

Computational Methods for Linguists

Ling 471

Olga Zamaraeva (Instructor)
Yuanhe Tian (TA)
04/01/21

What's your most/least favorite April 1st joke?

Class policies

(Forgot to mention last time)

- Academic integrity:
 - See the syllabus!
 - No plagiarism
- Assignments:
 - To be completed individually
 - No copying from websites
 - We will talk about how to use e.g. StackOverflow



Reminders

- Patas cluster
 - Request access
- “Assignment 0”
- Public website:
 - Assignment descriptions
 - Schedule
- Canvas
 - Announcements
 - Introductions and other discussion
 - Homework submission
 - Blog options



Ling471 & Ling472

Compared

- Ling472 (Intro to CompLing)
 - Usually harder for those without programming experience
 - But this year, flexible individual goals are allowed
 - Covers statistical and formal methods in CompLing
 - Organized around levels of linguistic structure
- Ling471 (Comp. Methods for Linguists)
 - Lots of programming etc. basics
 - Focuses only on statistical methods
 - Focuses on data science and how it is related to CompLing
 - Organized as incremental skills building
- Courses may be taken concurrently
 - There will be overlap, e.g. Machine Learning lectures
- Both instructors will collaborate on connecting the courses



Plan for today

Conceptual Overview

- Data science
 - Data
 - Language data
 - Statistical methods
 - Machine learning
 - Programming
 - Systems

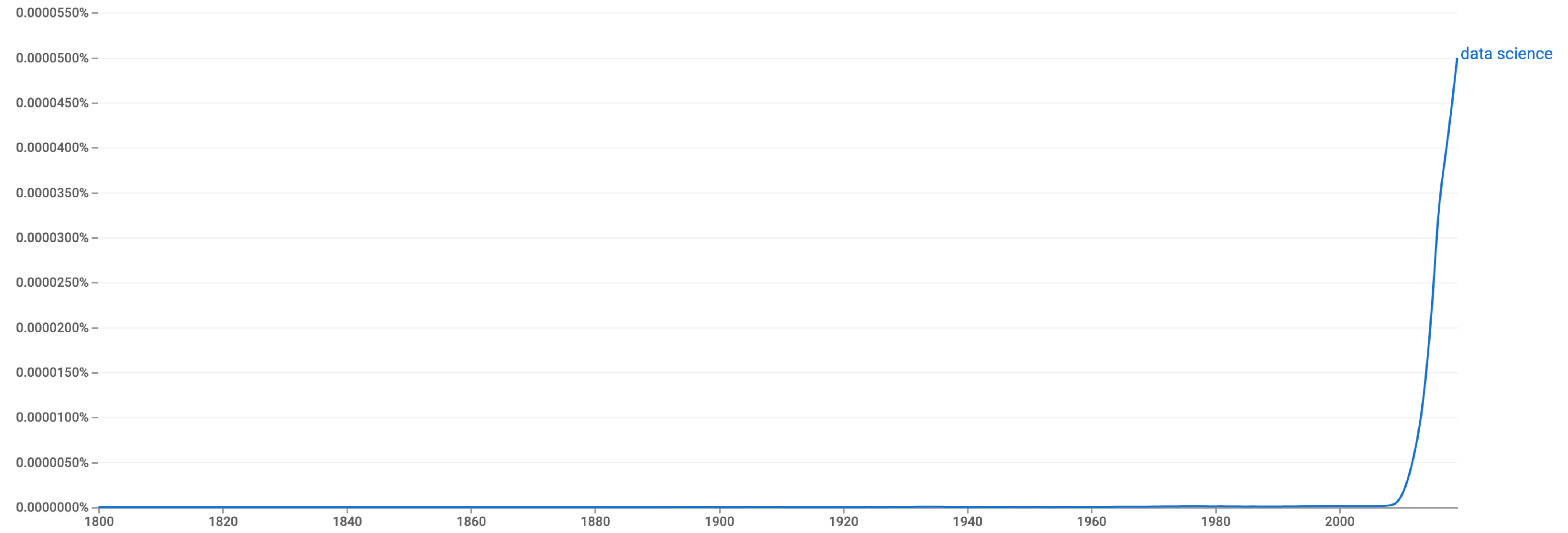


What is Data Science?

Google Books Ngram Viewer

data science

1800 - 2019 English (2019) Case-Insensitive Smoothing



(click on line/label for focus)

Activity

<https://PollEv.com/olzama> (NB: it splits by word!)

W


What is Data Science?




Total Results: 18

Wikipedia:

“Interdisciplinary field”





Data science 


Field of study


Data science is an inter-disciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from many structural and unstructured data. Data science is related to data mining, machine learning and big data.
[Wikipedia](#)

Ms salary: \$92,500 [forbes.com](#)

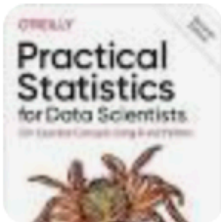




Mathematics 

Career path 

Skills 

Employment outlook 

Data science books View 45+ more

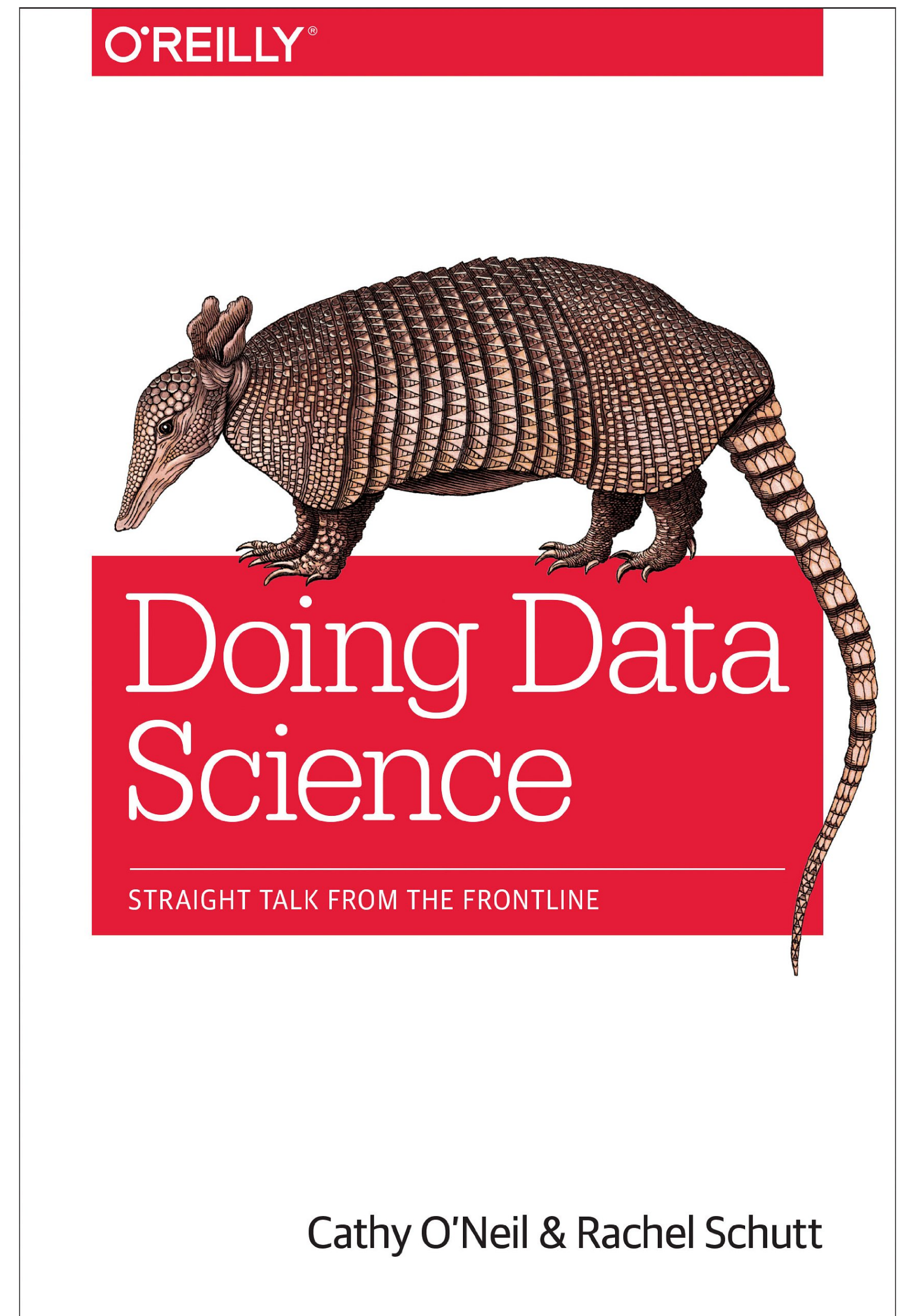
- 
Practical Statistics for Data...
- 
An Introducti... to Statisti...
- 
R for Data Science
- 
Python Data Science...
- 
The Elements of Statisti...

