

Quantitative Analysis of Text

genres are icebergs: with a visible portion floating above the water, and a much larger part hidden below, and extending to unknown depths.

Quantitative Analysis of text

“operationalize” a text

“concepts are transformed into a series of operations—which, in their turn, allow to measure all sorts of objects. Operationalizing means building a bridge from concepts to measurement, and then to the world”—Moretti

Quantitative Analysis of text

“operationalize” a text

Docuscope: Smart dictionary

vs

MFW: Most frequent words

Quantitative Analysis of text

Unsupervised: Cluster Analysis

Quantitative Analysis of text

Cluster Analysis of Folio Plays

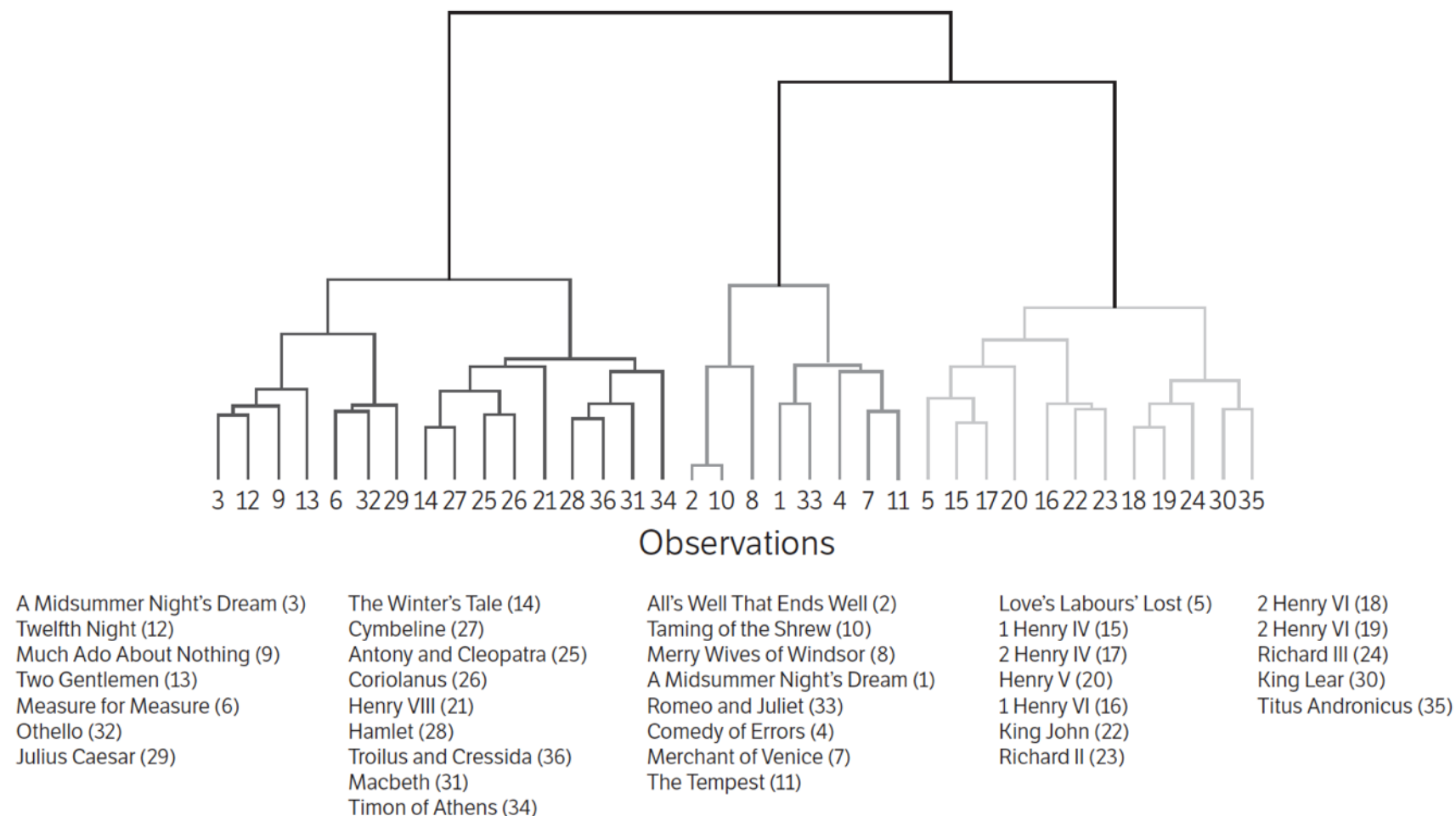
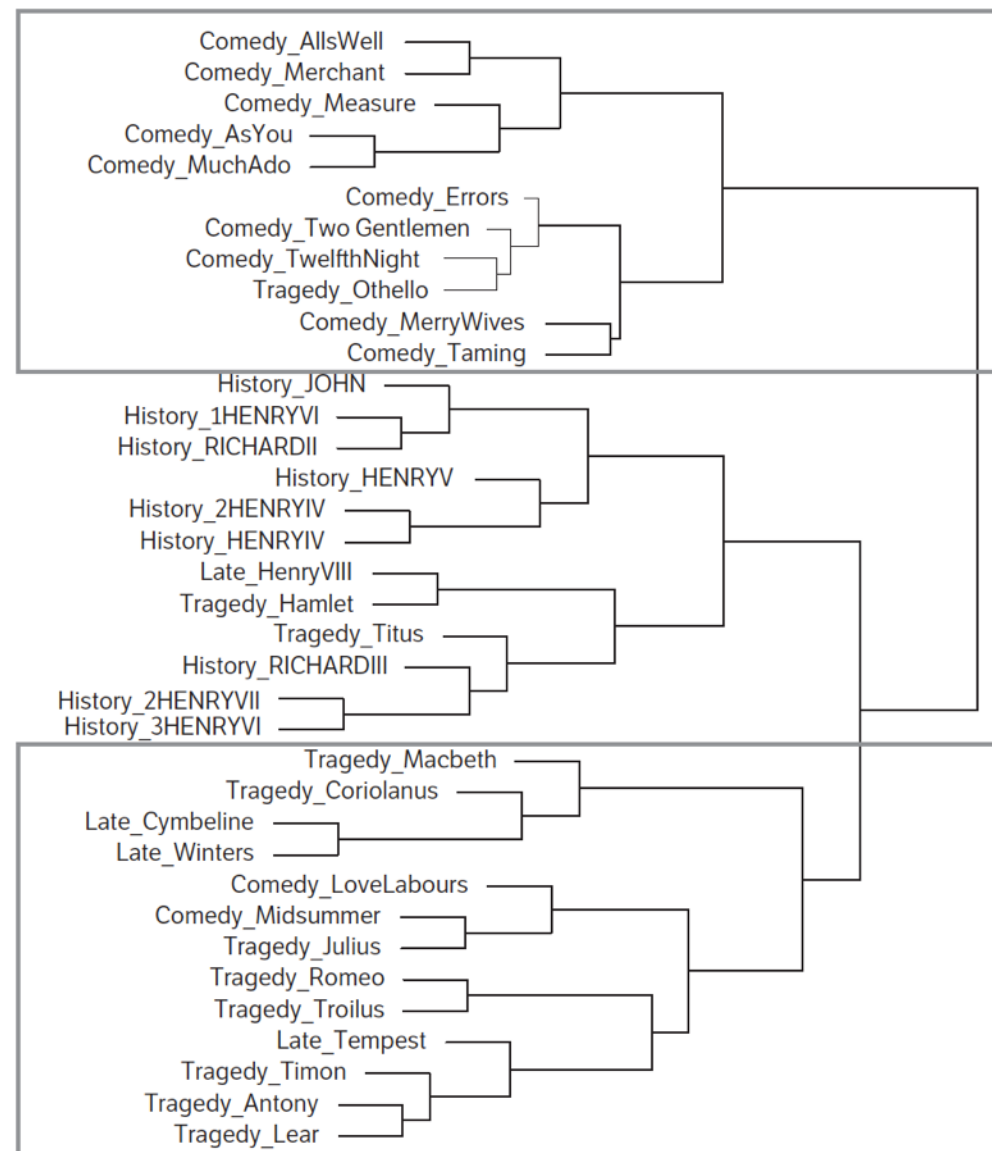


