### DIGILAB WORKSHOP SERIES: TEXT ANALYSIS 101

### CREATING YOUR OWN SOCIAL MEDIA CORPUS



12.



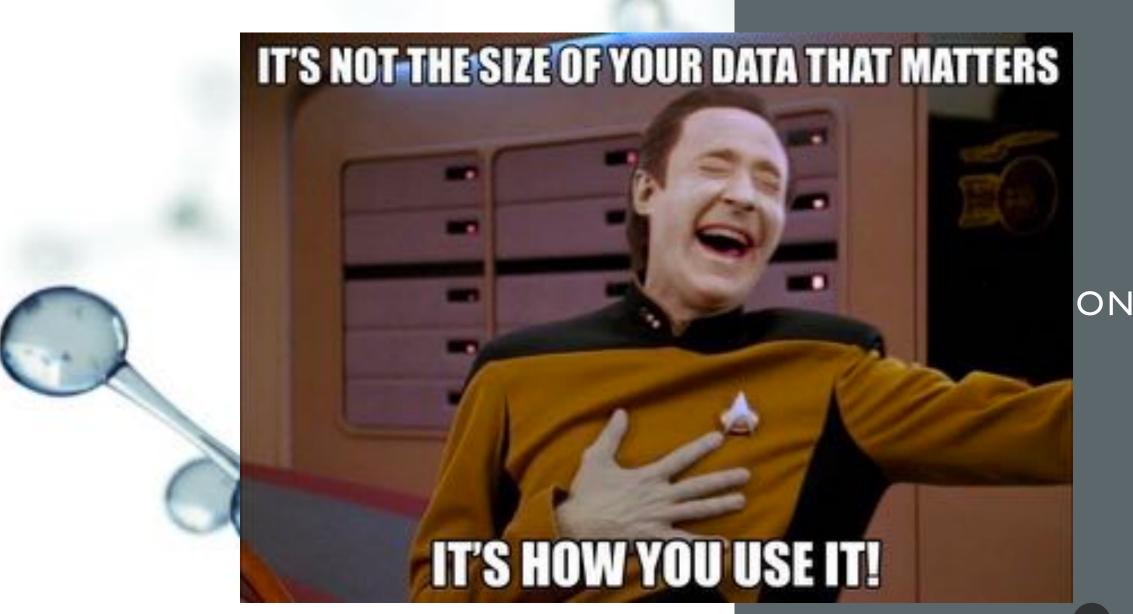






## SOCIAL MEDIA & TEXT ANALYSIS

- There are now many different types of social media that are useful in a variety of disciplines and areas for research.
- Social media web sites contain various types of services and thus create different formats of data, including text, image, video (Hu & Lui 2012).
  - Social networks (facebook, linkedin), media sharing (youtube, Instagram), discussion forums (reddit, quora), microblogging (Twitter, Facebook), review networks
- Applications include: business, research, event detection, linguistic and language change, network analysis, opinion mining, bioscience, insights into community behavior, and more!



## DATA SOURCES

- Linguistic Data Consortium
- The World Wide Web
- Kaggle.com: many pages of social media datasets, including tweets, and others: example: disaster tweets dataset, Instagram data, emojis, reddit, and many many others.
- <u>Stanford SNAP</u>: large network dataset collection, including data from amazon, social media, Wikipedia and others
- <u>Network Repository</u>: including social networks, biological, graph data and tools for analyzing and comparing available datasets

## METHODS: CORPUS CREATION





- <u>Tidytext</u>: helpful for data formatting and visualization; works well with other packages in the Tidyverse (Silge & Robinson 2016)
- <u>Quanteda</u>: incredibly useful package; includes preprocessing abilities, dtm function, as well as statistical analyses options like document classification and topic modeling
- <u>RedditExtractoR</u>: utilizes Reddit API to obtain posts, comments, and subreddit information
- <u>Rtweet</u>: useful package for getting Twitter data, with options for accessing followers, retweets, geolocation, and additional metadata.
- **Tokenizers**: useful options for text analytics, including tokenization and stemming.