O'Neil & Schutt:

Doing Data Science

*Data science is a **practice** where people with different types of expertise come together to solve a problem through analyzing large amounts of data.*

(NB: This is not our textbook or anything! You don't need to buy it.)



Compare:

- *Interdisciplinary field that uses ... algorithms...to extract knowledge*
- *A practice when different specialists solve a problem through data*
- How do these definitions differ? Which one do you like more? Why?



Another definition

(by Denis Altudov)

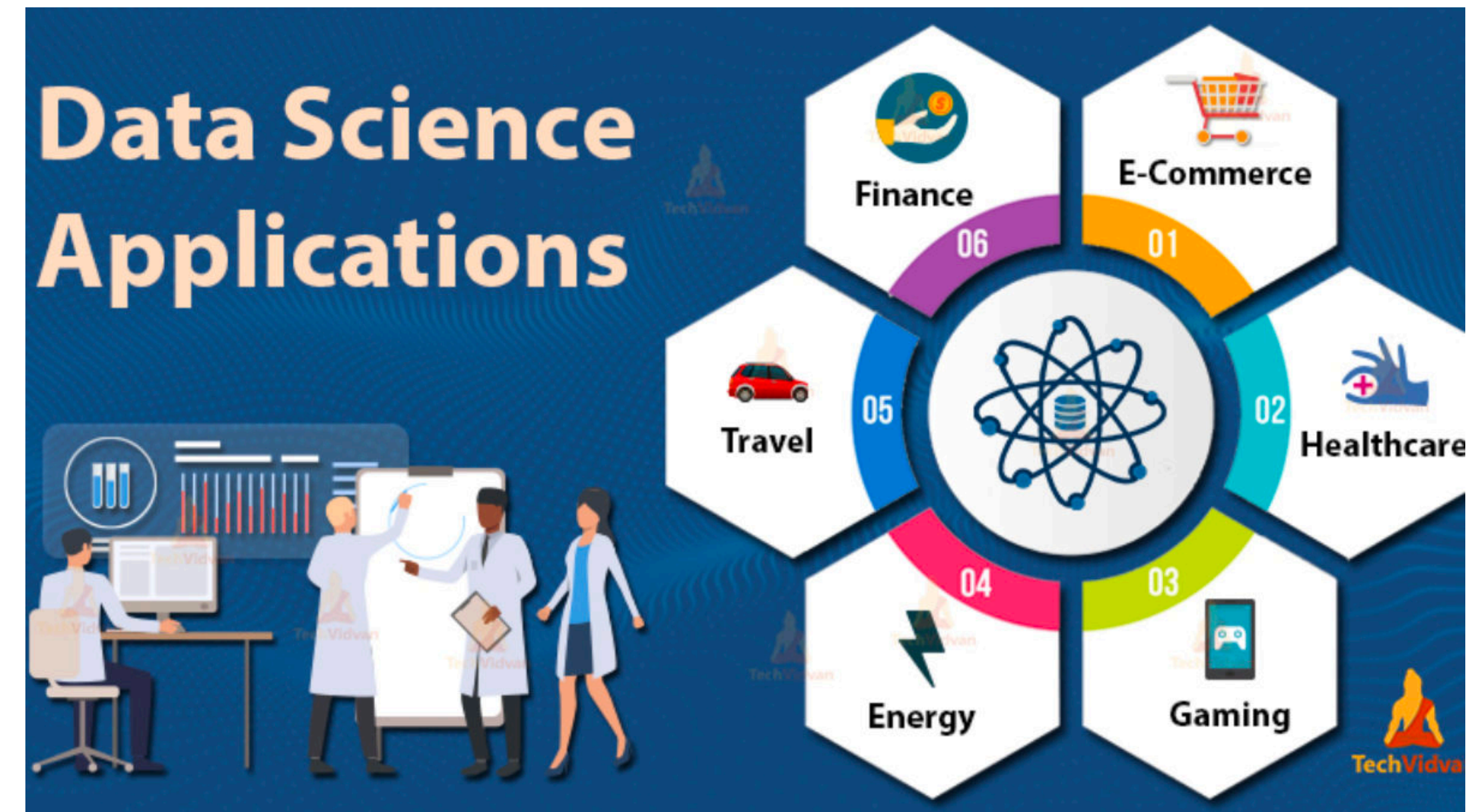
“Data science is application of rigorous statistical techniques to fuzzy objects of study”



Data Science

Examples

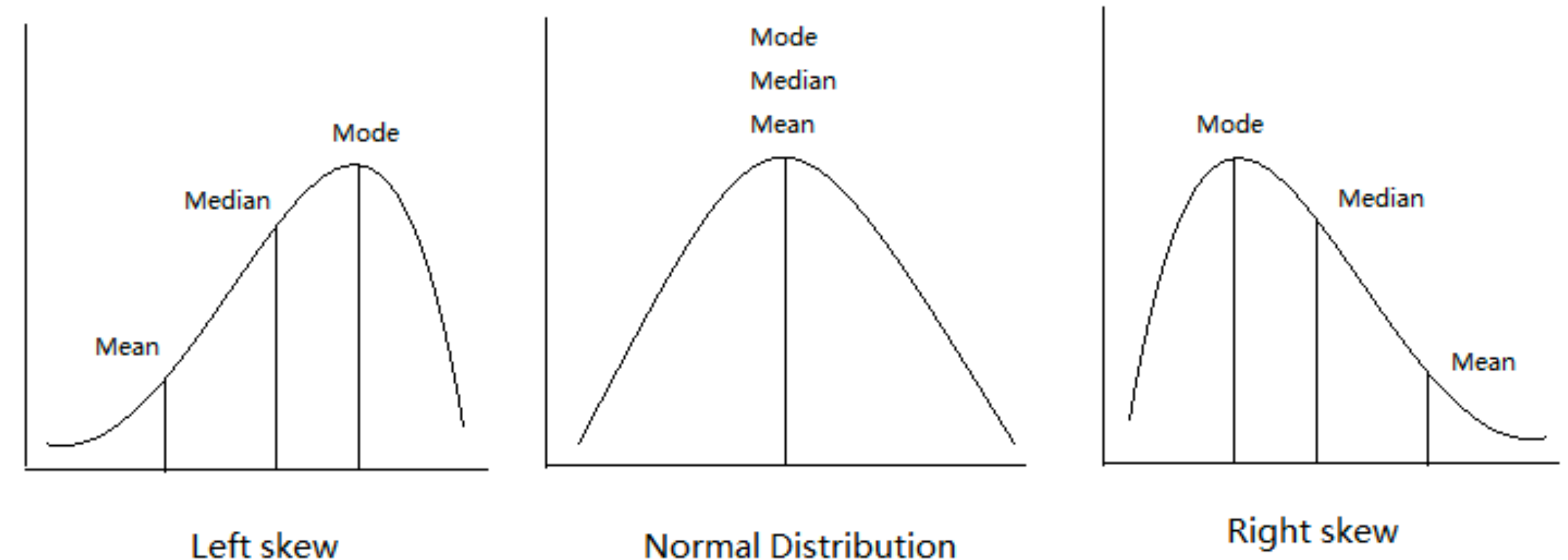
- Analyze data **statistically**
- Predict something
 - e.g. how people will behave
 - how something will impact people
- Retrospectively:
 - What happened in the past
 - What it means for the future



<https://techvidvan.com/tutorials/data-science-applications/>

Statistical analysis

- Data science is **inherently** statistical
- Statistics interprets quantitative data
 - Basic examples:
 - Mean and Median
 - “What’s the average salary in Seattle?”
 - “Are women more talkative than men?”
 - What about **qualitative** data?
 - Qualitative **methods** play a role
 - Data tends to be **digitized**
 - But the **nature** of data may remain **qualitative**



