



Purpose

Text analysis and text mining are powerful tools for automatically extracting patterns from unstructured texts, such as texts from open-ended survey questions. The open-source statistical programming language R includes modules for text analysis such as tm (Text Mining) and NLP/openNLP (Natural Language Processing) (<http://cran.r-project.org/web/views/NaturalLanguageProcessing.html>).

Shiny (<http://shiny.rstudio.com/>) is an application framework for creating interactive web interfaces to data analysis programs in R. Analysis outputs are wrapped in html markup for viewing in a web browser. Input forms allow the user to change analysis parameters interactively and immediately view the results.

A Shiny application is presented for exploratory analysis of open-ended text data, such as course evaluation comments. Texts are analyzed based on lexical n-gram frequency and grouped into “topics” using Latent Dirichlet Allocation (LDA). The user can change analysis parameters (such as minimum frequency, whether to use stemming, n-gram size, number of topics, etc.) interactively while viewing the results. Model results and parameters can be downloaded to an Excel file.

This tool could be useful for generating coding categories prior to doing a traditional content analysis, or automated exploration of a large set of open-ended responses. Source code is provided for use as a teaching tool or further development.

Method

The texts explored consisted of all student responses to open-ended questions about what they “most enjoyed” about my online courses from Fall 2009 to Spring 2014 (5 sections of Cognitive Psychology, 13 sections of Intro Psych I, and 5 sections of Intro Psych II). Each student response (295 total) was treated as one “text.”

Corpora are input as Excel files, one text per row. High-frequency English “stop” words were removed from the corpus, as were a small set of additional user-defined “stop” words in an Excel file.

R packages *tm* and *NLP* were used to process a corpus and create a term-document matrix. The *wordcloud* package was used to generate word clouds, *SnowballC* was used for stemming, and LDA topic models were created with the *topicmodels* package.

