

Assignment 3: Decision Trees and Naive Bayes Model

CS486/686 – Winter 2014

Out: March 4, 2014

Due: March 19 (11:59 pm), 2014. Submit an electronic copy of your assignment via LEARN.
Late assignments may be submitted within 24 hrs for 50% credit.

Be sure to include your name and student number with your assignment.

Text categorization is an important task in natural language processing and information retrieval. For instance, news articles, emails or blogs are often classified by topics. In this assignment, you will implement (in the language of your choice) a decision tree algorithm and a naive Bayes model to learn a classifier that can assign a newsgroup topic to any article. Download a training set and test set of articles with their correct newsgroup label from the course website. To simplify your implementation, these articles have been pre-processed and converted to the *bag of words* model. More precisely, each article is converted to a vector of binary values such that each entry indicates whether the document contains a specific word or not.

1. [60 pts] Decision Tree Learning

Implement a decision tree learning algorithm. Here, each decision node corresponds to a word feature. Instead of building a decision tree in a depth first manner as described in the lecture slides, use a priority queue to greedily expand the decision tree at the leaf that can be improved the most. The depth first approach assumes that there are only a few features to split on and therefore a full tree can be built. Since there are many word features in text categorization, it is not practical to build a complete tree. Instead you will greedily build a partial tree by splitting the leaf that can be improved the most at each iteration. Use a priority queue to decide which leaf to split at each step with a new word feature. Experiment with priority queues that rank leaves by

- (a) information gain and
- (b) information gain times number of training documents at the leaves.

To measure the information gain of a leaf use the best information gain achieved by any word feature when splitting this leaf. Grow the tree by splitting one leaf at a time until there are 100 interior nodes. Do this by training with the training set only.

Report the training and testing accuracy (i.e., percentage of correctly classified articles) of each tree (from 1 to 100 interior nodes) by producing two graphs (one for each type of queue) with two curves each (one curve for training accuracy and one curve for testing accuracy).

What to hand in:

- A printout of your code.
- A printout (or hand drawing) showing two decision trees (one tree for each type of queue) with the first 10 word features selected and their information gain measure (times the number of documents for the second tree).
- Two graphs (one for each type of queue) showing the training and testing accuracy as the number of nodes increases.
- Does overfitting occur? If yes, after how many nodes does overfitting occur?
- A brief discussion of the word features selected by the decision tree learning algorithm for each type of queue. In your opinion, did all the word features selected make sense? Which queue is better and why?

2. [40 pts] Naive Bayes Model

Learn a naive Bayes model by maximum likelihood learning with Laplace smoothing (a.k.a. add-one smoothing). More precisely, learn a Bayesian network where the root node is the label/category variable with one child variable per word feature. Learn the parameters of the model by maximizing the likelihood (with Laplace smoothing) of the training set only. Classify documents by computing the label/category with the highest posterior probability $\Pr(\text{label}|\text{words in document})$. Report the training and testing accuracy (i.e., percentage of correctly classified articles).

What to hand in:

- A printout of your code.
- A printout listing the 10 most discriminative word features measured by

$$\max_{\text{word}} |\log \Pr(\text{word}|\text{label}_1) - \log \Pr(\text{word}|\text{label}_2)|$$

Since the posterior of each label is multiplied by the conditional probability $\Pr(\text{word}|\text{label}_i)$, a word feature should be more discriminative when the ratio $\Pr(\text{word}|\text{label}_1)/\Pr(\text{word}|\text{label}_2)$ is large and therefore when the difference between $\log \Pr(\text{word}|\text{label}_1)$ and $\log \Pr(\text{word}|\text{label}_2)$ is large. In your opinion, are these good word features?

- Training and testing accuracy (i.e., two numbers indicating the percentage of correctly classified articles for the training and testing set).
- The naive Bayes model assumes that all word features are independent. Is this a reasonable assumption? Explain briefly.
- What could you do to extend the Naive Bayes model to take into account dependencies between words?
- Which approach performs best among decision trees and the naive Bayes model? Explain briefly why.