<https://medium.com/@nhan.tran/mean-median-an-mode-in-statistics-3359d3774b0b>

What we need to do data science:

- Datasets
 - BIG ones!
 - Why?
- Domain expertise
 - E.g. linguistics
 - This is **the whole point**
- Statistics (and probability theory!)
 - Statistics **interprets quantitative data**
- Programming
 - **Large** quantities of data cannot be processed with pen and paper
- Visualization and communication
 - Because **business**
- But wait! What about machine learning and deep learning?!
 - Yes, but they apply as **statistical methodologies**



Data science profiles

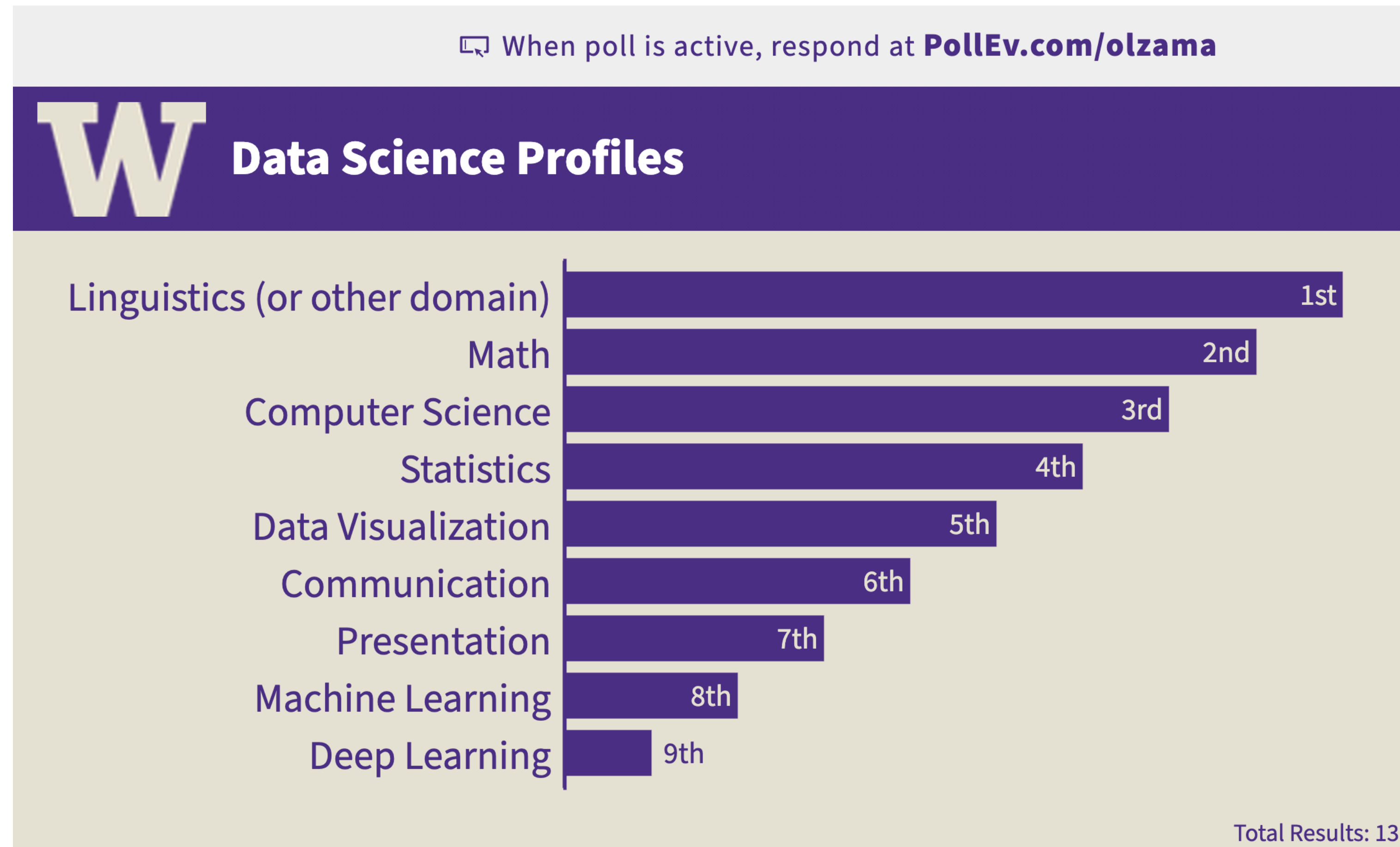
Activity

- Data science requires **different** skills
- Rank yours in relative terms of strength, to **visualize** your data science “profile”
- NB: This is your profile **the way it actually is**; not what you think it “should” be!



Our collective Data Science profile:

<https://PollEv.com/olzama>



Data and Datasets

- What counts as “data” in Linguistics?



Activity

<https://PollEv.com/olzama> (NB: it splits by word!)

W

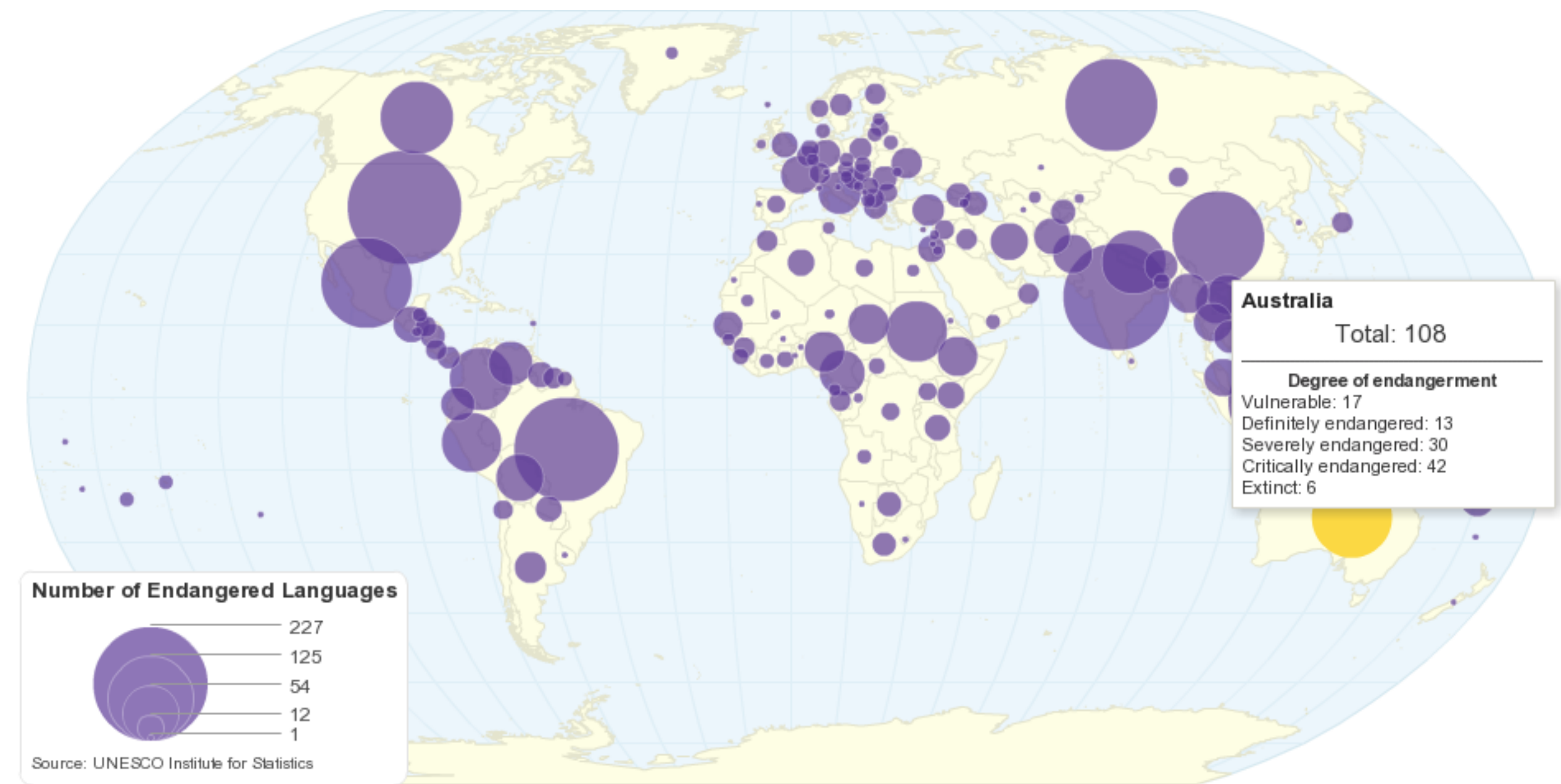
What counts as "data" in linguistics?

