

Workset 6: Macroanalysis

This week, rather than everyone try to do a cursory analysis we're going to start a sort of grouped sprint towards a preliminary finding. First, we'll divide into four groups of 5-6 around a specific corpus and problem; then we'll talk about it. The goal here is threefold:

1. Leverage teamwork so you can experience what it's like to undertake a text-analysis project, seeing where the major pitfalls and questions are without each person needing to do all the work that might entail on their own.
2. Engage in a discussion of what some interesting questions might be *from a particular corpus*.
3. Start trying out some forms of group collaboration before the project proposals are due in early November.

Instead of e-mailing the workset over the next week, we'll take forty minutes at the start of class on October 28th (in two weeks) to go over each of the four groups. So this will be a sort of a grouped sprint¹ towards a preliminary mini-works-in-progress conference.

The deliverables at the end will be:

1. Someone should be prepared to talk for eight minutes (no more), probably with slides. Make sure to choose someone from your group who can talk **succinctly**.
2. About 100-150 words on paper describing who did what.

Possible Corpora

Each group should choose a different method unless you've got a great justification

Possible corpora include:

1. Wordcount data from JSTOR data for research
2. Texts from the Internet Archive
3. Newspaper articles from the 19th century United States
4. Some smaller but *comprehensive* corpus.
 - For example, I've done some work with the State of the Union Addresses.
 - Or something tied to the group projects in public history.
 - Or build off of some of the documents scanned last week. You could try to generalize about OCR quality, say, using ngrams counts.

Possible methods

Each group should choose a different method unless you've got a great justification

- **Topic Modeling** You should be well familiar with it at this point. Actually running the model can be a bit of a pain, but there are detailed step-by-step instructions in a variety of places online, most notably the [Programming Historian](#). Does not require a large corpus.

1. In real life, the only actual grouped sprint in the world is the three-legged race. Ideally you'll look like the 1-mile relay team handing off perfectly, but don't worry if this ends up with your whole group face down in the mud, as long as you can explain how.

- **Text Vectors** This is similar to topic modeling in some ways: a hot topic from machine learning that identifies structural categories of language use in texts. But it principally focuses on just finding new elements of vocabulary that are similar to other words. You could use this to define a topic in the texts, you could use to uncover patterns in language use, or you could use it to find words that show a certain orientation: it provides a route towards answering weird questions like “where do foods fall on the salty-sweet spectrum in 19th century newspapers” or SAT-style analogies like “man:king::?:queen”. As discussed in class, this is relatively untraversed territory in Digital Humanities: high probability of failure, but a much greater chance of real novelty. Requires a very large corpus.
- **Named-entity extraction** This is slightly getting ahead of things, but you can use the Stanford Natural Language Toolkit to extract lists of names and places out of corpora of texts. Once you have the list, someone in the group can plot the locations on a map, as well.
- **N-gram analysis** Well chosen word counts can be the best data available. You might just want to look at numbers from Google Ngrams, from JStor, or you might want to tokenize yourself. (You could, say, compare the efficacy of different OCR software.)
- **Sentiment analysis** This is where some of Matt Jockers’ work has gone more recently: trying to track the existence of positive or negative language in a corpus. There are solid tools for sentiment analysis out there: one good one is Jockers’ own R package `syuzhet`. If you can think of a dataset where sentiment matters.
- **Other** Get in touch with me if you have a corpus that you think might benefit from some other method, whether one we saw in Macroanalysis or elsewhere.

Sample projects:

- Write a regular expression that captures 100 words around every use of the names of presidential candidates in elections between 1892 and 1920. Sentiment analyze them. Are there patterns of when positive language, negative language, or non-emotional language is used relative to the life-cycle of the campaign?
- Topic model a different field from JStor data than Underwood and Goldstone did.

Roles

You’ll have to figure out how to divide stuff up for this, which is why the set is not due for two weeks.

But some potential roles are:

- The project manager
- The corpus builder
- The corpus cleaner
- The algorithm runner
- The result inspector
- The talker.