

Intro to text mining in R

R workshops 2021

Baruch College

Víctor Peña, January 2021

Logistics

- Course website: <https://vicpena.github.io/workshops/2021/advanced>
- My email: victor.pena@baruch.cuny.edu
- I'll cover some sections of Text Mining with R, by Silge and Robinson and the Datacamp courses
 - Introduction to text analysis in R
 - Introduction to Natural Language Processing in R
 - String manipulation with stringr
- I'll post the code we write in our sessions on the course website

Topics

- Through case studies, we'll cover:
 - Analyzing word frequencies: word counts, tf-idf
 - Analyzing relationships between words: n-grams, correlations
 - Sentiment analysis
 - Topic modeling
 - Working with text data in R

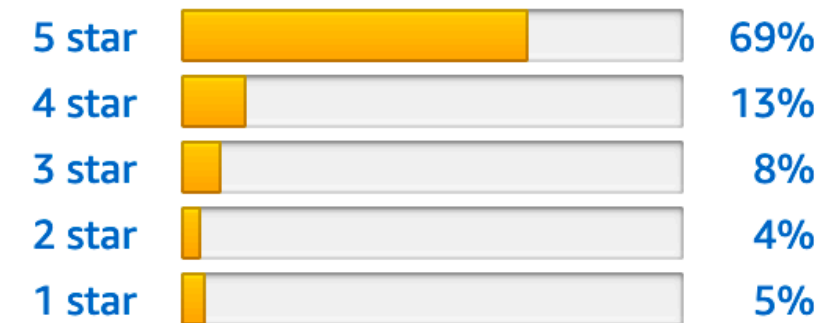
First case study: Roomba reviews

- We'll analyze a dataset which contains Amazon reviews for Roombas
- The dataset is clean and nicely formatted... which is not the usual starting point in text analysis
- The data can be found on the course website

Customer reviews

★★★★☆ 4.4 out of 5

7,078 global ratings



▼ [How are ratings calculated?](#)



iRobot Roomba 650 Robot Vacuum

by iRobot

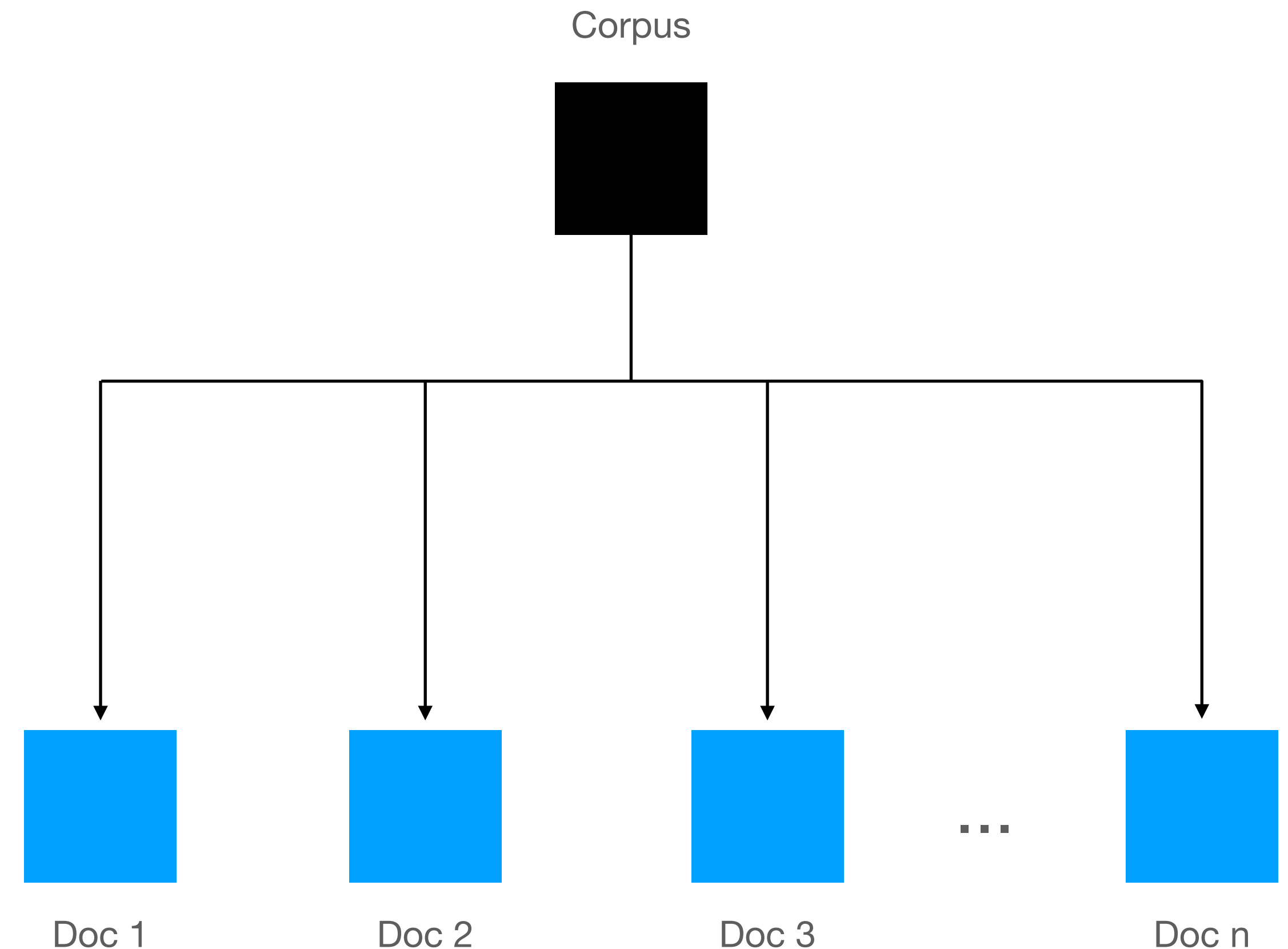
Size: R650 | [Change](#)