A word cloud of terms related to linguistics data collection. The words are arranged in a roughly circular pattern. The most prominent words are 'speech', 'corpus', 'interviews', 'text', and 'people'. Other visible words include 'eye-tracking', 'brain waves', 'sentences', 'eye-movement', and 'comprehension'. The words are in various shades of brown and tan.

Total Results: 16

Data In Linguistics

- Recorded speech and text
 - Why in this order?
 - Most languages of the world aren't written!
- Is text more structured than speech?
- Is it easier to work with text or with speech?
- What does it mean to focus on text vs speech?



<http://chartsbin.com/view/1339>

Data in linguistics and data science

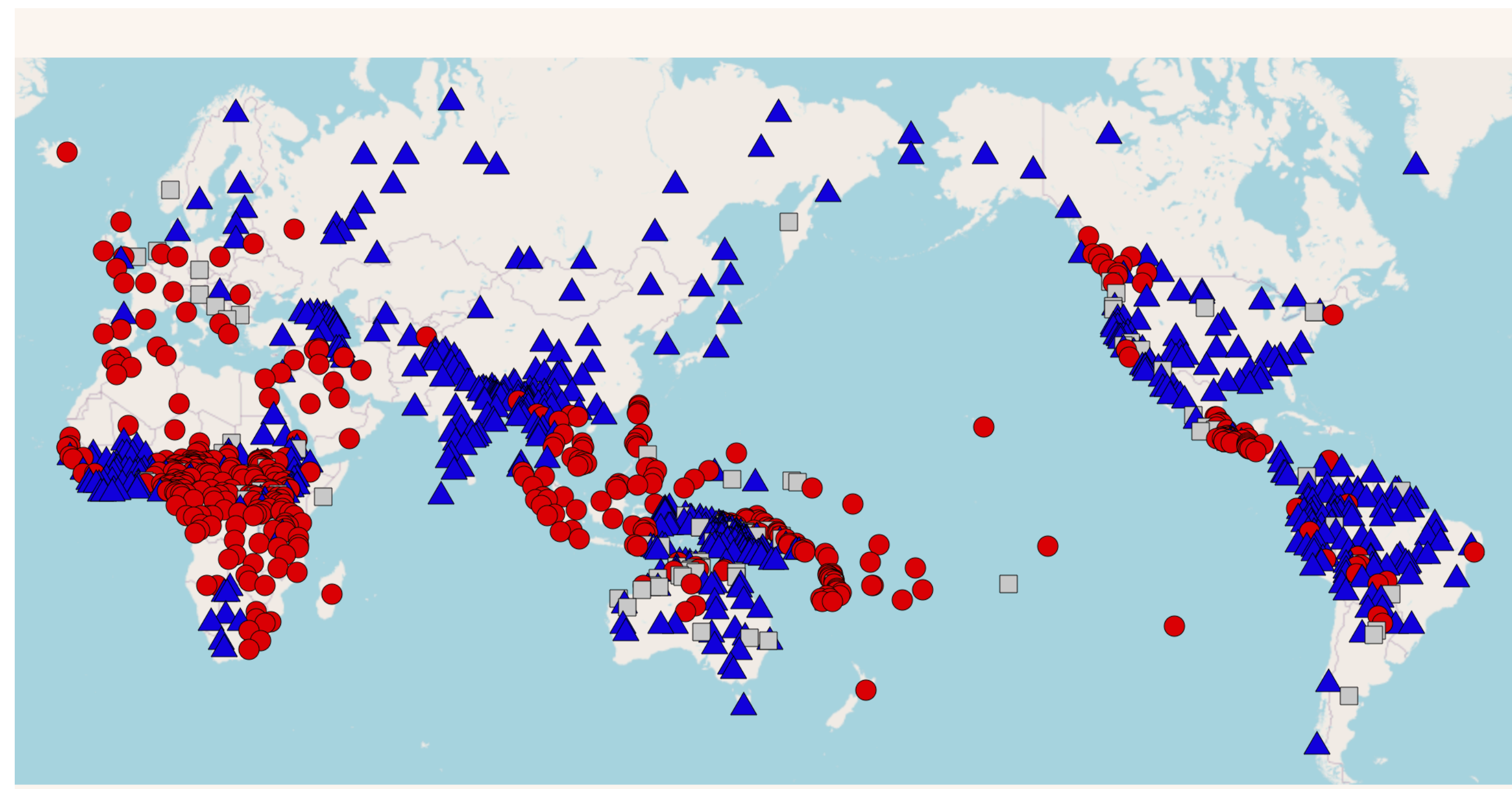
- Data science usually means LOTS of data
 - Why?
- Which areas of linguistics have LOTS of data?
 - **Most of them**
 - ...potentially



