



POCS
What's
The
Story?

Principles of Complex Systems, Vols. 1, 2, & 3D
CSYS/MATH 6701, 6713, & a pretend number
University of Vermont, Fall 2024

"Don't answer that, you'll get it wrong." [↗](#)

Assignment 13

[Frankie and the Dean](#) [↗](#), [Community](#), [Advanced Safety Features](#), [S6E07](#).

Episode links: [Wikipedia](#) [↗](#), [IMDB](#) [↗](#), [Fandom](#) [↗](#), [TV Tropes](#) [↗](#).

Due: Monday, February 10, by 11:59 pm

<https://pdodds.w3.uvm.edu/teaching/courses/2024-2025pocsverse/assignments/13/>

Some useful reminders:

Deliverator: Prof. Peter Sheridan Dodds (contact through Teams)

Office: The Ether and/or Innovation, fourth floor

Office hours: See Teams calendar

Course website: <https://pdodds.w3.uvm.edu/teaching/courses/2024-2025pocsverse>

Overleaf: \LaTeX templates and settings for all assignments are available at

<https://www.overleaf.com/read/tsxfwwmwdgxj>.

Some guidelines:

1. Each student should submit their own assignment.
2. All parts are worth 3 points unless marked otherwise.
3. Please show all your work/workings/workingses clearly and list the names of others with whom you ~~conspired~~ collaborated.
4. We recommend that you write up your assignments in \LaTeX (using the Overleaf template). However, if you are new to \LaTeX or it is all proving too much, you may submit handwritten versions. Whatever you do, please only submit single PDFs.
5. For coding, we recommend you improve your skills with Python, R, and/or Julia.
Please do not use any kind of AI thing unless directed. The (evil) Deliverator uses (evil) Matlab.
6. There is no need to include your code but you can if you are feeling especially proud.

Assignment submission:

Via Brightspace (which is not to be confused with the death vortex of the same name).

Again: One PDF document per assignment only.

Please submit your project's current draft in pdf format via Brightspace four days after the due date for this assignment (normally a Friday). For teams, please list all team member names clearly at the start.

Semester goal: A paper per text studied, building through assignments.

Three stories to analyze:

- **Pride and Prejudice**

<https://www.gutenberg.org/ebooks/1342>

- **Frankenstein; or the Modern Prometheus**

<https://www.gutenberg.org/files/84/84-h/84-h.htm>

- **Moby Dick; Or, The Whale**

<https://www.gutenberg.org/ebooks/2701>

Data:

For this assignment, the novels have been processed into 1-grams with an attempt to capture all elements including punctuation.

You can use the data below, or what you produced in the previous assignment.

The basic data format is as a time series with one 1-gram per line (links below).

For each story, also linked to below are the rank distributions of 1-grams by counts.

https:

[//pdodds.w3.uvm.edu/permanent-share/pride_and_prejudice_narrativetimeseries.txt](https://pdodds.w3.uvm.edu/permanent-share/pride_and_prejudice_narrativetimeseries.txt)

https://pdodds.w3.uvm.edu/permanent-share/pride_and_prejudice_1grams.txt

https://pdodds.w3.uvm.edu/permanent-share/frankenstein_narrativetimeseries.txt

https://pdodds.w3.uvm.edu/permanent-share/frankenstein_1grams.txt

https://pdodds.w3.uvm.edu/permanent-share/moby-dick_narrativetimeseries.txt

https://pdodds.w3.uvm.edu/permanent-share/moby-dick_1grams.txt

Instrument—Shifterator:

- For word shifts, use either:

The Python package described in Ref. [1].

Or the D3.js shifterator created by Andy Reagan here:

<https://observablehq.com/@andyreagan/d3-shifterator-v4-preloaded>.

Links to paper versions (arXiv is always best), Github repository, and an exhilarating Twitter feed can be found here:

<https://pdodds.w3.uvm.edu/research/papers/gallagher2021a/>.

Various Matlab versions made by the Unreliable Deliverator do exist and need to be shared on Gitplaces. A more sophisticated map+list version has long been in development.

Github repository: <https://github.com/ryanjgallagher/shifterator>

Very Important Notes for the Python version:

- The default setting for reference happiness score for the labMT lexicon in the Shifterator **is wrong**. It's set at 5. Please make sure to use the average of the reference text. Again: **It's wrong as is**, and this is due to a misunderstanding that cannot be righted. There is only so much the Deliverator can do.
- You will need Python 3.8. Apparently it all goes to pieces otherwise.

