

Alok, ~~De~~ Deytia, Tyler, Kapil TEAM 1

Problems with data

↳ Say
↳ Says
↳ Said

- incorrect regex leads to missing ^{or extra} words
- vernacular words
- Extra data (title and author randomly placed in text)
- inadequate translations
- time frame from the original
-

Tz

1. Issue w/ translatability - context
2. Issues of Time Period, gender - context
3. genre of works - context
4. semantics and meaning of two words
With same spelling but different meanings - Context
5. legal constraints for providing a book online for free - missing data

Chunyan Tang

John Duggan

R.J. Lusk

Team 3 - DJ, Mark, Sam

~~Data~~ Issues with Data:

- Empty space as a word (stopwords + html cleaning didn't catch)
- Books not on Project Gutenberg
- Extraneous words (e.g., "Gutenberg", "copyright")
- Words that don't carry significant meaning alone ("would", "could", "will")
- Words that are the same but have different meanings based on context ("set" can be a verb, adjective, or noun)

Problems with Data in Assignment 1

multiple contexts

- different meanings of words
- similar words not being counted as the same

Missing Context

- missing pages
- table of context index etc not relevant to analysis
- words in a vacuum, no information on word frequency in relation to other words

Incorrect tampered / filtered

- ~~copyright~~ copyright /
- legalise words appearing too frequently, formatting words that have nothing to do with context of text modifying results.
- incorrect digital copies
- empty characters

Data Clones

- Words with :
- says, said ?
- Capitalization - dult
- non-author text → Gutenberg press

why
some lines
after
said
not?

→ wrong author

→ wrong unit → intro, chapter?

→ translated → is it the original the same

→ why exclude stopwords

→ what about handwriting?

lost word order

→ look at infrequent words
or word distribution?

same with data