Figure 1: Dendrogram illustrating clustering of Shakespeare plays rated on DocuScope's Language Action Types (LATs) produced in 2003. Clustering method: complete linkage, Euclidean distances. Notice the presence of comedies in the first and third columns, late plays and tragedies in the second, and histories in the fourth and fifth. "Incorrect classifications" such as *Othello* and *Love's Labours' Lost* are discussed on Witmore's blog, www.winedarksea.org.

Quantitative Analysis of text

Shakespeare Plays
Using Euclidean Distance with Complete Linkage and 37 Features



Plot Created: Feb. 4, 2009
By: mjockers

Figure 3.2: Dendrogram of Shakespeare First Folio plays using Most Frequent Words with major clusters highlighted. Here Jockers used the 37 features from the Shakespeare plays that had a mean relative frequency of greater than or equal to .03%. Note the similarity between this tree and Docuscope's diagram in fig. 1.1, with the close pairings of *Winter's Tale* and *Cymbeline*; *2 Henry VI* and *3 Henry VI*, and the proximity of *Coriolanus* to the *Cymbeline-Winter's Tale* pair.

As soon as school was over, we met again.

Quantitative Analysis of text

PCA: Principal components analysis

Quantitative Analysis of text

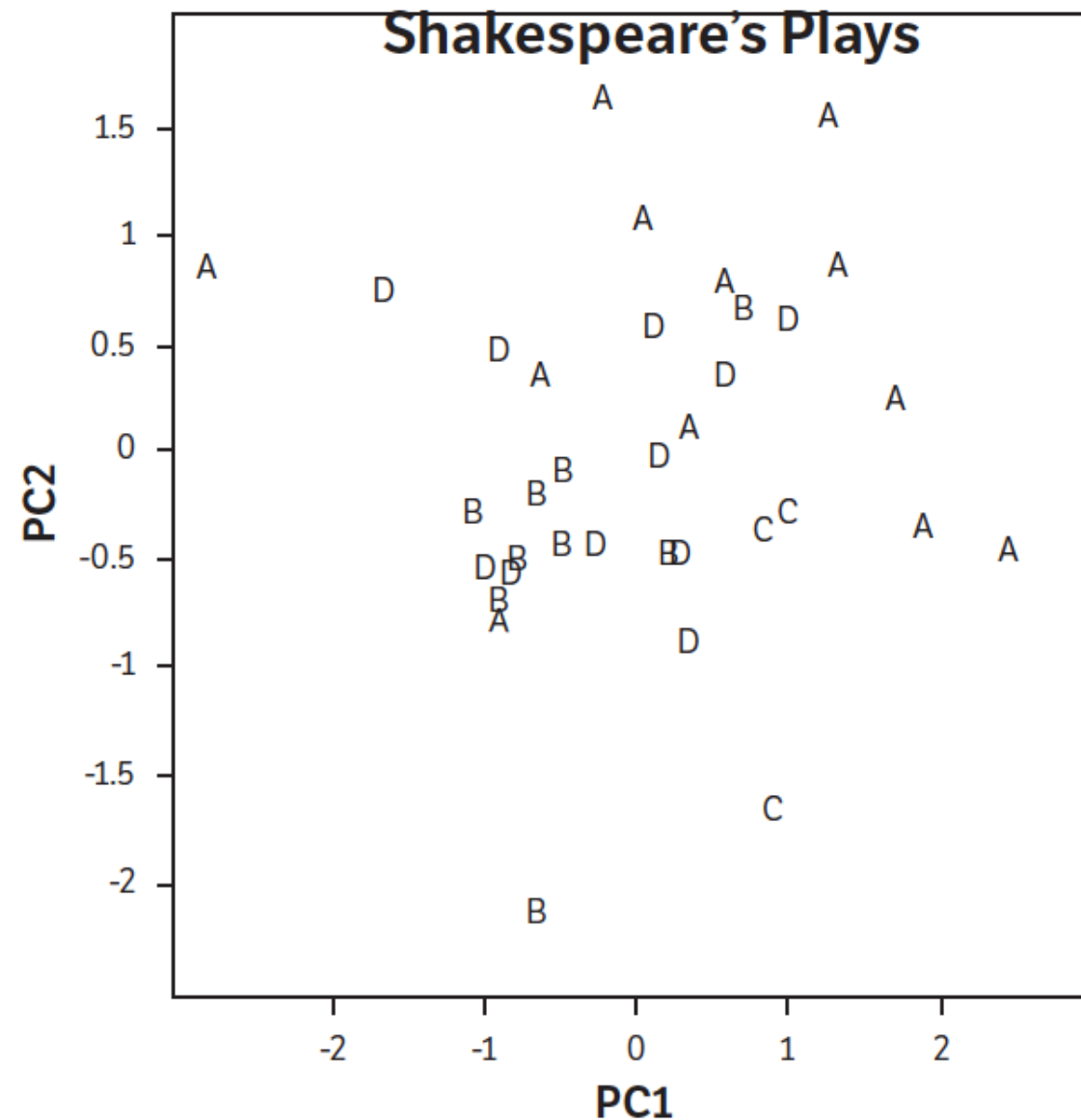


Figure 4.2: Scatterplot matrix in which Shakespeare's plays are rated on their first two principal components after having been counted by DocuScope and analyzed in terms of aggregates of LATs. PCA performed on the covariance matrix, unscaled data. Item key: A = comedy, B = History, C = Late Plays, D = Tragedies. Note how the two components place comedies in the upper right quadrant, histories in the lower left, and several late plays in the lower right (whereas tragedies, for some reason, are dispersed all over the field).

Quantitative Analysis of text

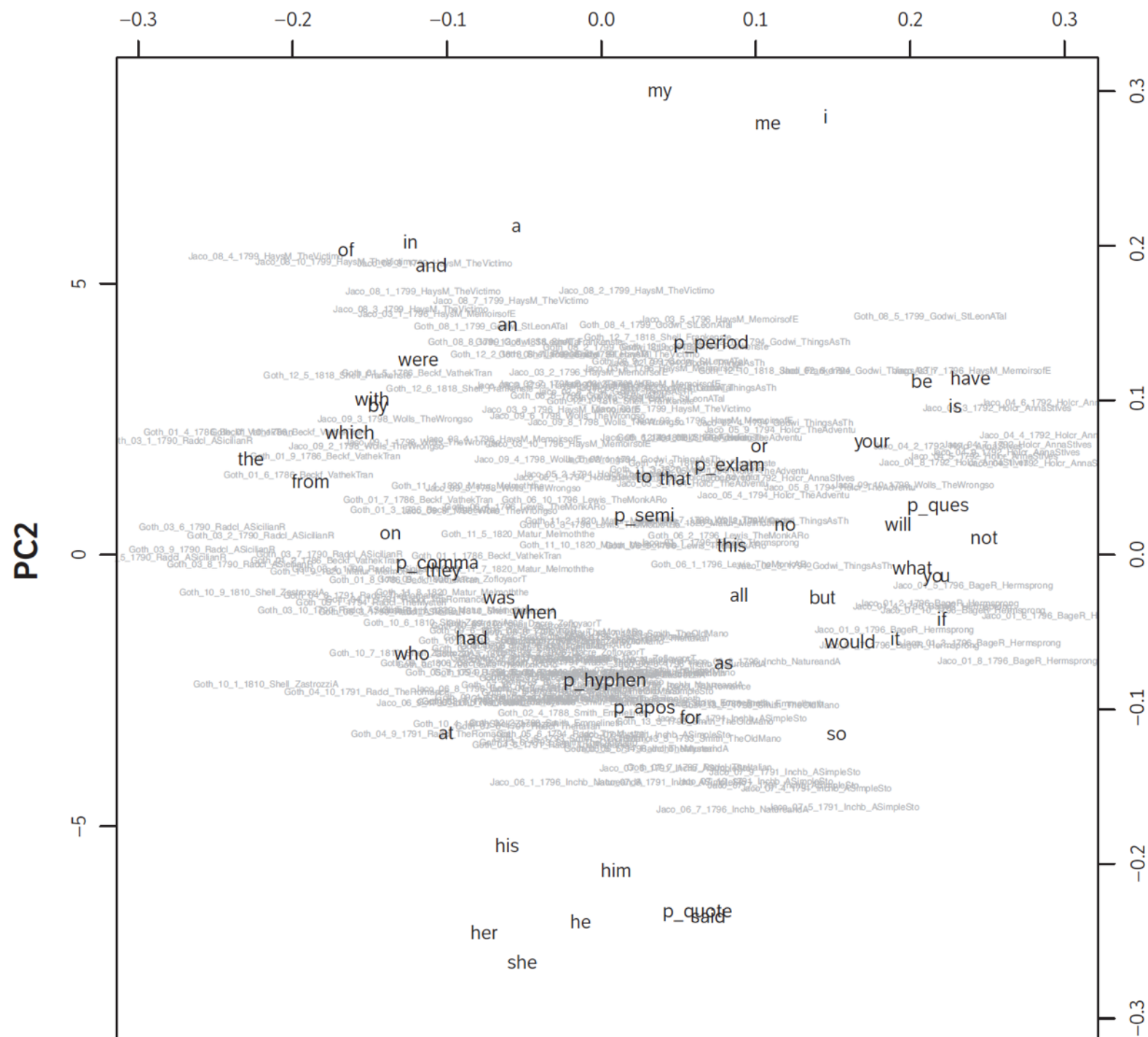
Category: genre

genres, like buildings, possess distinctive features at every possible scale of analysis: mortar, bricks, and architecture

Quantitative Analysis of text

Category: genre — 19th-century novel

Quantitative Analysis of text



Quantitative Analysis of text

Category: genre

Did we think they had produced new knowledge?