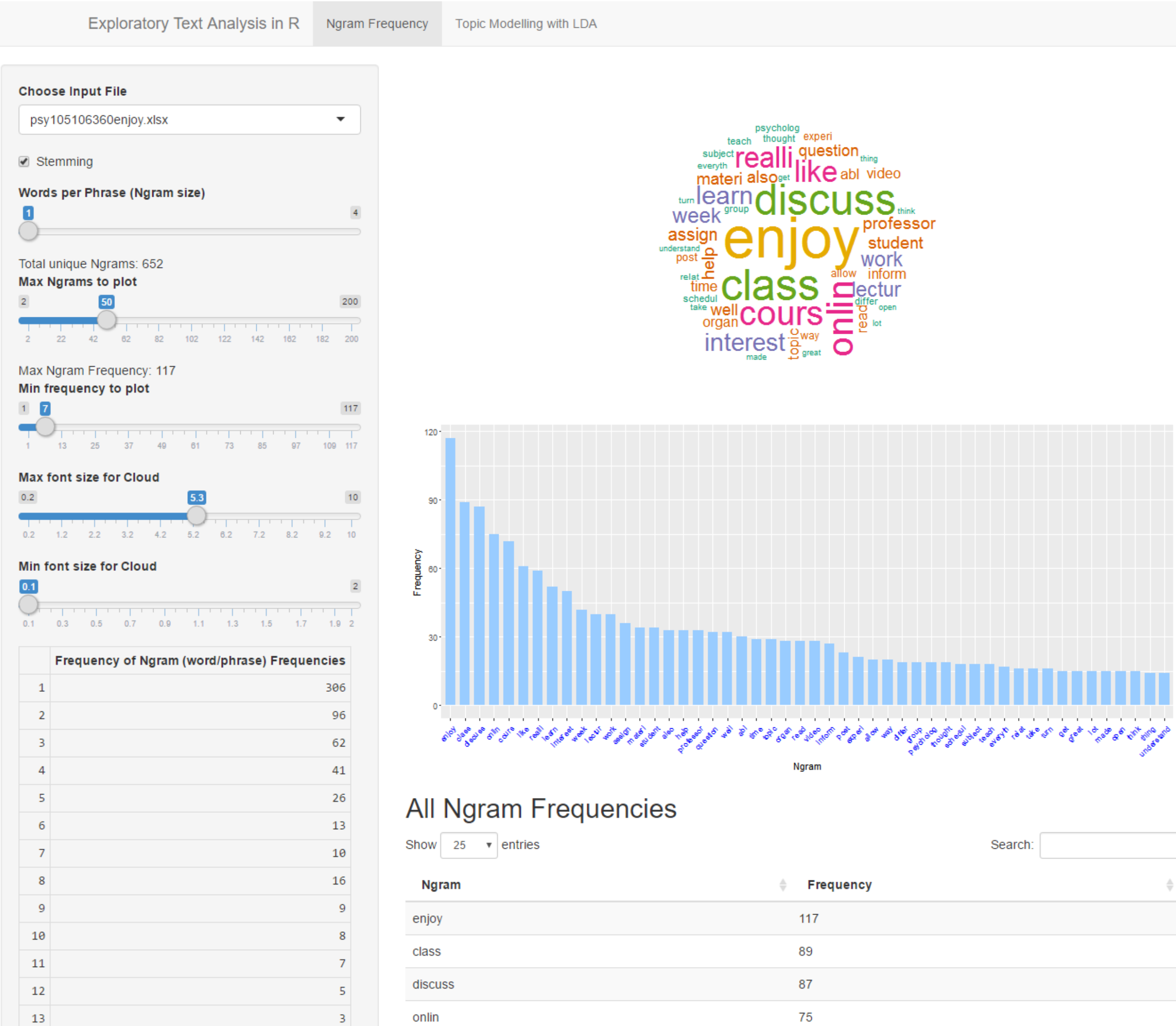
Results

Because the optimal n-gram size, number of topics, and other parameters may vary from corpus to corpus, it is helpful to be able to interactively change them and immediately see the results. In the example shown here, one of the topics appears to focus on flexibility (although in general this corpus did not produce any set of obviously interpretable and distinct topics).

Conclusions

Text mining can be a useful tool for exploring open-ended survey responses. R has packages that make text mining available to non-experts, and Shiny is a useful platform for making those capabilities accessible to a wider audience.

Ngram Frequency Analysis Tab



Downloaded Results File

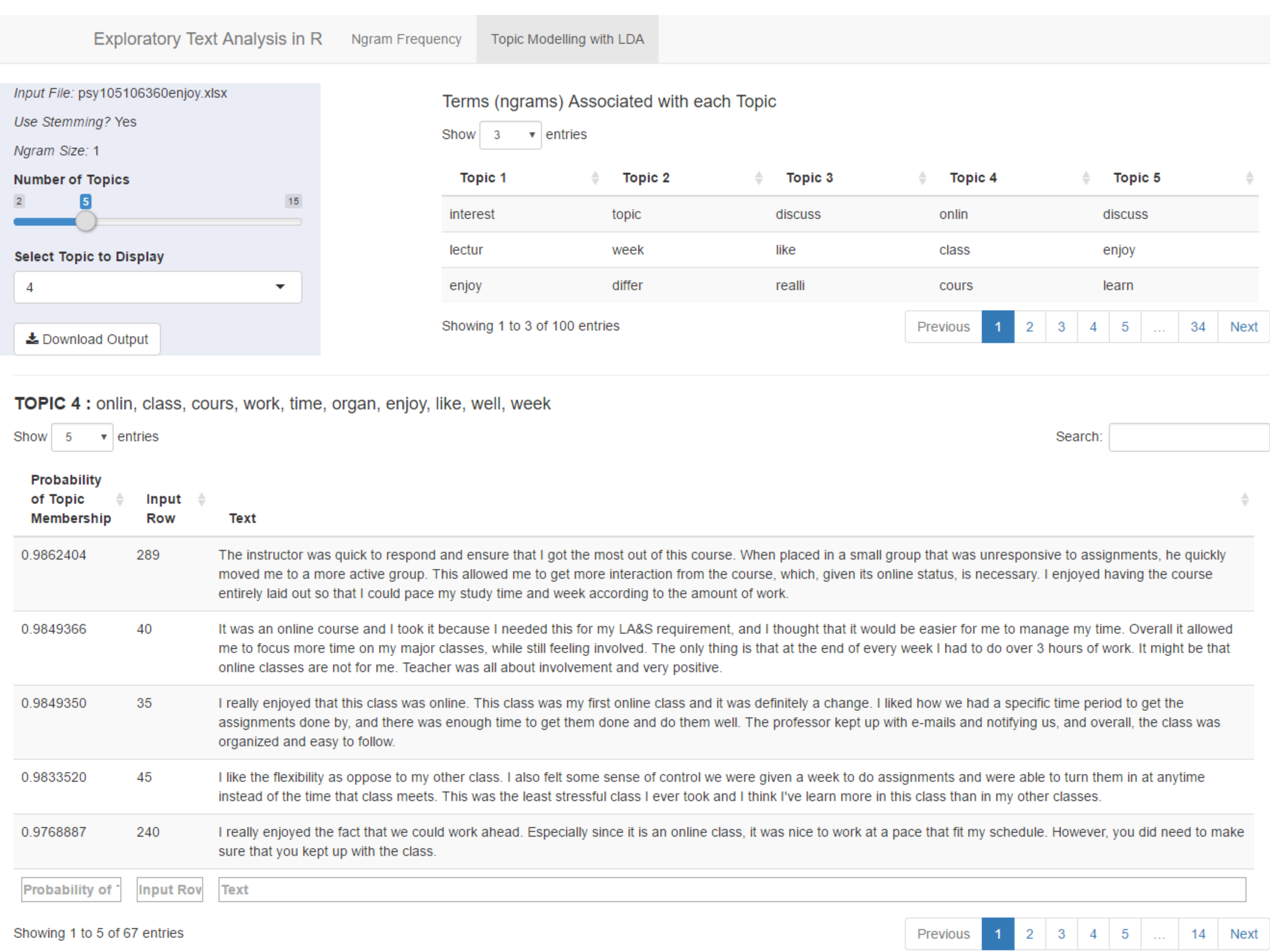
LDA_model_psy105106360enjoy (13).xlsx - Excel