<https://depts.washington.edu/ldplab/>

Linguistics and data science

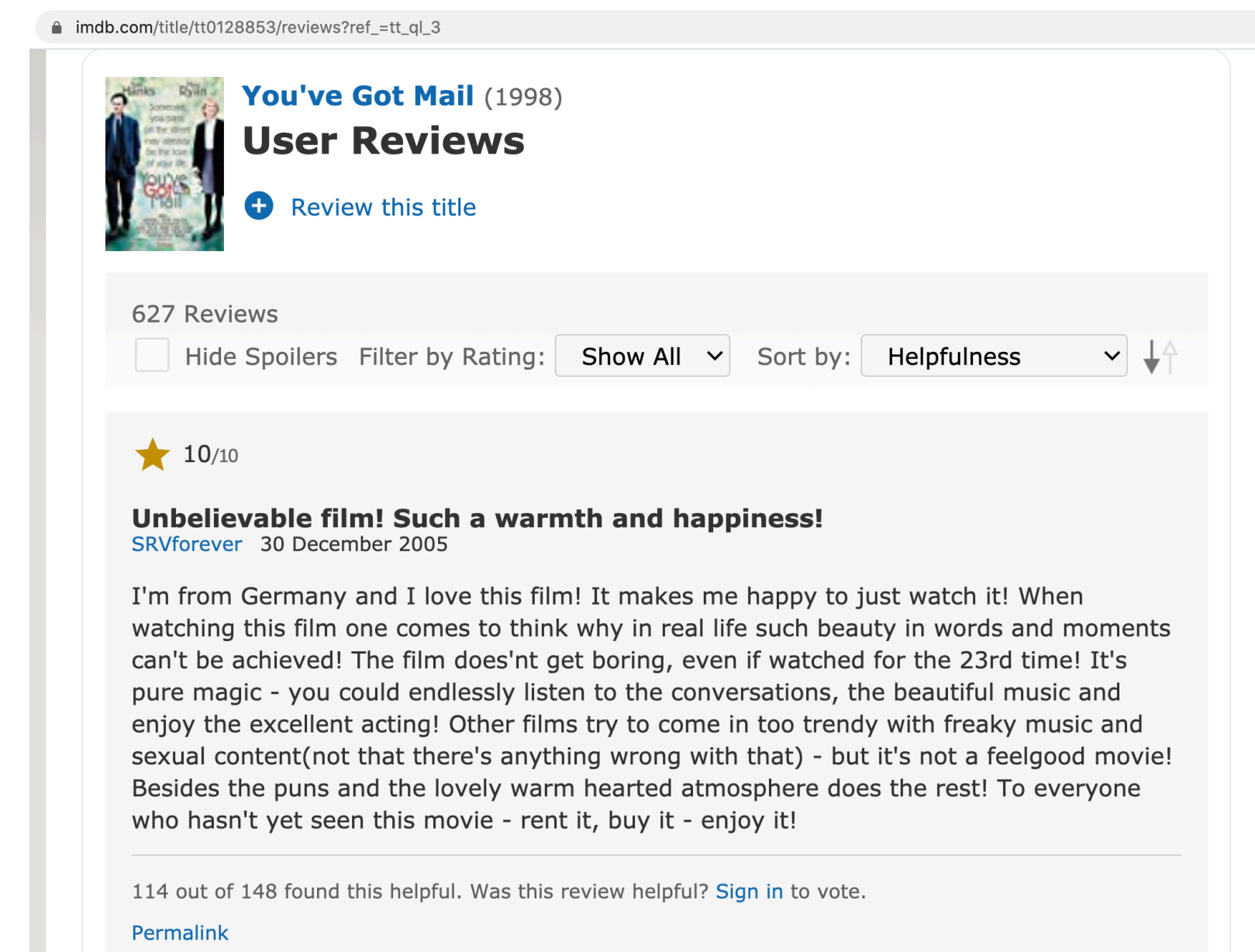
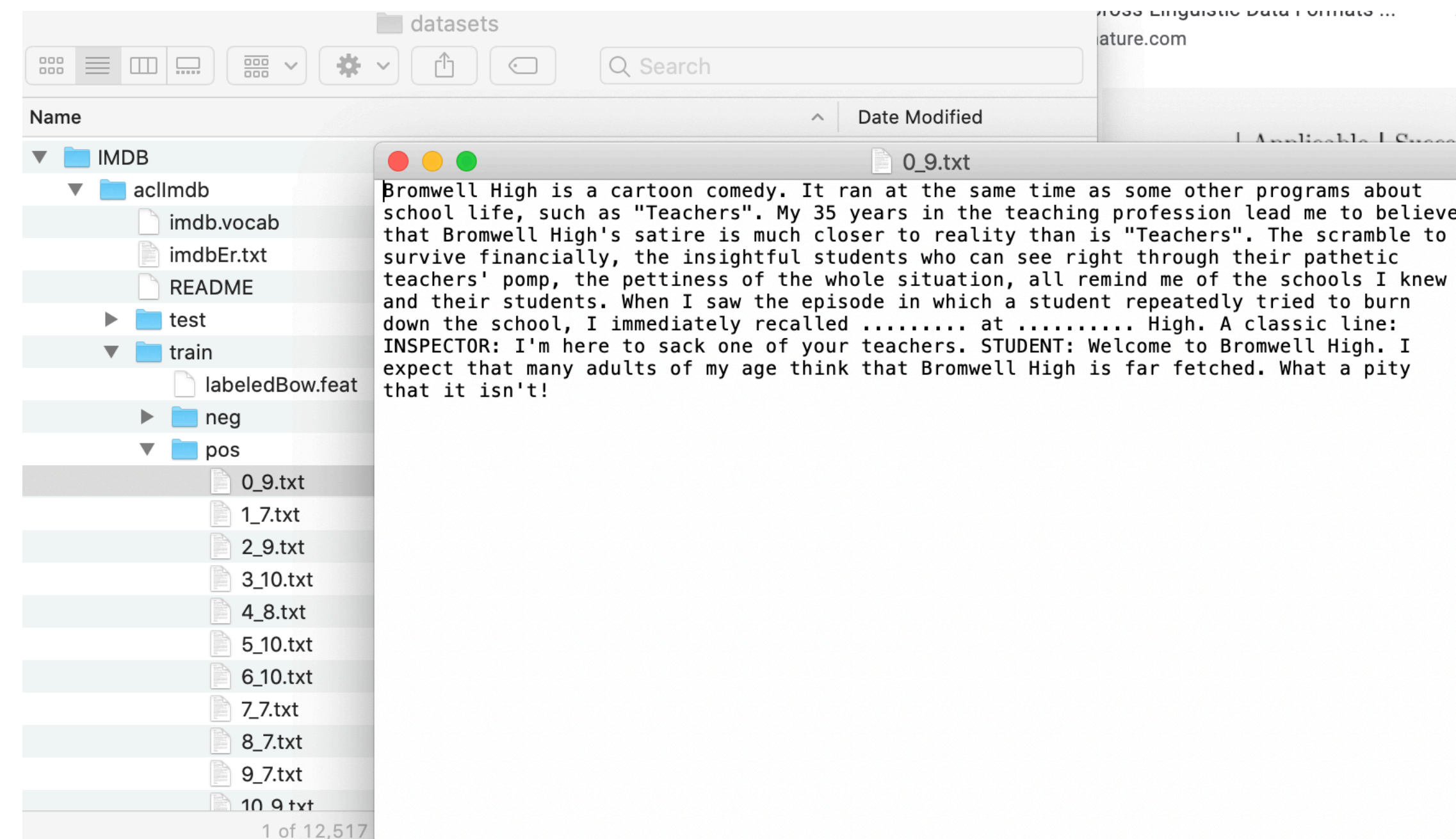
- “Corpus linguistics”
 - Various subfields; statistical analysis over large texts
- Sociolinguistics
 - Statistically significant correlations between sociolinguistic variables
- Historical (“dyachronic”) linguistics
- Linguistic typology
- What else?
 - Almost **everything, potentially**
 - So long as the data can be **managed**



<https://wals.info/feature/86A?v1=t00d&v3=sccc#2/21.0/152.9> (Dryer, 2005. WALS. Order of Genitive and Noun)