The answer, of course, was NO

Quantitative Analysis of text

Category: genre

Nothing in common in terms of units of analysis, everything in common in terms of results...

Quantitative Analysis of text

Category: genre, “the great unread”

One could give DocuScope and MFW thousands of texts of unknown generic affiliation, and see where they would fall in the gravitational field of better-known genres. One could envisage generation-by-generation maps of the literary universe, with galaxies, supernovae, black holes...

Quantitative Analysis of text

Category: genre

Roughly speaking, we found that the gothic novel averages less talk and more action than the Jacobin.

Trying to understand how a computer reads, reading into the results,

Rough speech is the explanatory version of looking for constellations the random stars

Quantitative Analysis of text

Category: genre breaks down

Separation of the genres, not genres themselves

Quantitative Analysis of text

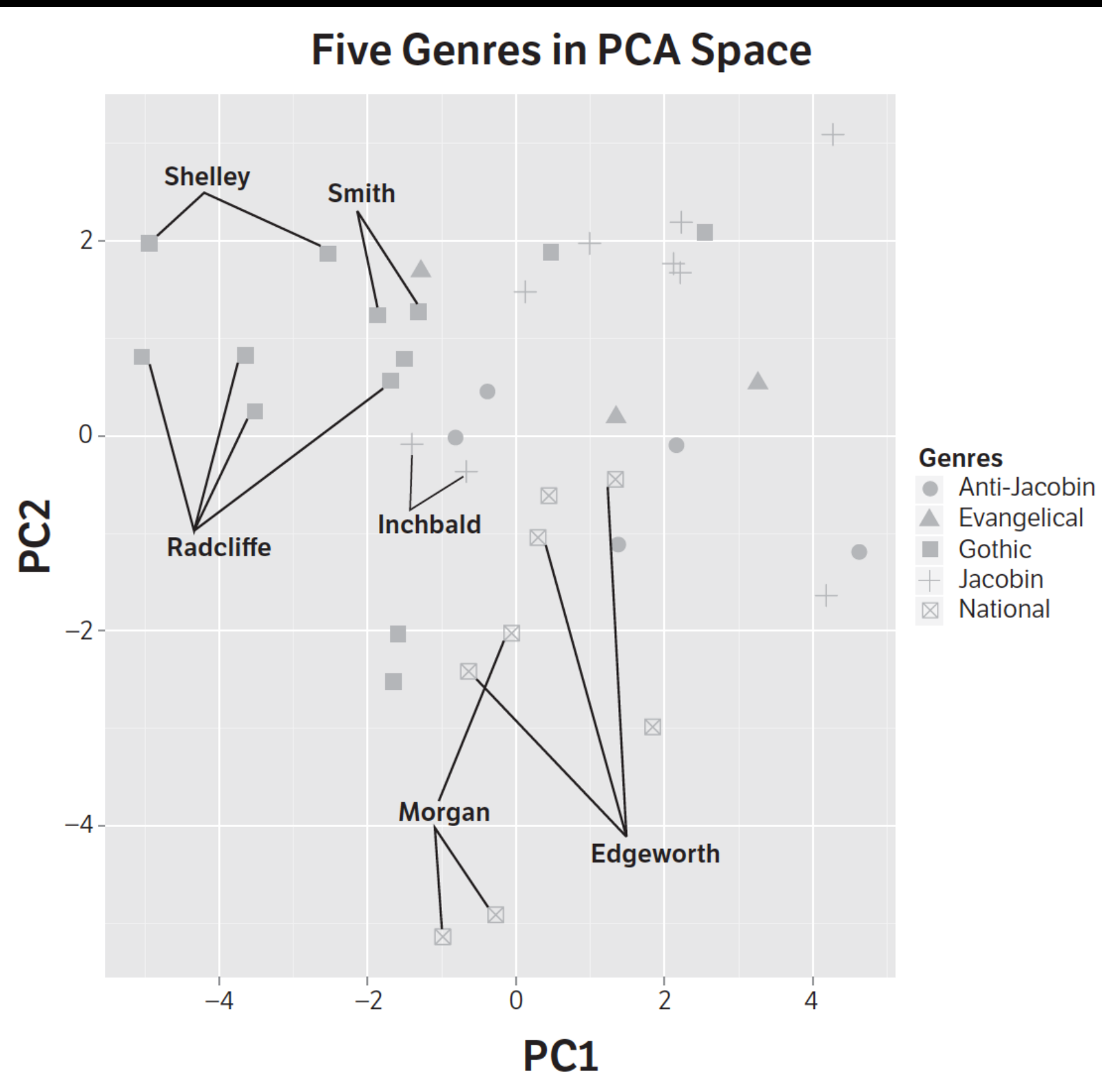
Category: genre breaks down

Why should authors be so much more recognizable than genres?

Probably, because Docuscope and MFW are very good at capturing something all writers do, whether they know it or not: using imperceptible linguistic patterns that provide an unmistakable stylistic “signature”

Quantitative Analysis of text

Category: genre breaks



Quantitative Analysis of text

Category: genre breaks down

*Why did it do so “well” with
Shakespeare and so poorly with
19th-century novels?*

Quantitative Analysis of text

Category: genre breaks down

*Why did it do so “well” with
Shakespeare and so poorly with
19th-century novels?*

Also, authors exist, but do genres?

*Oxford Shakespeare: whether a
society should value imaginative
reading*

Quantitative Analysis of text

Question of domain: taxonomy

*The question of what you were
studying--genres of the single author,
authors across time, genres across
time*

Quantitative Analysis of text

Domains: Chronology, Author, Genre

Earlobes, fingernails ... “involuntary signs”

There is a problem with earlobes and fingernails: good as they might be at identifying the author of a painting, they are worthless at explaining its **meaning**.