## ADDITIONAL RESOURCES

iScience Maps: web-based option for getting Twitter data, with options for sorting and analyzing the data

Naoyun: software for connecting Twitter data with Gephi, with options for visualizing "live Twitter activity"

Netlytic: uses APIs to collect public data from Twitter, YouTube, and RSS feeds. Includes free and paid user options, with network and text analytics

Socioviz: get and analyze Twitter data in this web-based environment The Chorus Project: free web-based option for analyzing and obtaining Twitter data; based out of the UK

Webometric Analyst: free Windows-based program for gathering data, including Social media, from the Statistical Cybermetrics Research Digital Footprints: obtain and analyze Facebook data; web-based service available for researchers, based out of Aarhus University InfoExtractor: no longer maintained, but offers options for getting data from different URLs Snoopreport: free for researchers; focus on obtaining Instagram data



- <u>streamR</u>: Access to Twitter Streaming API via R
- <u>twittR:</u> also useful for getting twitter data in R
- <u>Rfacebook:</u> Rfacebook: Access to Facebook API via R
- **instaR:** access Instagram data via the Instagram API; an approved developer account is required



## ADDITIONAL TOOLS



- Python libraries:
  - Facebook SDK: Facebook data scraper
  - <u>Twitter scraper</u>: for use with Python 3.6+; can get tweets based on user or other search terms
  - <u>Reddit scraper</u>: interacts with Reddit API and PRAW library to obtain Reddit data
  - <u>Tweepy</u> in Python will interact with Twitter API
- <u>URS</u>: Universal Reddit Scraper; command line tool to obtain Reddit data
- <u>MOZDEH</u>: Windows based programming for gathering social media data
- Webscraping: <u>Chrome plugin</u>
  - <u>Beautiful Soup</u>: useful python library for webscraping; better for smaller amounts of data
  - <u>Scrapy</u>: python library; best for larger datasets
  - <u>Selenium</u>: flexible, also beginner friendly library
  - R: <u>xml2</u> and <u>rvest</u> work well in conjunction to harvest web data
  - <u>Rcurl</u> & <u>RSelenium</u>



## RECOMMENDED RESOURCES

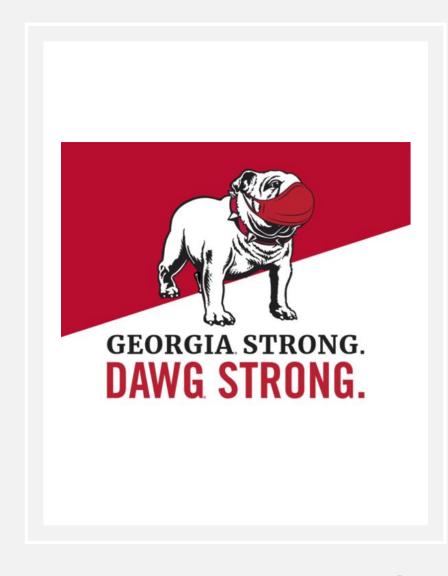
- Silge and Robinson's Text Mining with R
- Beckman et al.'s <u>Intro to Statistical</u> <u>Programming with R</u>
- Social Media Research using R
- Python 3 tutorials
- <u>Social Media Analytics</u>: helpful overview of options and types of analyses

12

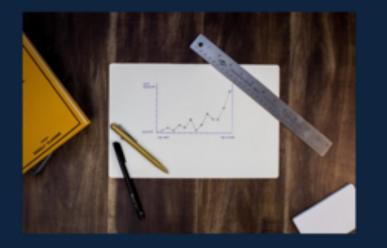
- <u>Text Analysis Glossary</u>
- Corpus Approaches to Social Media (Rüdiger & Dayter 2020)
- Linked-in Learning Tutorials
  - \_\_\_\_\_

## COURSES AT UGA

- This Fall 2021:
- Natural Language Processing: LING 4570/6570
- Style: ENGL/LING 4826/6826
- American English: ENGL/LING 4010/6010
- Note: These all count toward the Digital Humanities Undergraduate certificate!



# **DATA OFFICE HOURS**



#### CONSULTATIONS FOR DATA CLEANING, STRUCTURING, AND VISUALIZING

Whether just starting your work, or trying to make sense of your research, schedule an appointment for our Data Office Hours and bring your data (text, archival information, numerical data, etc.) for advice and guidance on your project. Expertise in corpus linguistics, Excel, and R, among other tools for data structuring and visualization.

#### TUESDAYS • 4:00-5:00 WEDNESDAYS • 2:00-3:00

To schedule an appointment visit DIGI.UGA.EDU/RESOURCES

WILLSON

## UP NEXT...

#### 22 April: Text Analysis Applications: Social Media

#### THANKS FOR LISTENING!

KATHERINE.KUIPER25@ UGA.EDU

PLEASE FILL OUT THIS <u>SURVEY</u>.