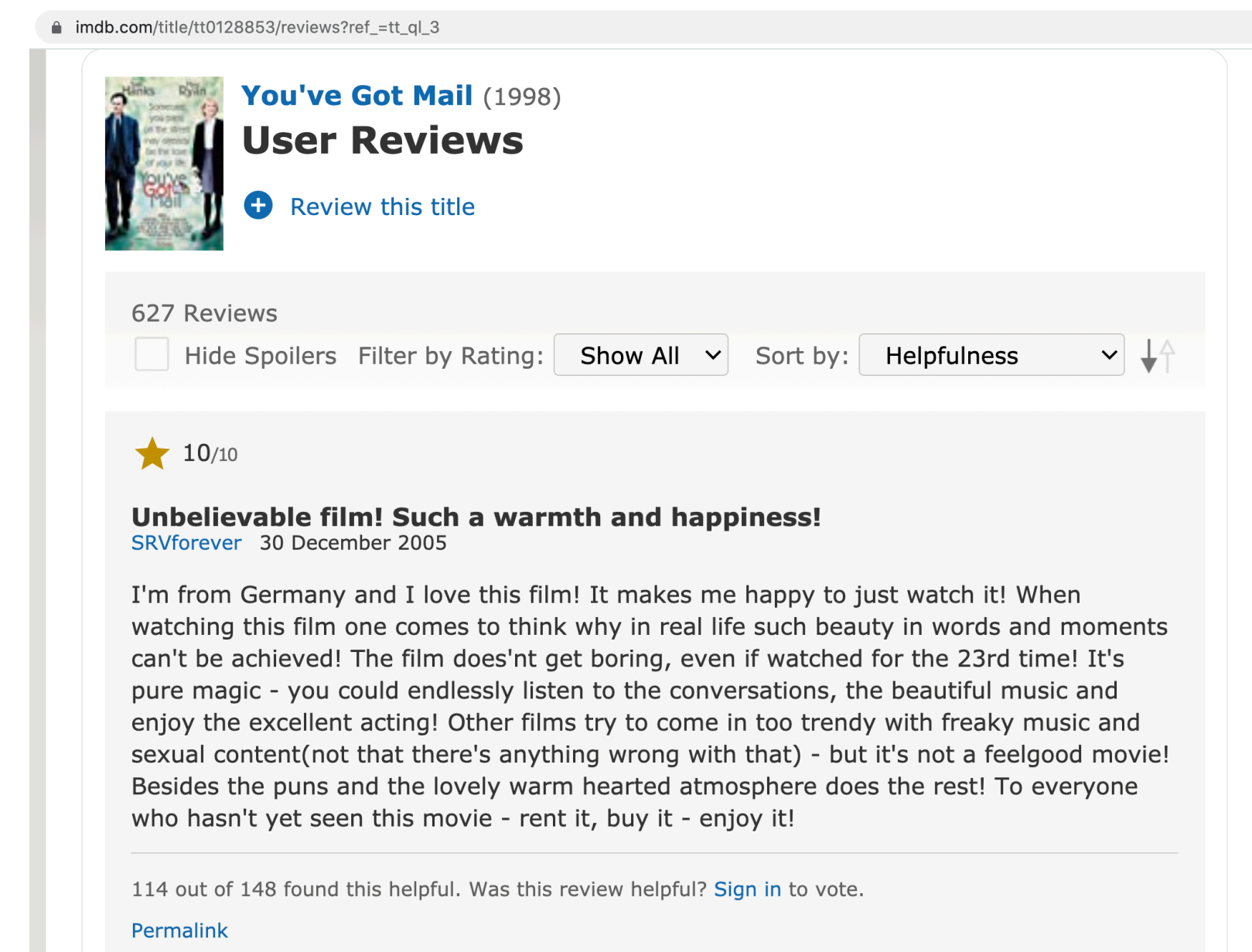
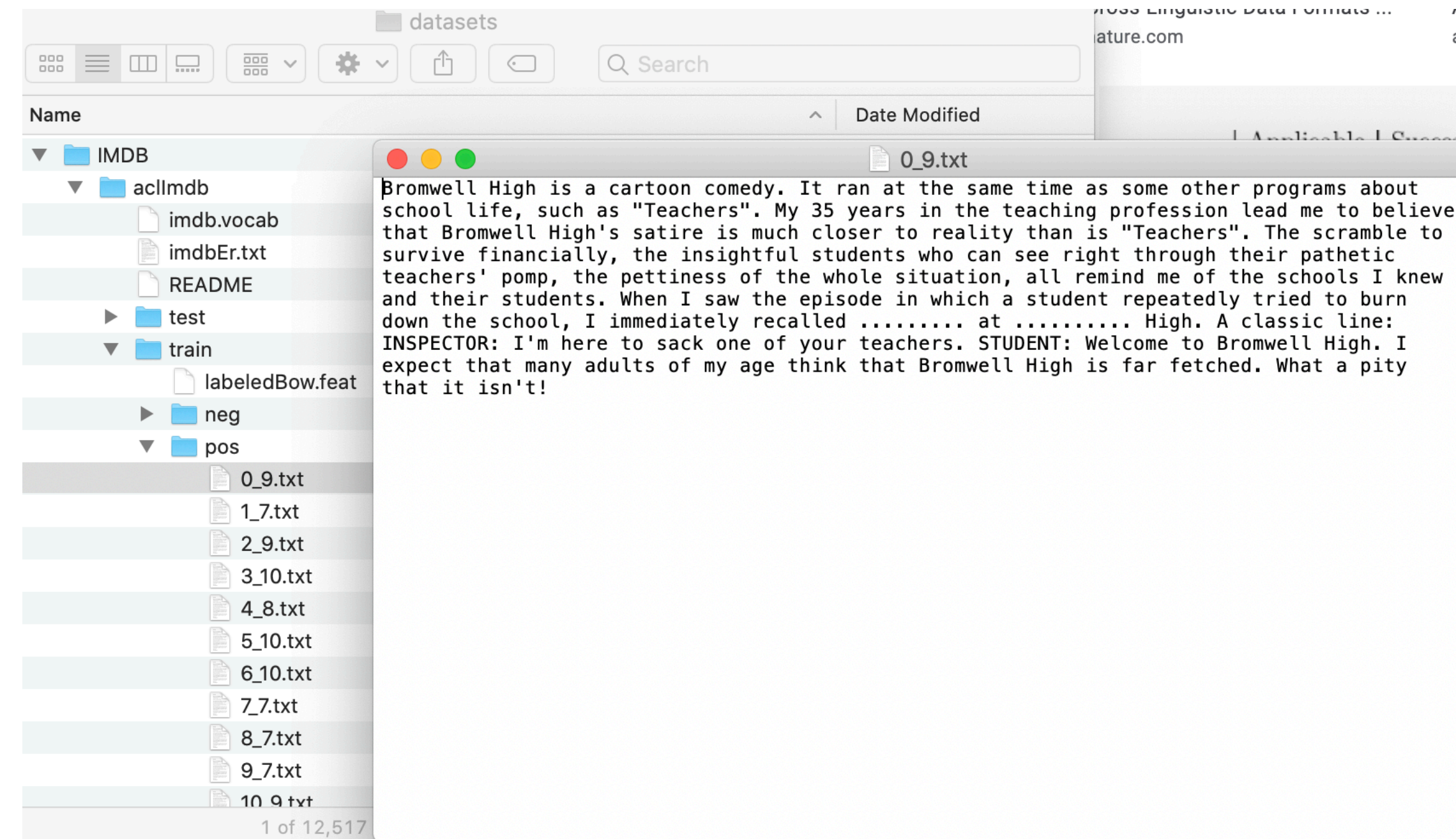
Datasets

- Datasets are sets of data :)
- No really! Anything that's data is a dataset
- But, in practice, datasets are **workable** data
 - Size, annotation, structure, consistency, noise
 - Some datasets are used again and again in research
 - The IMDB reviews dataset!
 - That's why, datasets should be and are recognized as a **contribution to science**



Datasets

- Datasets take many forms
- But they all have **provenance**
- **Language** datasets are produced by people
- People have age and a variety of identities
 - ...influencing **how** they talk and write
- The provenance of the dataset biases the system that's based on it
 - Can introduce **additional** biases to "balance it out"
 - But **no such thing** as an unbiased system



Data statements

Bender & Friedman 2019

- Reading assigned “for April 8”
 - But will discuss now (not on April 8)
 - Read it soon (to do Assignment 1)
- A **practice** of stating facts about the dataset used to “train” the system
 - “Training” is a machine learning concept
 - The system “trains” by getting feedback on how well it did something **with respect to the data**

Data statements alone won't 'solve' bias, but if we do not make a commitment to data statements or a similar practice for making explicit the characteristics of datasets, then we will single-handedly undermine the field's ability to address bias.



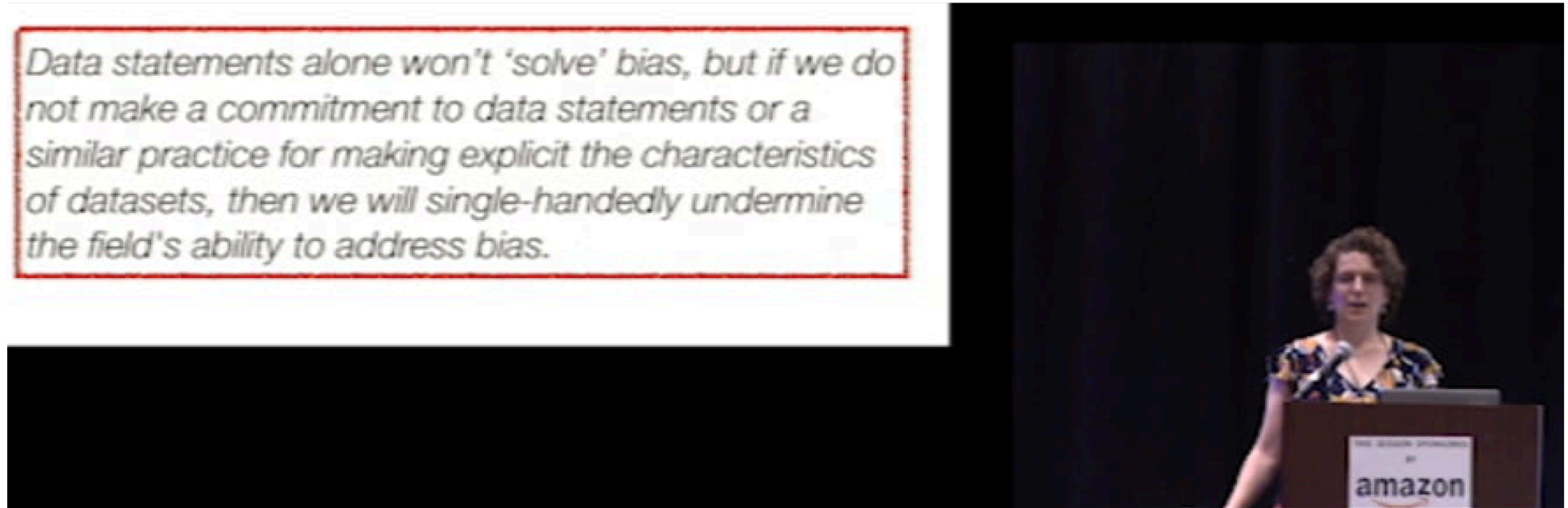
Emily M. Bender giving a talk on Data Statements for NLP in 2019 at NAACL

Data statements

Bender & Friedman 2019

- Schema:
 - Curation rational
 - How were data selected and why
 - Speaker demographics
 - Annotator demographics
 - Many datasets are **annotated** in some way
 - e.g. positive or negative review?
 - Could the judgment here differ?!
 - Speech situation
 - Text genre etc.
 - Recording quality

Data statements alone won't 'solve' bias, but if we do not make a commitment to data statements or a similar practice for making explicit the characteristics of datasets, then we will single-handedly undermine the field's ability to address bias.

