ROC curve would you prefer to classify your spam mail?

- **Topic two:**

- Explain what a confusion matrix is and how each cell is calculated.
- Using the confusion matrix you calculated in question 1d to answer the following questions: What percentage of countries with “good life expectancy” that were classified as having such

actually had “good” life expectancy according to the majority rules cutpoint (i.e., 50%) based on each of the two classifiers?

- What are other terms used to describe the percentages you calculated above?
- How do the two classifiers compare? Does this fit your expectations based on the description of each classifier?

- **Topic three:**

- Summarize the classification tree from Practice Problem 1b. Make sure to include *at least* the following points: how the splits on each variable were selected, how a new observation would be predicted by this classification tree.
- In part c, you considered more factors. Do you think there may be other important factors to consider? Explain how including these might impact the accuracy of your tree.

# Group information

Group #	Member 1 (First Last Names)	Member 2	Member 3	Member 4
1				
2				
3				
4				
5				