Quantitative Analysis of text

The experiment then turned into an **exploration**

There is something paradoxical in these traits that classify so well, and explain so little.

Assigning values/ quantities to language the only way to run computer read processes

Final projects: building a data set

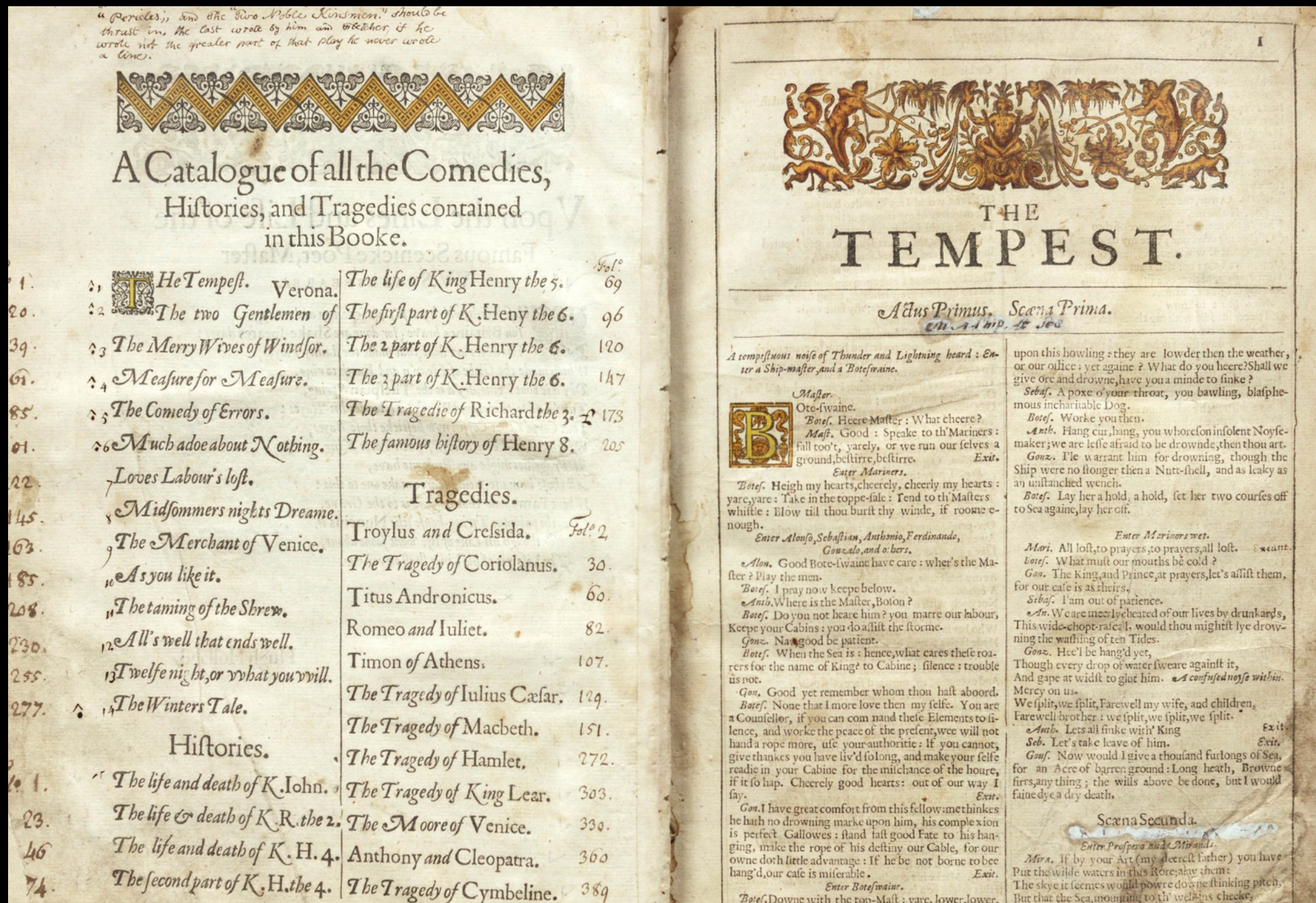
The experiment then turned into an **exploration**

