

# DATA.STAT.840 Statistical Methods for Text Data Analysis

## Exercises for Lecture 6: Topic models

### Exercise 6.1: Latent semantic analysis

- (a) Download the 20 Newsgroups data set from <http://qwone.com/~jason/20Newsgroups/>.
- (b) In this exercise we consider only four of the newsgroups: rec.autos, rec.motorcycles, rec.sport.baseball, and rec.sport.hockey. Process the documents of the four newsgroups using the pipeline described on the lectures, including vocabulary pruning.
- (c) Create a TF-IDF representation for the documents, using Length-normalized frequency (TF) and Smoothed logarithmic inverse document frequency (IDF).
- (d) Apply latent semantic analysis to the TF-IDF matrix, to find 10 underlying factors.
- (e) Describe the resulting factors: list the 10 words with highest (absolute) weight in each factor. Do the factors seem related to individual newsgroups? Does their content seem meaningful?
- (f) Do the same with 15 factors (the first 10 factors will be the same). Do the new 5 factors seem more or less meaningful?

Report your analysis results and your code.

### Exercise 6.2: Probabilistic latent semantic analysis

- (a) Using the same data as in Exercise 6.1 (four newsgroups), create a term frequency matrix of raw term counts for the documents.
- (b) Apply PLSA to the term frequency matrix to find 10 underlying factors.
- (c) Describe the resulting factors: list the 10 words with highest probability in each factor.
- (d) Find, for each factor, the document (message) with highest probability of that factor, and print its 100 first words.
- (e) Do the factors seem related to individual newsgroups? Does their content seem meaningful?

Report your analysis results and your code.

### Exercise 6.3: Latent Dirichlet allocation

- (f) Using the same data as in Exercise 6.1 (four newsgroups), create a term frequency matrix of raw term counts for the documents.
- (g) Apply Latent Dirichlet Allocation to the term frequency matrix to find 10 underlying topic.
- (h) Describe the resulting factors: list the 10 words with highest probability in each topic.
- (i) Find, for each topic, the document (message) with highest probability of that topic, and print its 100 first words.
- (j) Do the topics seem related to individual newsgroups? Does their content seem meaningful?
- (k) Carry out (g)-(i) with 15 topics instead. Are some of the topics the same as before? What are the differences?

Report your analysis results and your code.