

4F13 Machine Learning: Coursework #3: Latent Dirichlet Allocation

Carl Edward Rasmussen

Due: 12:00 noon, Nov 30th, 2018 online via moodle

Your answers should contain an explanation of what you do, and 2-4 central commands to achieve it (but complete listings are unnecessary). You must also give an *interpretation* of what the numerical values and graphs you provide *mean* – why are the results the way they are? **Each question should be labelled and answered separately and distinctly.** Total combined length of answers must not exceed 1000 words; clearly indicate the actual total number of words in your coursework. All questions carry approximately equal weight.

In this assignment, we will give you two short pieces of matlab code, which implement the main ingredients of Gibbs sampling for a Mixture of Multinomials `bmm.m` and for LDA `lda.m`. Before you start answering questions, you should spend some time understanding in detail, what this code does. This will enable you to answer all the questions with very little programming effort on your part.

The data is in the file `kos_doc_data.mat`. The word counts are in the matrix variables `A` and `B` for training and testing respectively, both matrices with 3 columns: document ID, word ID and word count. The words themselves are the variable `V`, such that eg. `V(841) = 'bush'`.

- a) Using the training data in `A`, find the maximum likelihood multinomial over words, and show the 20 largest probability items in a histogram. You may use the `barh` command followed by the command `set(gca, 'YTickLabel', V(s), 'Ytick', 1:20)`, where `s` is an array of appropriate indices. Using this multinomial model, what will the test set probability and/or log probability be if the test set `B` contains a word which is not contained in the training set `A`? Explain the implications of this.
- b) Instead of the maximum likelihood fit in question a), do Bayesian inference using a symmetric Dirichlet prior with a concentration parameter $\alpha = 0.1$ on the word probabilities. Compare the expressions for the predictive word probabilities for these two types of inference, and explain the implications, both for common and rare words.
- c) For the Bayesian model, what is the log probability for the test document with ID 2001? Explain whether, when computing the log probability of a test document, you would use the multinomial with or without the “combinatorial factor”. What is the per-word perplexity for the document with ID 2001? What is the per-word perplexity over all documents in `B`? Explain why the perplexities are different for different documents? What would the perplexity be for a uniform multinomial?
- d) The `bmm.m` script implements Gibbs sampling for a mixture of multinomials model. Use and modify the script to plot the evolution of the mixing proportions as a function of the number of Gibbs sweeps up to 20 iterations. The mixing proportions are the posterior probabilities of each of the mixture components. Rerun with different random seeds. Explain carefully, how would you determine whether the Gibbs sampler converges to and explores the stationary distribution (the posterior), does it?
- e) Use and modify `lda.m`. Plot topic posteriors for $K = 20$ as a function of the number of Gibbs sweeps, up to 20 sweeps. Comment on these. Compute the perplexity for the documents in `B` for the state after 20 Gibbs sweeps, and compare to previously computed perplexities. Are 20 Gibbs sweeps adequate? Plot the word entropy (what units do you use?) for each of the topics as a function of the number of Gibbs sweeps. Explain what you see.

Note that the performance and learning time for LDA depends a lot on the number of topics K .