Final projects: building a data set

exploratory analysis

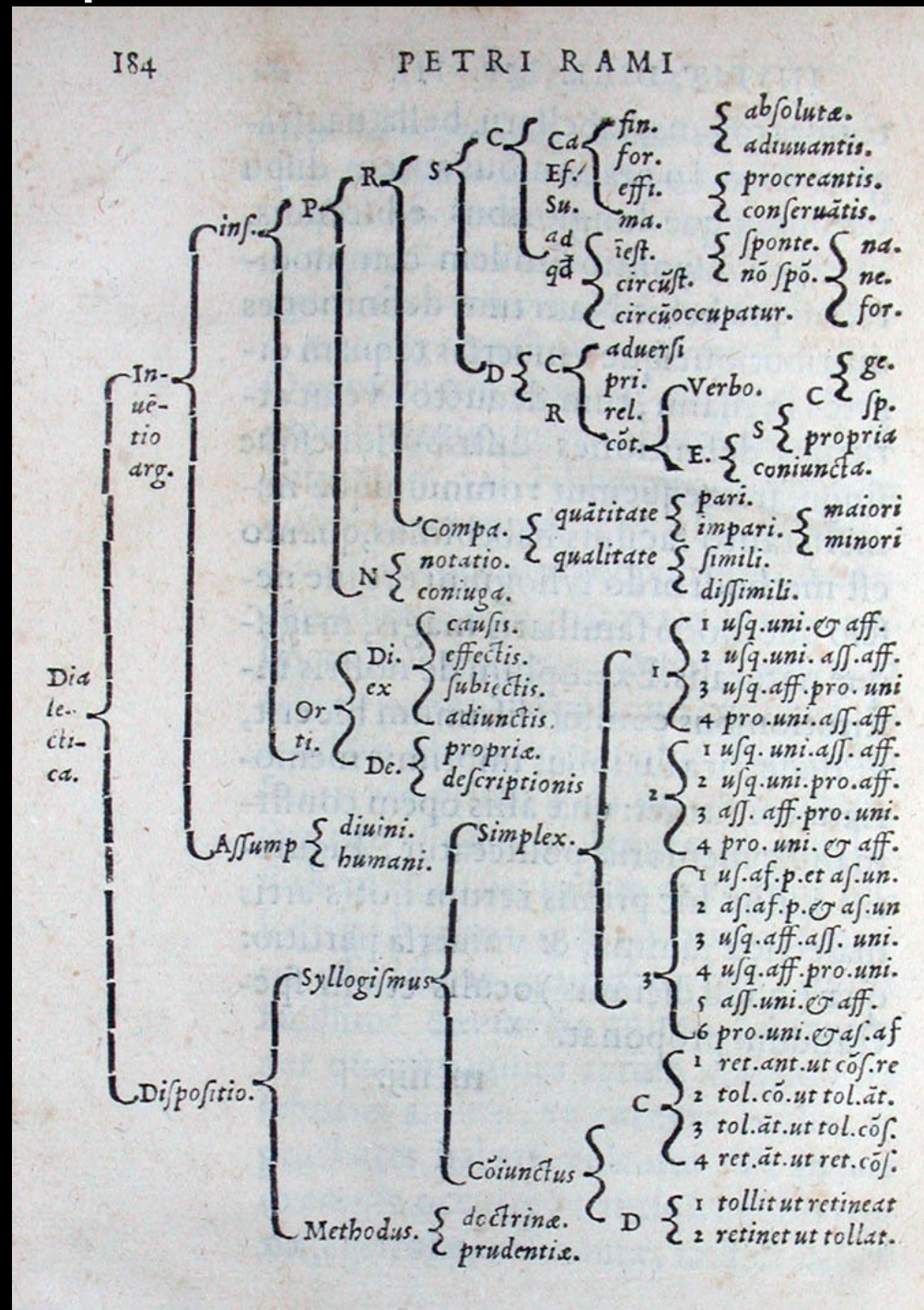
Measurements —> Meaning

1. Reading / Research: choosing your subject/domain

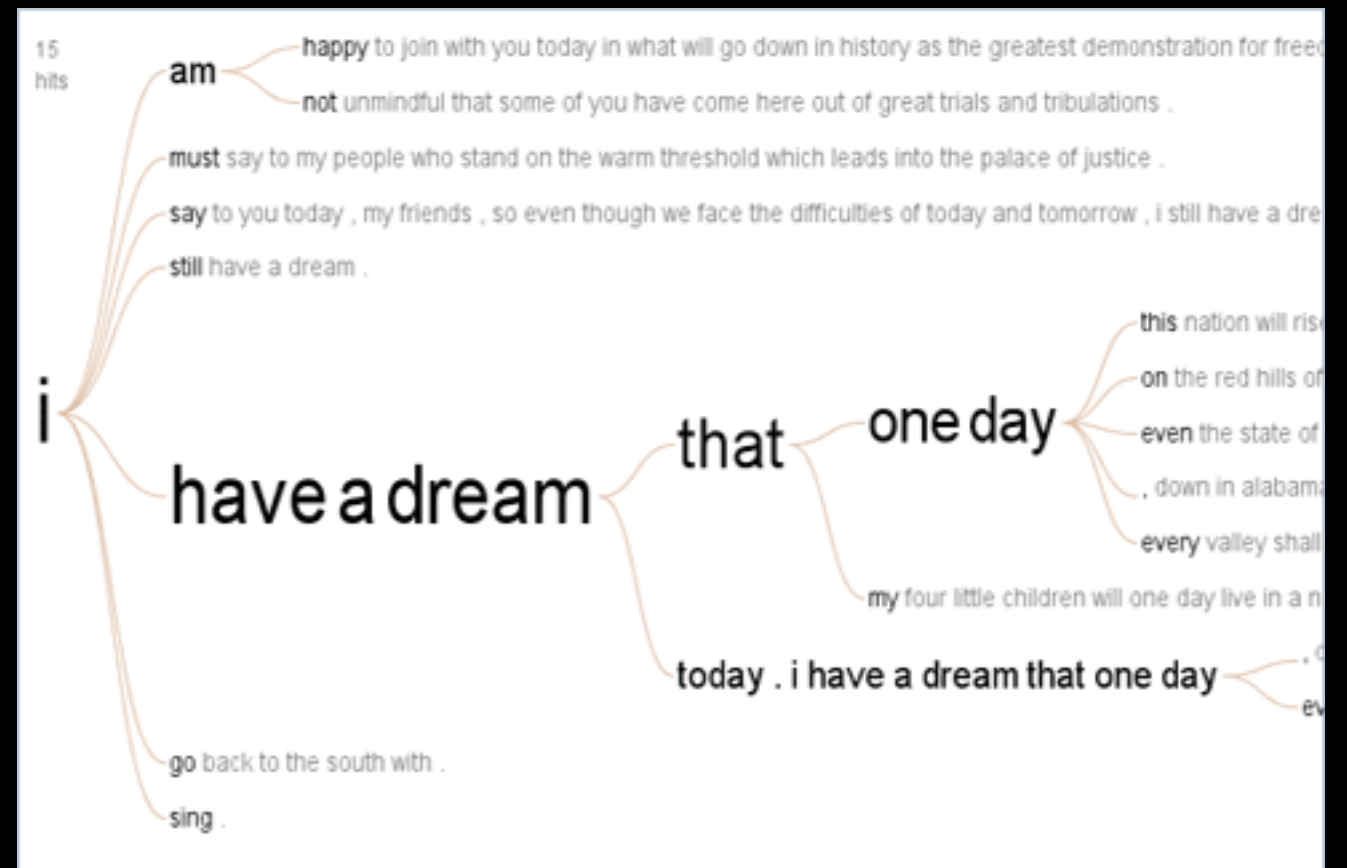


Measurements —> Meaning

2. Organizing / Taxonomies: **finding the data and imagining a shape for it**



Digitization



Left: Petrus Ramus, classification structure from *Dialectique*, 1555

Right: Fernanda Viegas & Martin Wattenberg, Word Tree 2007

Measurements —> Meaning

3. Building / Data Structures: **designing your data architecture**

```
mondial=# SELECT country, max(population)
mondial=#       FROM city
mondial=#       WHERE population IS NOT NULL
mondial=#       GROUP BY country
mondial=#       ORDER BY country
mondial=#       LIMIT 15;
```

country	max
A	1761738
AFG	2435400
AG	22219
AL	418495
AND	22256
ANG	2107648
ARM	1066264
AUS	4605992
AZ	2150800
B	507911
BD	7423137
BDS	88529
BEN	665100