#### WORKS CITED

- Batrinca, Bogdan & Philip Treleaven. 2014. Social media analytics: a survey of techniques, tools, and platforms. Al & Society.
- Beckman, Matthew, Stéphane Guerrier, Justin Lee, Roberto Molinari, Samuel Orso & legor Rudnytskyi. 2020. An Introduction to Statistical Programming with R. <u>https://smac-group.github.io/ds/index.html</u>
- Benoit, Kenneth, Kohei Watanabe, Haiyan Wang, Paul Nulty, Adam Obeng, Stefan Müller, Akitaka Matsuo, William Lowe. (2018). "quanteda: An R package for the quantitative analysis of textual data." *Journal of Open Source Software*, 3(30), 774. doi: 10.21105/joss.00774, https://quanteda.io.
- Bird, Steven, Ewan Klein, and Edward Loper. 2019. Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit
- Brezina, Vaclav. 2018. Statistics in Corpus Linguistics.
- Brown, Simon. 2016. Tips for Computational Text Analysis. https://matrix.berkeley.edu/research/tips-computational-text-analysis
- Bussiere, Kirsten. 2018. Digital Humanities A Primer.
- Evert, Stefan. 2007. Corpora and collocations. http://www.stefan-evert.de/PUB/Evert2007HSK\_extended\_manuscript.pdf
- Feinerer, Ingo. 2020. Introduction to the tm Package: Text Mining in R. <u>https://cran.r-project.org/web/packages/tm/vignettes/tm.pdf</u>
- Freelon, Deen. http://socialmediadata.wikidot.com/
- Han, Na-Rae. Python 3 tutorials. http://www.pitt.edu/~naraehan/python3/.
- Jockers, Matthew. 2020. Introduction to the Syuzhet Package. <u>https://cran.r-project.org/web/packages/syuzhet/vignettes/syuzhet-vignette.html</u>
- Kearney, Matthew. 2018. R: Collecting and Analyzing Twitter Data: featuring {rtweet}. NiCAR 2018. https://mkearney.github.io/nicar\_tworkshop/#1
- Kearney, Matthew, Andrew Heiss, and Francois Briatte. 2020. Package 'rtweet'. https://cran.rproject.org/web/packages/rtweet/rtweet.pdf
- Kuiper, Katie Ireland. 2021. Text Analysis Glossary. DigiLab.
- Laudun, John. Text Analytics 101. <u>https://johnlaudun.org/20130221-text-analytics-101/</u>
- Lincoln, Mullen, 2018. https://cran.r-project.org/web/packages/tokenizers/vignettes/introduction-to-tokenizers.html

- Machlis, Sharon. 2020. How to search Twitter with rtweet and R. infoworld.com
- 2020.Modern Perl:Why Perl Rules for Text. <u>https://somedudesays.com/2020/02/modern-perl-why-perl-rules-for-text/</u>
- https://monkeylearn.com/text-analysis/
- Millot, Thomas. Photo. Unsplash
- Morikawa, Rei. 2019. 12 Best Social Media Datasets for Machine Learning. https://lionbridge.ai/datasets/12best-social-media-datasets/
- Nordquist, R. 2019. "Definition and Examples of Text in Language Studies. <u>https://www.thoughtco.com/text-language-studies-1692537</u>
- Norman, Jeremy. Thomas Mendenhall Issues One of the Earliest Attempts at Stylomtery. Historyofinformation.com <a href="https://www.historyofinformation.com/detail.php?id=4120">https://www.historyofinformation.com/detail.php?id=4120</a>
- O'Connor, Brendan, David Bamman, and Noah Smith. 2011. Computational Text Analysis for Scoial Science: Model Assumptions and Complexity.
- Parlante, Nick. 2002. Essential Perl. http://cslibrary.stanford.edu/108/EssentialPerl.html.
- Rivera, Ian. 2019. package RedditExtractoR. <u>https://cran.r-project.org/web/packages/RedditExtractoR/RedditExtractoR.pdf</u>
- Rüdiger, Sophia, and Daria Dayter. 2020. Corpus Approaches to Social Media. In Studies in Corpus Linguistics.
- Silge, Julia, and David Robinson. 2017. Text Mining with R: A Tidy Approach. <u>https://www.tidytextmining.com/</u>
- Verma, Abhishek. 2021. Inspirational Quotes from GoodReads website. https://www.kaggle.com/abhishekvermasg1/goodreads-quotes/metadata
- Wickham et al., (2019). Welcome to the tidyverse. Journal of Open Source Software, 4(43), 1686, https://doi.org/10.21105/joss.01686
- Wiedemann, Gregor & Niekler, Andreas. 2017. Hands-on: A five day text mining course for humanists and social scientists in R. Proceedings of the 1st Workshop on Teaching NLP for Digital Humanities (Teach4DH@GSCL 2017), Berlin.
- Witten, Ian. 2004. Text mining. https://www.cms.waikato.ac.nz/~ihw/papers/04-IHW-Textmining.pdf