[Write a review](#)

Analyzing word/document frequencies

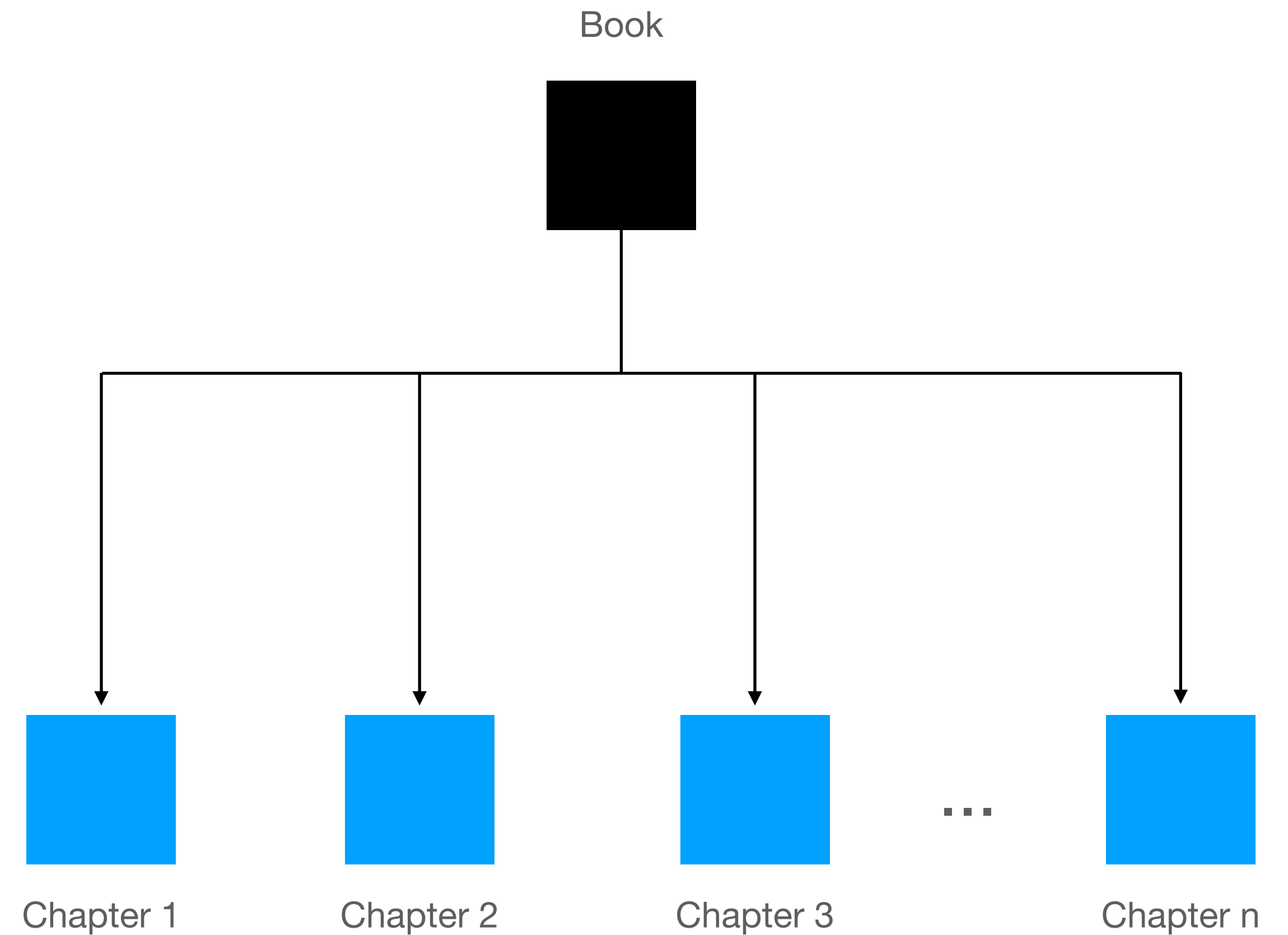
Conceptual framework

- We have a corpus of documents
- We want to know how different or similar the different documents are
- For example, are there certain words that are used more often in certain documents?
- Can we quantify "distance" or "correlation" between documents, in some way?



