

Assignment 3: Decision Trees and Naive Bayes Model

CS486/686 – Spring 2015

Out: June 23, 2015

Due: July 8 (11:59 pm), 2015. 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. [50 pts] Decision Tree Learning

Implement a decision tree learning algorithm. Here, each decision node corresponds to a word feature, which is selected by maximizing the information gain. Instead of designing a decision tree learning algorithm that builds a full tree, which may overfit, design your algorithm to take as input a maximum depth. Experiment with your algorithm by building trees with increasing maximum depth until a full tree is obtained. Report the training and testing accuracy (i.e., percentage of correctly classified articles) of each tree by producing a graph with two curves (one curve for training accuracy and one curve for testing accuracy as a function of the maximum depth). Report also the tree that achieved the highest testing accuracy.

What to hand in:

- A printout of your code.
- A graph showing the training and testing accuracy as the maximum depth increases.
- Does overfitting occur? If yes, after what maximum depth does overfitting occur?
- A printout (or hand drawing) showing the decision tree that achieved the highest testing accuracy. At each leaf, show the class and at each internal node, show the word feature with its information gain.
- A brief discussion of the word features selected by the decision tree that achieved the highest testing accuracy. In your opinion, did all the word features selected make sense?

2. [50 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.