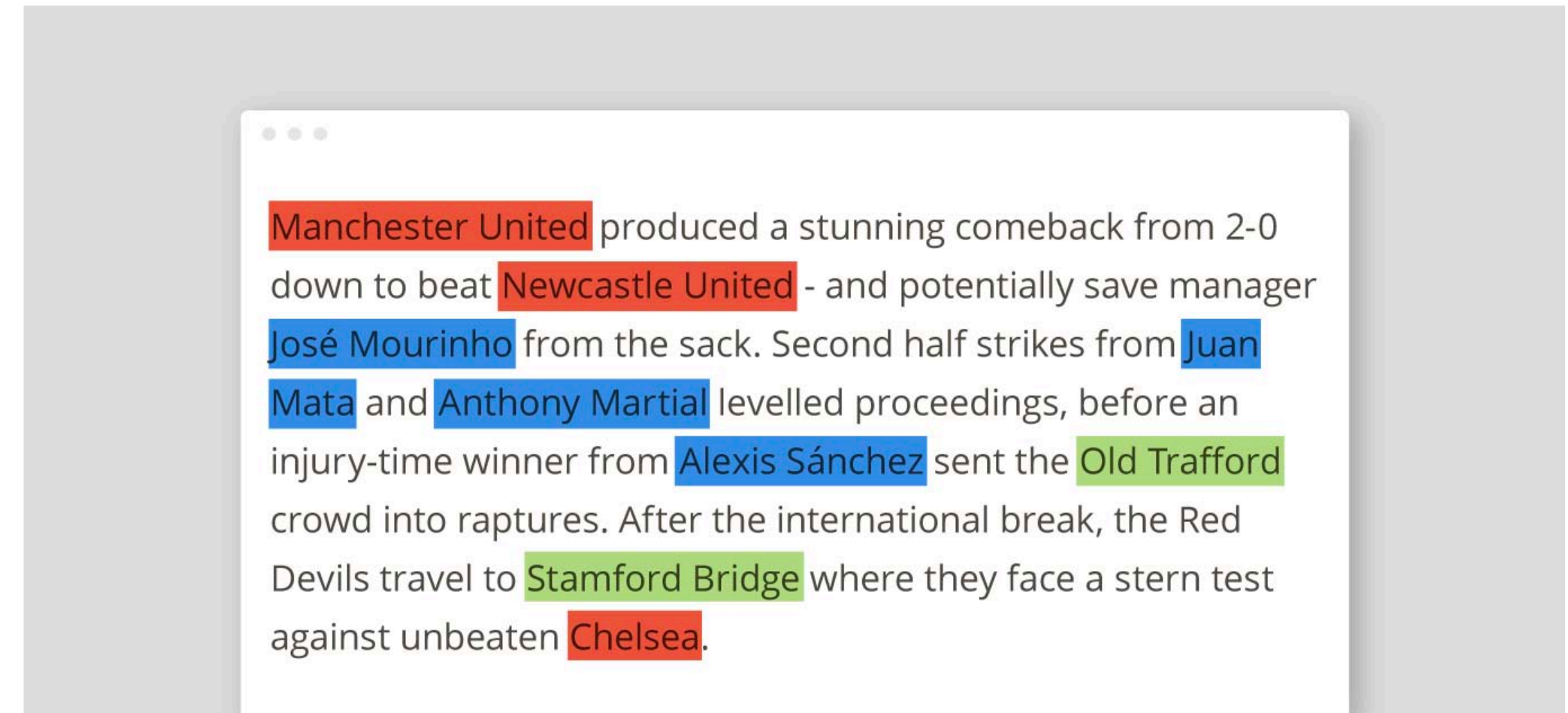


Emily M. Bender giving a talk on Data Statements for NLP in 2019 at NAACL

Annotated data

- Speaker characteristics (e.g. region)
- L1/L2 speech
- Part of speech
- Named entity
- Sentence structure
- Discourse structure

- Some of the above are what we want the system to predict; other things are useful **downstream**



<https://lionbridge.ai/articles/training-data-is-key-for-natural-language-processing-algorithms/>

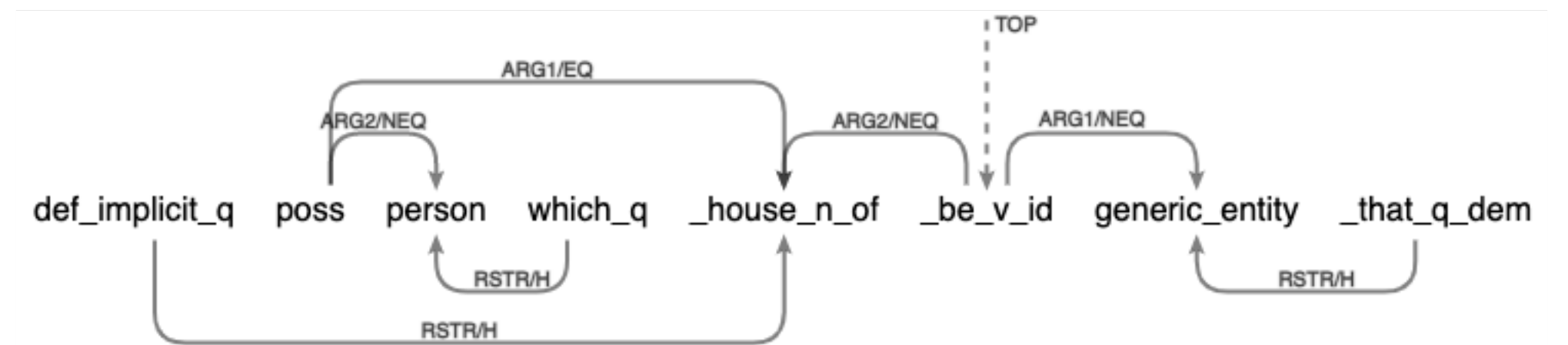
Annotated data

In linguistics

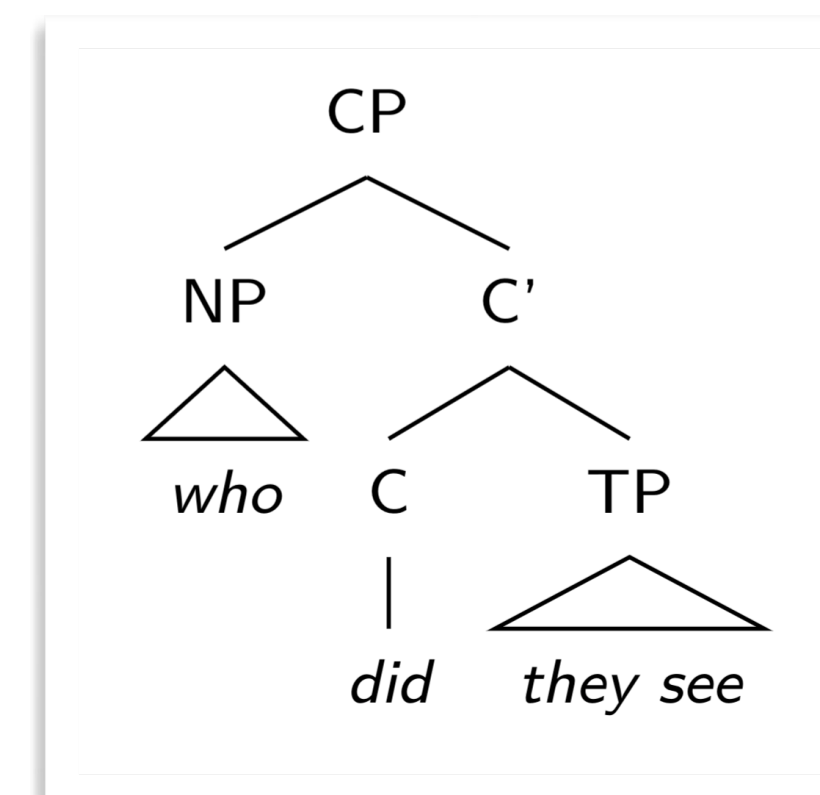
- Recorded speech and text associated with sociolinguistic variables:
 - Gender, age, geographic region...
- Interlinearized Glossed Text
 - Linguistic **analysis** and **annotation**
- Structural** annotations
- What about syntax trees?
- In **computational** linguistics?
 - In **NLP**, emphasis on **raw** data
 - Why?
 - NLP is a computer science discipline
 - Deep learning

mā?ā-nĩ *sàá* = \emptyset *nè* =*V*
 who-INDEP house =be there =Q
 ‘Whose house is that?’ [bxI]