Example: book

- **Corpus:** book itself
- **Documents:** chapters
- **Questions we'll answer:**
 - Are there certain words that appear more often in certain chapters?
 - Which chapters are most similar?



TF-IDF is TF times IDF

- **TF:** term frequency of a word in a document

$$\text{TF}(\text{word, document}) = (\# \text{ times word appears in document}) / (\# \text{ words in document})$$

- **IDF:** log "inverse document frequency" of a word in the corpus

$$\text{IDF}(\text{word}) = \log[(\# \text{ documents}) / (\# \text{ documents containing word})]$$

- If a word appears in all documents $\Rightarrow \text{IDF}(\text{word}) = \log(1) = 0$
- If a word does not appear in any document $\Rightarrow \text{IDF}(\text{word}) = \log(\# \text{ docs} / 0) = \text{Infinity}$
- If $\text{IDF}(\text{word})$ is high, the word only appears in a few documents; if $\text{IDF}(\text{word})$ is low, it appears in most documents
- **TF-IDF:** simply multiply TF times IDF

$$\text{TF-IDF}(\text{word, document}) = \text{TF}(\text{word, document}) * \text{IDF}(\text{word})$$

$TF(\text{word}, \text{document}) = (\# \text{ times word appears in document}) / (\# \text{ words in document})$

$IDF(\text{word}) = \log[(\# \text{ documents}) / (\# \text{ documents containing word})]$

- **Document 1**

- cat, dog, cat, cat

- **Document 2**

- dog, dog, dog, mouse

- **Document 3**

- beetle, mouse, dog, aye-aye

TF-IDF(cat, doc 1)?

- $TF(\text{cat}, \text{doc 1}) = 3/4$
- $IDF(\text{cat}) = \log(3/1)$
- $TF-IDF \sim (3/4) * \log(3/1) \sim 0.8239$

$TF(\text{word}, \text{document}) = (\# \text{ times word appears in document}) / (\# \text{ words in document})$

$IDF(\text{word}) = \log[(\# \text{ documents}) / (\# \text{ documents containing word})]$

- **Document 1**

- cat, dog, cat, cat

- **Document 2**

- dog, dog, dog, mouse

- **Document 3**

- beetle, mouse, dog, aye-aye

TF-IDF(cat, doc 2)?

- $TF(\text{cat}, \text{doc 2}) = 0/4 = 0$
- $IDF(\text{cat}) = \log(3/1)$
- $TF-IDF(\text{cat}, \text{doc 2}) = 0$

If a term doesn't appear in a document
its TF-IDF is going to be 0

$TF(\text{word}, \text{document}) = (\# \text{ times word appears in document}) / (\# \text{ words in document})$

$IDF(\text{word}) = \log[(\# \text{ documents}) / (\# \text{ documents containing word})]$

- **Document 1**

- cat, dog, cat, cat

- **Document 2**

- dog, dog, dog, mouse

- **Document 3**

- beetle, mouse, dog, aye-aye

TF-IDF(dog, doc 2)?

- $TF(\text{dog}, \text{doc 2}) = 3/4$
- log-odds: $\log(TF/(1-TF))$
- $IDF(\text{dog}) = \log(3/3) = 0$
- $TF-IDF(\text{dog}, \text{doc 2}) = (3/4) * 0 = 0$

If a term appears in all documents, then
TF-IDF is going to be 0

$TF(\text{word}, \text{document}) = (\# \text{ times word appears in document}) / (\# \text{ words in document})$

$IDF(\text{word}) = \log[(\# \text{ documents}) / (\# \text{ documents containing word})]$

- **Document 1**

- cat, dog, cat, cat

- **Document 2**

- dog, dog, dog, mouse

- **Document 3**

- beetle, mouse, dog, aye-aye



TF-IDF(aye-aye, document 3)?

- $TF(\text{aye-aye}, \text{doc 3}) = 1/4$
- $IDF(\text{aye-aye}) = \log(3/1)$
- $TF-IDF = (1/4) * \log(3) \sim 0.275$

Odds and ends

- If a word in a document has a high TF-IDF, it means that that word appears *"weirdly often"* in that document
 - Another way of saying the same thing: *the word is "highly specific to that document"*
- In R, we won't have to do these computations by hand
 - We can use `bind_tf_idf` in `library(tidytext)`

Let's find TF-IDFs in a real example

- Let's work with the text of Animal Farm, by George Orwell