Things to do:

1. Lexical calculus (a subset of type calculus):

Derive the word shift equation for simple additive lexical instruments.

You will have the derivation per class.

The idea is to simply work through it yourself.

There are no advanced mathematics here.

But over and over, people do not understand what's going on.

Word shifts are a kind of discrete derivative (difference) with words on the inside.

Per lectures, the goal is to derive.

$$\delta h_{\text{avg},i} = \frac{100}{|h_{\text{avg}}^{(\text{comp})} - h_{\text{avg}}^{(\text{ref})}|} \underbrace{[h_{\text{avg}}(w_i) - h_{\text{avg}}^{(\text{ref})}]}_{+/-} \underbrace{[p_i^{(\text{comp})} - p_i^{(\text{ref})}]}_{\uparrow/\downarrow}$$

Performed in class and in numerous papers [2, 3, 4].

2. (9 points, 3 points for each novel)

Take the UTF-8 text versions of each of these three novels and parse them into narrative time series of 1-grams which include all punctuation, numbers, and words.

See below for instructions.

For each novel, present your output for the first paragraph, rendered as a single, wrapped line.

3. (9 points, 3 points for each novel)

For each novel, generate size rank distributions of 1-grams with size being counts.

Display the 1-grams and counts for the first 100 1-grams.

4. (9 points, 3 points for each novel)

For each novel, plot the size rank distribution for 1-grams.

Outputs provided for comparison:


The basic data format is as a time series with one 1-gram per line (links below).

For each story, also linked to below are the rank distributions of 1-grams by counts.

https://pdodds.w3.uvm.edu/permanent-share/pride_and_prejudice_narrativetimeseries.txt

https://pdodds.w3.uvm.edu/permanent-share/pride_and_prejudice_1grams.txt
https://pdodds.w3.uvm.edu/permanent-share/frankenstein_narrativetimeseries.txt
https://pdodds.w3.uvm.edu/permanent-share/frankenstein_1grams.txt
https://pdodds.w3.uvm.edu/permanent-share/moby-dick_narrativetimeseries.txt
https://pdodds.w3.uvm.edu/permanent-share/moby-dick_1grams.txt

General instructions for processing texts:

- Using an editor, remove the start and finish material for each of the three novels that Gutenberg adds to books somewhat inconsistently.
- Using a judicious selection of regex operations, create a script* (Python is strongly recommended but you can use whatever you like) to break up the text into meaningful 1-grams that may be words, punctuation, and numbers.
Regex = Regular Expression. 
* If you want to use tokenizer, you can. But the idea is to be careful and really understand what's happening the text as you smash it into pieces (storyons).
- See Perl code below for an example.

```
## for Frankenstein
$text =~ s/D--n/Damn/g;

## separate out some basic punctuation
$text =~ s/([\!?\,\.\])/ \1 /g;
$text =~ s/:/ : /g;
$text =~ s;/ / \; /g;

## remove underscores used for emphasis
$text =~ s/_//g;

## isolate parentheses
$text =~ s/\(/ ( /g;
$text =~ s/\)/ ) /g;

## dash madness
$text =~ s/----/ --- /g; ## long dash
$text =~ s/--/ --- /g; ## em dash
$text =~ s;/- / --- /g; ## em dash
$text =~ s/- / --- /g; ## em dash
```

```

## handle specific salutations
$text =~ s/Mr \./Mr./g;
$text =~ s/Mrs \./Mrs./g;
$text =~ s/Dr \./Dr./g;

## clean up white space duplication
$text =~ s/\s+ /g;

## separate quotes
## opening quotes should have a space before them (except for em dashes, treat
$text =~ s/\s"/ " /g;
## closing quotes should be what's left:
$text =~ s/"/ " /g;
$text =~ s"/ " /g;
$text =~ s"/ " /g;

## separate off opening single quotation mark
$text =~ s/'/ ' /g;

## dysfunctional catapostrophes:
## clean up apostrophes and opening and closing single quote mark
## closing single quote mark should generally be isolated
## leave alone to preserve contractions
$text =~ s/'/ ' /g;

$text =~ s'/ '/g;
## opening quote mark
## \p{L} stands for any UTF-8 letter
$text =~ s/(\s)'(\p{L})/\1' \2/g;
## closing quote mark (will be a problem with contractions)
$text =~ s/(\p{L})'(\s)/\1 '\2/g;

## split off possession indicator
$text =~ s/'s/ 's/g;

## remove any white space at the front
$text =~ s/^\s+//;

## add new line at the end

```

```

$text = $text."\n";

## now create time series text by replacing spaces with returns
($timeseriestext = $text) =~ s/ /\n/g;

## last: remove any white space redundancies
$timeseriestext =~ s/\s+/\n/ms;

```

5. Measure the average happiness of each text using the labMT word list with the lexical lens:

$$\mathcal{L} = \{\tau \in \Omega \mid h_{\text{avg}}(\tau) \leq 4 \text{ or } h_{\text{avg}}(\tau) \geq 6\} \quad (1)$$

where Ω is the labMT lexicon and τ is a word in Ω .