Topic	Input Row	Probability for Topic 1	Probability for Topic 2	Probability for Topic 3	Probability for Topic 4	Probability for Topic 5	Text
1	5	0.0036282	0.0036282	0.11163456	0.0036283	0.8774808	initially, i thought the discussion board posts were tedious, but they're actual
3	3	0.0089842	0.0089834	0.64840588	0.0089832	0.3246433	i enjoyed the discussion questions because they made you really think about t
4	5	0.3137238	0.0294139	0.02941532	0.0294117	0.5980353	The discussion boards were nice
5	5	0.0158774	0.0158773	0.01587791	0.0158767	0.9364906	i liked people's examples in your turn to teach assignments
6	4	0.0410936	0.0410936	0.04109595	0.8356231	0.0410936	Videos and the experiments
7	3	0.0108712	0.0108701	0.81289491	0.1544925	0.0108714	Very interesting course, would have taken even if it were not a requirement. I
8	1	0.9565195	0.0108697	0.01087025	0.0108704	0.0108703	i enjoyed that this class was online and allowed me to work at my own pace. I

LDA_model_psy105106360enjoy (13).xlsx - Excel

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1		0.011236	0.019847	0.004727	0.007256	0.006252									
2	abl	7.71E-13	1.5E-137	0.001587	4.65E-53	9.1E-150									
3	access	9.5E-174	0.004179	2.17E-60	2.1E-187	1.3E-175									
4	accommod	2.43E-09	1.4E-133	0.003175	3.1E-154	1.7E-153									
5	accord	5.8E-192	6.1E-179	2.54E-58	0.001429	1.4E-191									
6	account	1.7E-161	3.6E-150	0.001587	3.6E-159	9.4E-161									
7	accumul	2.74E-32	0.002089	0.004717	4.08E-05	1.58E-40									
8	activ	1.37E-09	9.69E-65	0.012433	1.58E-12	0.000231									
9	actual	9.3E-129	3.27E-94	0.002432	5E-120	0.000644									
10	add														

LDA Topic Modeling Tab



Resources and Downloads

A working demo of the Shiny app is at <https://shiny.is.depaul.edu>

Source code, sample input files, and this poster are available on github.

Links to both can be found at:

<http://davidallbritton.com>

To install and use the app on your own computer:

1. Download and install R <http://www.r-project.org/>
2. Download and install the Desktop version of Rstudio <http://www.rstudio.com/>
3. Download the R script and sample input file (click “Download Zip” if you do not use GitHub) <https://github.com/davidallbritton/>
4. Place the R script and sample input files in a folder on your computer (after extracting from the zip archive)
5. Open the ui.R script in RStudio and click “run app”
6. Select “Run External” in the drop-down menu to the right of “Run App” for best results