Heath 2017. *A grammar of Jalkunan (Mande)*

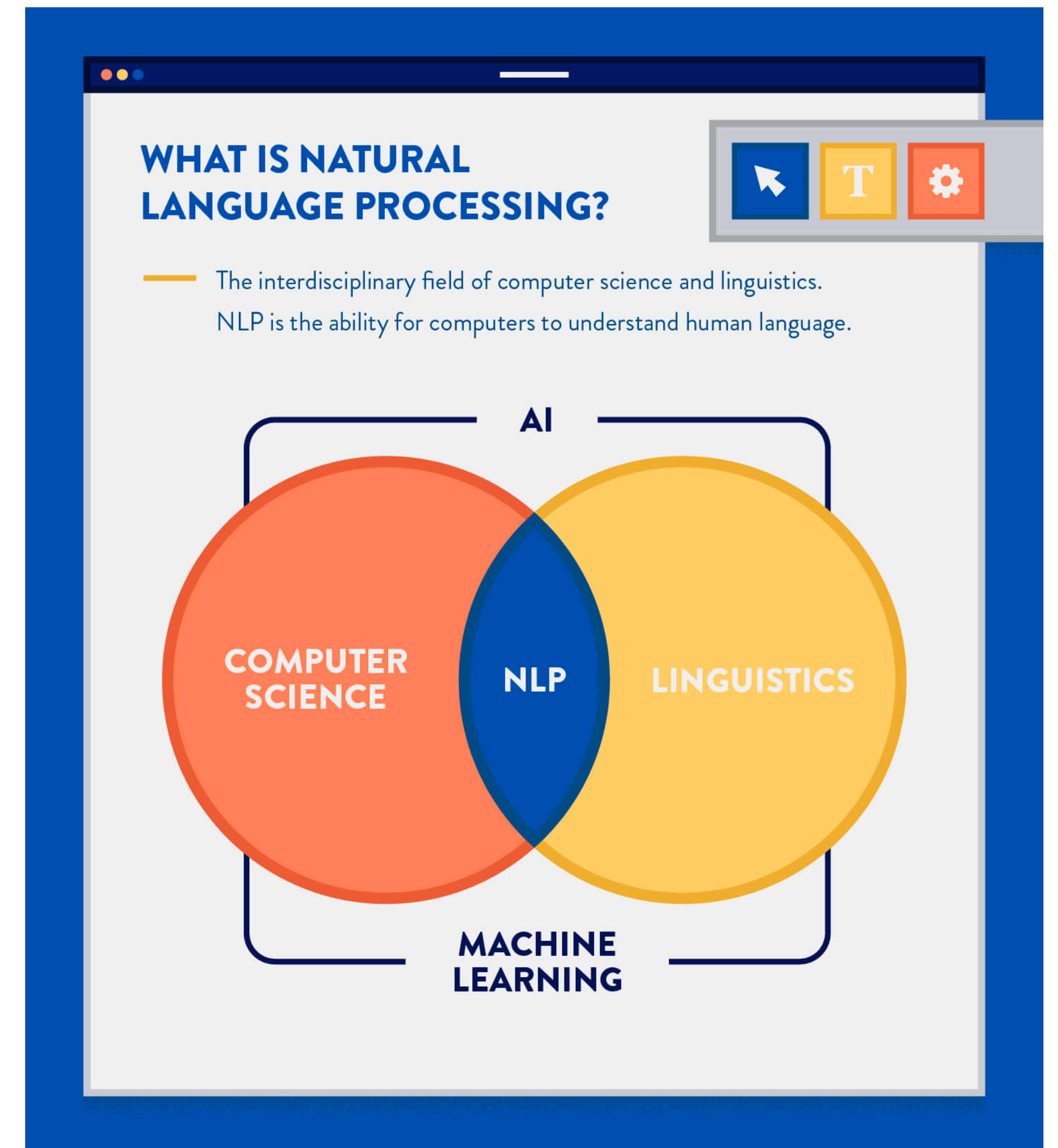


A dependency graph for the above sentence
 One of the **most** useful data formats in NLP!



Computational linguistics and natural language processing

- Is it the same thing?
 - For many people, yes! Interchangeable use...
 - For some:
 - **CL:** What can we learn about human language with computer aid?
 - E.g. How did pronunciation in Wales changed over the last 50 years?
 - Is this data science?
 - **NLP:** What can we learn about the world through the lens of language data? (Mostly text...)
 - E.g. How likely is the person to enjoy one movie based on **how they talk** about another movie?
 - Is this data science?



<https://clevertap.com/blog/natural-language-processing/>

NB: Olga **doesn't** think computers can "understand" human language!
(But that's a philosophical debate.)

Natural language processing and data science

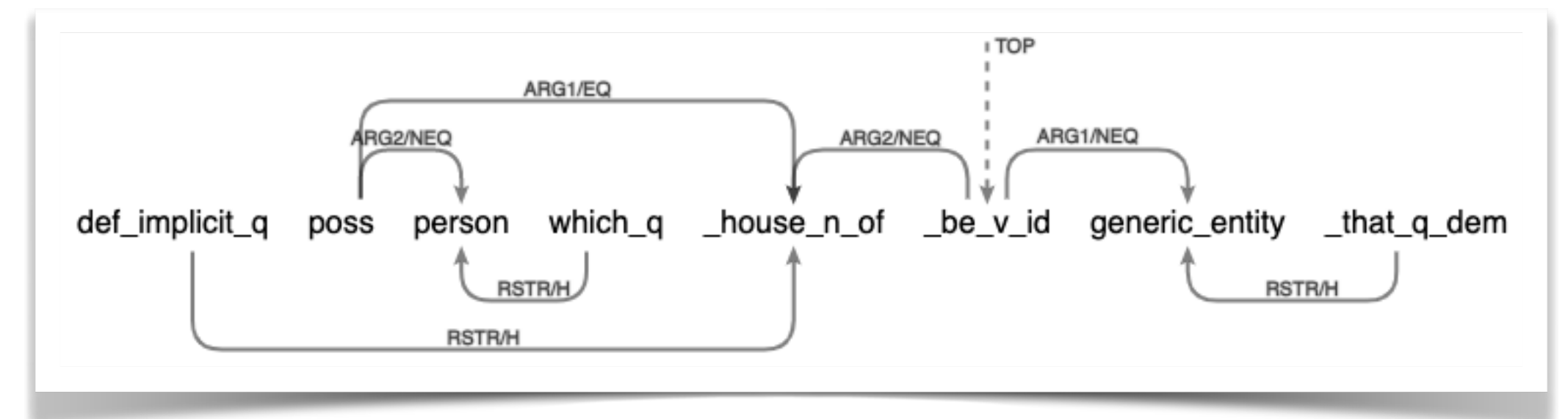
- Predicting people's behavior based on how they talk and/or write:
 - Advertisement, news, **movie** suggestions, etc.
 - IMDB reviews dataset
 - Each review is **labeled** as "positive" or "negative"
- "NLP is a branch of data science..."
 - maybe... **if** data science is a discipline!
 - Automatic speech recognition?
 - Machine translation?



<https://www.translatemedia.com/translation-blog/machine-translation-multilingual-sentiment-analysis-projects/>

Raw data in NLP

- Goal: **Automation**
 - Reduce/eliminate the need in annotation
 - Pros and Cons?
 - Cheaper products and services
 - Fewer jobs
- Despite the goal, NLP **relies** on annotation:
 - For training some machine learning algorithms
 - For **evaluation**



A dependency graph

Questions?