6. (3 points each)

For the following T_{ref} and T_{comp} for each novel, generate a collection of word shifts as described below.

Continue to use the same lexical lens \mathcal{L} as above.

Pride and Prejudice:

- Pride and Prejudice:
 T_{ref} = the first 40% of the book,
 T_{comp} = 70% to 75% of the book.
- Frankenstein:
 T_{ref} = the first 20% of the book,
 T_{comp} = the last 10% of the book.
- Moby Dick:
 T_{ref} = the whole book,
 T_{comp} = 80% to 90% of the book.

- (a) For each novel, produce word shifts comparing text T_{comp} relative to text T_{ref} . Important: Use the average happiness of text T_{ref} as the baseline (this is not the default in the Python package).
- (b) Interpret the word shifts. Does what you see make sense? Are there any surprises? Are some words being used in what the average person might not think is their primary meaning? For example, “crying” in Moby Dick means yelling, and “sick” can mean “awesome.”
- (c) Reverse the comparison: Produce word shifts comparing text T_{ref} relative to text T_{comp} . Now use the average happiness of text T_{comp} as the baseline.

- (d) Comment on any asymmetries you see (the basic word shifts we use are asymmetric).
 - (e) Produce word shifts comparing text T_{ref} relative to text T_{comp} . Now use 5 as the baseline reference score (neutral on the happiness-sadness spectrum of 1–9 for labMT).
 - (f) Compared to your first word shifts, how interpretable are these ones?
- 7.
- Please use Overleaf for writing up your project.
 - Build your paper using:
<https://github.com/petersheridandodds/universal-paper-template>
 - Please use Github and Gitlab to share the code and data things you make.
 - For this first assignment, just getting the paper template up is enough.

Some of you may of course have projects underway from PoCS, Vol. I.

If not, please begin formulating project ideas.

See storyology slides.

General suggestion: Come up with some rich, text-based set of stories for analysis.




For example: One (longish) book, or a book series, or a TV series.

Data would be the original text (books), subtitles, screenplay, or scripts (TV series).




- You must be able to obtain the full text.
- You will want something with at least around 10^5 words. More than 10^6 would be great.
- Transcripts of shows may be good for extracting temporal character interaction networks.

Please talk about possibilities with others in the class.

For this assignment, simply list at least one possibility (which may be your existing project), noting the approximate text size in number of words, or whatever measure of size makes sense.

- Explore and find inspiration in the Pudding.
 Example: Film dialogue analysis.
<https://pudding.cool/2017/03/film-dialogue/> 
- [ConvoKit](#) : Cornell Conversation Analysis Toolkit
- [Open Psychometrics](#) : Plenty of stuff here including overhauling OCEAN for regular people.

References

- [1] R. J. Gallagher, M. R. Frank, L. Mitchell, A. J. Schwartz, A. J. Reagan, C. M. Danforth, and P. S. Dodds. Generalized word shift graphs: A method for visualizing and explaining pairwise comparisons between texts. *EPJ Data Science*, 10:4, 2021. Available online at <https://arxiv.org/abs/2008.02250>. pdf 
- [2] P. S. Dodds and C. M. Danforth. Measuring the happiness of large-scale written expression: Songs, blogs, and presidents. *Journal of Happiness Studies*, 11(4):441–456, 2009. pdf 
- [3] P. S. Dodds, K. D. Harris, I. M. Kloumann, C. A. Bliss, and C. M. Danforth. Temporal patterns of happiness and information in a global social network: Hedonometrics and Twitter. *PLoS ONE*, 6:e26752, 2011. pdf 
- [4] P. S. Dodds, E. M. Clark, S. Desu, M. R. Frank, A. J. Reagan, J. R. Williams, L. Mitchell, K. D. Harris, I. M. Kloumann, J. P. Bagrow, K. Megerdumian, M. T. McMahon, B. F. Tivnan, and C. M. Danforth. Human language reveals a universal positivity bias. *Proc. Natl. Acad. Sci.*, 112(8):2389–2394, 2015. Available online at <http://www.pnas.org/content/112/8/2389>. pdf 