BF	1475223
BG	1270284

(15 rows)

```
mondial=#
```

Measurements —> Meaning

4. Searching / Aggregating

```
File Edit View Insert Cell Kernel Help Trusted Python 3
[Icons: Save, Add, Cut, Copy, Paste, Up, Down, Previous, Stop, Refresh, Code, Keyboard]

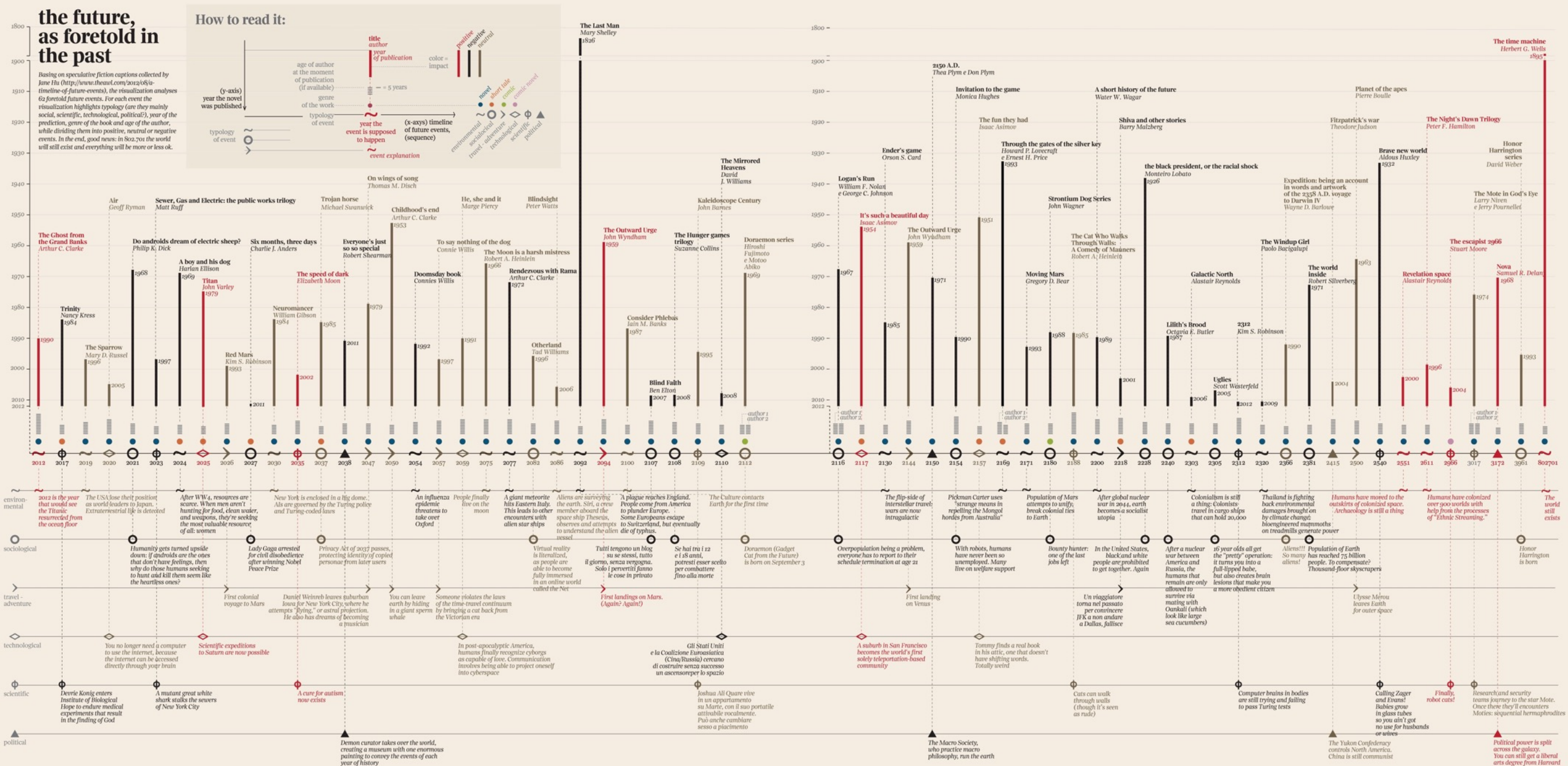
In [5]: urls = []
url_failed = []
for director in directors:
    search_d = director.lower().replace(" ", "+")
    url_string = 'http://www.imdb.com/find?&q=' + search_d + '&s=all'
    search_name = director.lower()
    raw_html = urlopen(url_string).read()
    soup_doc = BeautifulSoup(raw_html, "html.parser")
    next_url = soup_doc.find(class_='result_text')
    if next_url is not None:
        good_url = next_url.a['href']
        cleaner_url = good_url.split("?ref")
        if cleaner_url[0].startswith('/name'):
            urls.append(cleaner_url[0])
        else:
            url_failed.append(search_name)
    print(cleaner_url)
```


