2024-03-20

## Data in the Humanities

## Wk. 07: Scraping & Curating Structured Data

## Welcome, <u>Filipa Calado</u>! 👋

## Happenings This Week

### *Workshops*

– Thu 3/21 <u>Data Visualization with Tableau</u> (PUL)
– Tue 3/26 <u>Legal Issues in Computational Research Using Text and Data Mining</u> (PUL)
– Wed 3/27 <u>Getting Started with LaTeX</u> (PUL)

### *Events*

– <u>Syriac Transcribathon</u> (CDH / NES), 3/25 – 3/28
– LLM Reading Group: <u>AI & Industry</u> (CDH), 3/27

### *CFPs*

– <u>CDH Graduate Fellowship</u>, Fall 2024, due 3/31
– <u>DH for Hellenic Studies</u> CFP: Athens in June! Now due 3/31

## Survey Results

Let's review next week. Please answer if you haven't already!

## Final Project options

1. Seminar Paper
2. Speculative Data Curation
3. Speculative Project Proposal
4. Group collaboration
for next week

## install QGIS

https://www.qgis.org

# QGIS

**Readings** 📖

– Katie Rawson and Trevor Muñoz, "Against Cleaning," in *Debates in the Digital Humanities 2019* (University Of Minnesota Press, 2019).

– Henk Alkemade et al., "Datasheets for Digital Cultural Heritage Datasets" 9, no. 1 (2023)

– *recommended but not required*

  – Jessica Marie Johnson, "Markup Bodies: Black [Life] Studies and Slavery [Death] Studies at the Digital Crossroads"

  – Daniel Rosenberg, "Data Before the Fact"

  – Marisa Elena Duarte and Miranda Belarde-Lewis, "Imagining: Creating Spaces for Indigenous Ontologies"

Gebru et al., "Datasheets for Datasets" (2018)

contra The Pile &c.

### Datasheets for Datasets

Timnit Gebru [1]  Jamie Morgenstern [2]  Briana Vecchione [3]  Jennifer Wortman Vaughan [1]  Hanna Wallach [1]
Hal Daumé III [1 4]  Kate Crawford [1 5]

**Abstract**

The machine learning community has no standardized way to document how and why a dataset was created, what information it contains, what tasks it should and should not be used for, and whether it might raise any ethical or legal concerns. To address this gap, we propose the con-

We therefore propose the concept of datasheets for datasets. In the electronics industry, every component is accompanied by a datasheet describing standard operating characteristics, test results, and recommended usage. By analogy, we recommend that every dataset be accompanied with a datasheet documenting its motivation, creation, composition, intended uses, distribution, maintenance, and other information. We anticipate that such datasheets will increase transparency

https://datanutrition.org

https://github.com/jackbandy/bookcorpus-datasheet

## Dataset Facts

**Dataset** BookCorpus
**Instances Per Dataset** 7,185 unique books, 11,038 total

---

### Motivation

| | |
|---|---|
| **Original Authors** | Zhu and Kiros et al. (2015) [39] |
| **Original Use Case** | Sentence embedding |
| **Funding** | Google, Samsung, NSERC, CIFAR, ONR |

### Composition

| | |
|---|---|
| **Sample or Complete** | Sample, ≈2% of smashwords.com in 2014 |
| **Missing Data** | 98 empty files, ≤655 truncated files |
| **Sensitive Information** | Author email addresses |

### Collection

| | |
|---|---|
| **Sampling Strategy** | Free books with ≥20,000 words |
| **Ethical Review** | None stated |
| **Author Consent** | None |

### Cleaning and Labeling

| | |
|---|---|
| **Cleaning Done** | None stated, some implicit |
| **Labeling Done** | None stated, genres by smashwords.com |

### Uses and Distribution

| | |
|---|---|
| **Notable Uses** | Language models (e.g. GPT [29], BERT [9]) |
| **Other Uses** | List available on HuggingFace [12] |
| **Original Distribution** | Author website (now defunct) [39] |
| **Replicate Distribution** | BookCorpusOpen [13] |

### Maintenance and Evolution

| | |
|---|---|
| **Corrections or Erratum** | None |
| **Methods to Extend** | "Homemade BookCorpus" [21] |
| **Replicate Maintainers** | Shawn Presser [12] |

---

| Genres | % of BookCorpus* |
|---|---|
| **Romance** 2,881 books | 26.1% |
| **Fantasy** 1,502 books | 13.6% |
| **Vampires** 600 books | 5.4% |

| | |
|---|---|
| Horror 4.1% | • Teen 3.9% |
| Adventure 3.5% | • Literature 3.0% |
| Historical Fiction 1.6% | |

Not a significant source of nonfiction.

\* Percentages based on directories in books_txt_full. Some books cross-listed.

https://menus.nypl.org/