

In this mini-project, you will implement a decision-tree algorithm and apply it to molecular biology, a Promoter is a region of DNA that facilitates the transcription of a particular gene. Promoter compilations and analyses have led to computer programs which predict the location of promoter sequences on the basis of homology either to the consensus sequence or to a reference list of promoters. Such programs are of practical significance in searching new sequences.

The task for you is to develop a decision tree algorithm, learn from data, and predict for unseen DNA sequences whether they are promoters or non-promoters.

Our data was provided by UCI Machine Learning Repository and can be found in Molecular Biology (Promoter Gene Sequences) Data Set. It has 106 instances and 57 features. We randomly split the data set into the training (71 instances) and validation (35 instances) sets, which are both well balanced.

1. Download the data set which contains the training data and validation sets. Each DNA sequence is represented by one line, with the first 57 characters (one of 'a', 'g', 'c' and 't') representing the sequence, and the last character (separated by a space from the sequence) indicating the class ('+' for promoter, '-' for non-promoter).
2. Implement the ID3 decision tree learner, as described in Chapter 3 of Mitchell or in your Web link <http://www.cise.ufl.edu/~ddd/cap6635/Fall-97/Short-papers/2.htm>. You may program in any language that you prefer. Your program should assume input in the above format. For initial debugging, it is recommended that you construct a very simple data set (e.g., based on a boolean formula) and test your program on it.
3. Implement both accuracy (misclassification impurity) and information gain (entropy impurity) for evaluation criterion. Also, implement split stopping using chi-square test. The test-statistic is:

$$\sum_{i=1}^v \frac{(p_i - p'_i)^2}{p'_i} + \frac{(n_i - n'_i)^2}{n'_i}$$

where p_i is the number of promoters in the i -th child, p'_i is the expected number of promoters in the i -th child, n_i is the number of non-promoters in the i -th child and n'_i is the expected number of non-promoters in the i -th child. The degree of freedom is $v-1$

4. Use your algorithm to train a decision tree classifier using the training data and report accuracy on the validation data. Compare accuracies by varying the evaluation criteria and confidence level in determining split stopping. For the latter, use 99%, 95% and 0% (i.e., you always grow the full tree).
5. **Turn in the following:**
 - o Your code. Submit your file through UNM Learn. Your code should contain appropriate comments to facilitate understanding. If appropriate, your code must contain a Makefile or an executable script that receives the paths to the training and testing files and a README file
 - o A report of at most 4 pages (letter size, 1 inch margins, 12pt font) that describes:
 - A high-level description on how your code works.
 - The accuracies you obtain under various settings.
 - Explain which options work well and why.
 - If all your accuracies are low, tell us what you have tried to improve the accuracies and what you suspect is failing.

You can use any programming language of your choice

You need to implement all the **algorithms** by yourself (do not use libraries or already established functions to calculate entropy, information gain, chi-square test, etc)

This is an individual project. Only if you are register in the 429 (undergraduate) section you can work in pairs. Do not share your code or your report with your classmates, you can however ask questions in piazza providing small snippets of your code for discussion.

Rubric:

- Your code is thoroughly commented (5 pts)
- You provided a README file (5 pts)
- Your code executes correctly and is platform independent and you provided a Makefile or script for execution and clear instructions for execution are provided in the README file (20 pts)
- Implementation of:
 - ID3 implementation (10 pts)
 - Misclassification impurity (5 pts)
 - Information gain (5 pts)
 - Split stopping using chi-square (10 pts)
 - Accuracies for split stopping at 99, 95, and 0 (10 pts)
- Discussion in the report about:
 - A high-level description on how your code works.(5 pts)
 - The accuracies you obtain under various settings. (5 pts)
 - Explain which options work well and why. (5 pts)
 - If all your accuracies are low, tell us what you have tried to improve the accuracies and what you suspect is failing. (5 pts)
- Your report is clear, concise and well organized 10 pts
- TOTAL: 100 pts (10 pts of your final grade)