Measurements —> Meaning

Design/Aesthetics

5. Displaying / Generating Knowledge

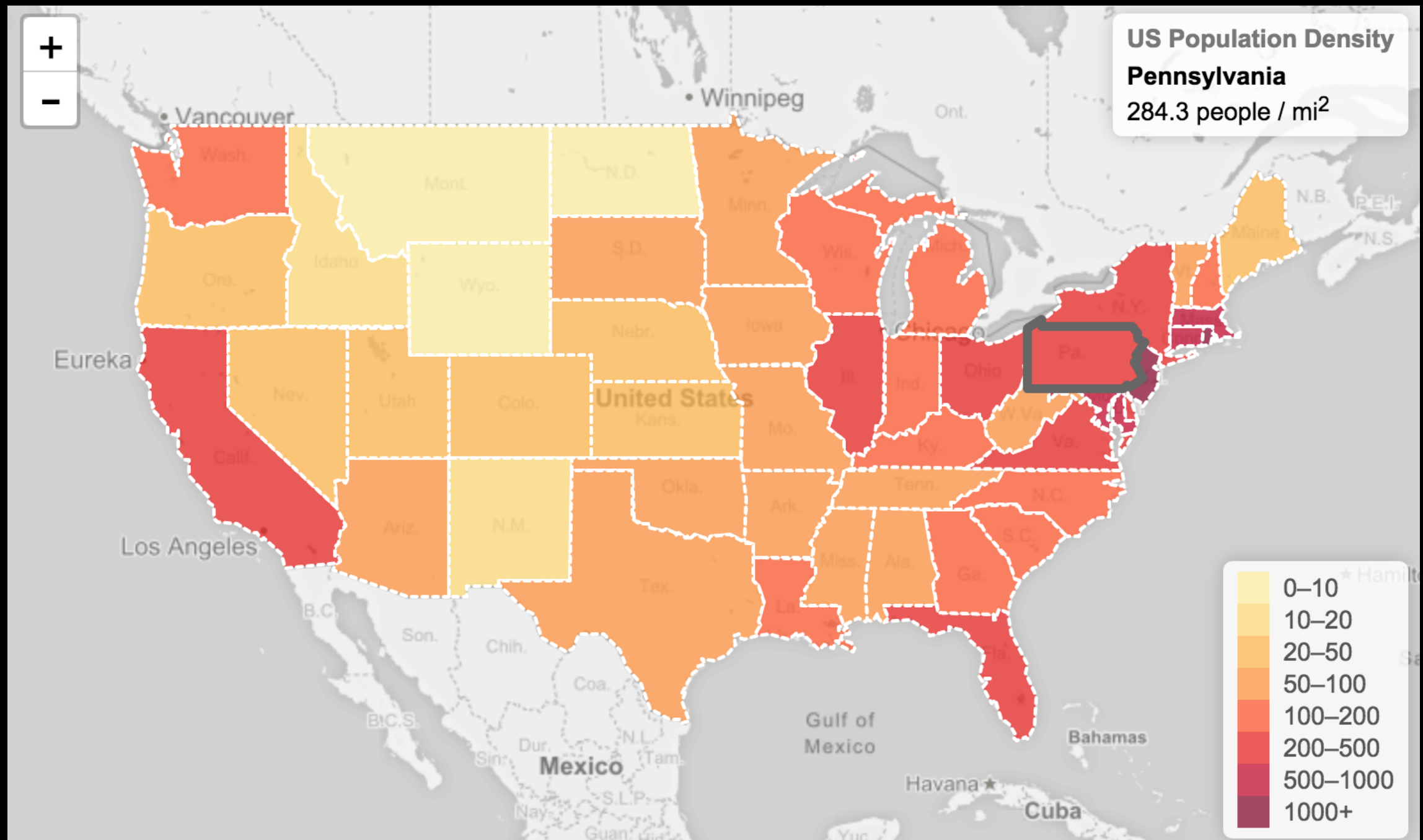
Interactivity/Screen Space



Measurements —> Meaning

Interactive map in leaflet.js

5. Displaying / Generating Knowledge



Final projects: building a data set

exploratory analysis > **map as visual presentation**

Maps are an artificial category that you must adhere to...

There may be a visual form more well suited to your subject, but we all agree the world exists...

Interrogating the form screen-based maps...

All screen/visual representations of data are a type of mapping...

Your focus will be building output for the map, not on the design

Final projects: building a data set

Choosing a subject

1. what is the subject of the project?
2. what is your main research question?
3. what is/are your data source(s)?
4. how will you transform the data into your own data set: scraping, regex, etc?
5. what will your architecture be?
6. how can this data set be geocoded?
7. on the map, what level of study would be displayed: Country, State, City, neighborhood, etc?
8. what information would be displayed when you click on/rollover a country (city, etc)?

Final projects: building a data set

Choosing a subject

1. what is the subject of the project?
2. what is your main research question?
3. what is/are your data source(s)?
4. how will you transform the data into your own data set: scraping, regex, etc?
- 5. what will your architecture be?**
6. how can this data set be geocoded?
7. on the map, what level of study would be displayed: Country, State, City, neighborhood, etc?
8. what information would be displayed when you click on/rollover a country (city, etc)?

Final projects: building a data set

Levels of information the map can contain (properties)

Geometry--Points/Shapes: these are combinations of latitudes and longitudes that leaflet will translate to the screen for you.

Colors/size: these are the obvious ways to visually separate different points or shapes.

Groups: Groups allow you to have different layers of shapes/ points visible.

Headline/Lede: when you roll over a point or shape this is the most immediate information that you want to convey

Article: this is for deeper reading/analysis-- on click you can get as much information as you want to output on the page

*See sample simple maps and geojson output
—note they are intentionally unimpressive*