CORNELL CENTER for
SOCIAL SCIENCES

*Accelerates, enhances, and amplifies*
*social science research at Cornell.*

# Natural Language Processing

Remy Stewart

CCSS-ResearchSupport@cornell.edu

# CCSS Research Support Code of Conduct

The Cornell Center for Social Sciences provides a welcoming environment for everyone embracing all backgrounds or identities. All instructors and attendees agree to abide by our community norms. We encourage the following behaviors in our workshops:

- Respect differing viewpoints and ideas
- Share your own perspectives and ask any questions
- Accept constructive criticism
- Use welcoming and inclusive language
- Show courtesy and respect for all instructors and attendees

If you believe that an instructor or attendee has violated the code of conduct, please report the violation to [CCSS-ResearchSupport@cornell.edu](mailto:CCSS-ResearchSupport@cornell.edu). We take all reported incidents seriously.

# Land Acknowledgement

Cornell University is located on the traditional homelands of the Gayogo̱hó:nǫ' (the Cayuga Nation). The Gayogo̱hó:nǫ' are members of the Haudenosaunee Confederacy, an alliance of six sovereign Nations with a historic and contemporary presence on this land. The Confederacy precedes the establishment of Cornell University, New York state, and the United States of America. We acknowledge the painful history of Gayogo̱hó:nǫ' dispossession, and honor the ongoing connection of Gayogo̱hó:nǫ' people, past and present, to these lands and waters.

Here are additional links for more on the history of Cornell's violent, colonial formation, the movement to return native lands, and about the AIISP program at Cornell.

Consider donating to the Gayogo̱hó:nǫ' sovereignty initiative here.

# Text as Data



- Text has been a core data source within social sciences

- Often previously limited to qualitative methods

- Digital Age has rapidly shifted these circumstances
  - Wide range of innovative computational methods
  - Powerful & cheaper compute resources
  - Open-source software that is accessible to learn

# Text Data in the Digital Age

- Text data is rapidly created as written behavior is recorded online

- Linked to the growth of social media & online communities

- Digital conversion of printed books, historical records, government documents, etc.

# Natural Language Processing (NLP)



- Interdisciplinary field using computational tools to interpret text & language

- While not exclusively ML based, ML is highly influential within NLP

- Powers the AI behind language translation, audio-to-text transcription, speech recognition (Siri, Alexa)

- NLP as a natural computational approach to incorporate within social sciences

# Word embeddings from Kozlowski, Taddy, and Evans (2019)

# Structural Topic Models from Wilkerson & Casas (2017)

# Core Terminology

**1**. A **corpus** is the total collection of text being analyzed.

→ **All 400k of our tweets.**

**2**. A **vocabulary** is all of the unique words in a corpus.

→ **All words used at least once across the tweets.**

**3**. A **document** is an individual record within the corpus.

→ **A user tweets "This brings me so much joy."**

**4**. **Tokens** are the smaller text units within a document. Most commonly refers to words, but can also be sentences, paragraphs, or text characters.

→ **"This", "brings", "me", "so", …**

**5**. **N-grams** considers multiple word tokens at once with "n" as a specified number.

→ **Bi-gram of "This_brings"**
**Tri-Gram of "This_brings_me"**

# Coding Demo Outline

## Supervised Text Classifier

| Text-Based ML Considerations | Baseline Logistic Classifier | Adjusted Logistic Classifier | Model Interpretability |
|---|---|---|---|
| - Vocabulary<br>-Word Frequency Distributions | - Multilabel Task<br>- Reviewing Performance | - Regularization<br>- Class Balancing | - Top Keywords<br>- Misclassified Text |

## Unsupervised Topic Model

| Document-Term Matrix | Latent Dirichlet Allocation (LDA) | Interpreting Topics | Visualizing Topics |
|---|---|---|---|
| - Vectorization & tokenization<br>- Custom stop words | - k parameter | - Words across topics<br>- Topics across documents | - pyLDAvis |

# Let's head to the Github for our final coding demo!

Alternatively, you can access this session's Colab file directly via the following link:

https://colab.research.google.com/drive/1VwetZHUP75J I6rKcmKI26skQaJ2J-zvR?usp=sharing

# Additional NLP Methods

## Word Embeddings

Mapping word similarity within vector spaces



## Text Networks

Node/edge relationships between text

# Additional NLP Methods

**Text & Causality**

Text as treatment vs. text as outcome



A gap in views of the availability of jobs and 'good jobs'

% saying _____ in their community

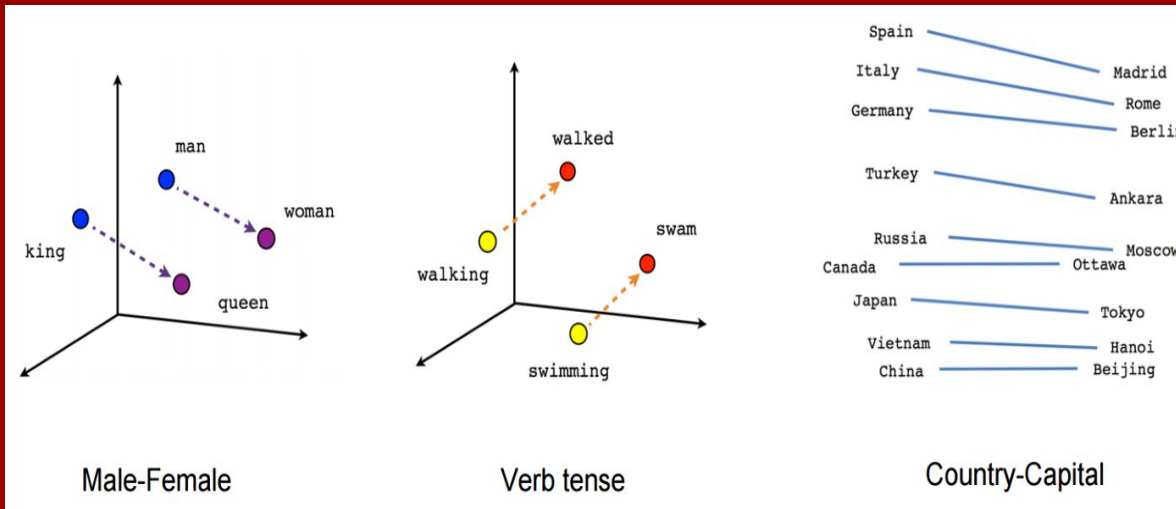| | Jobs are difficult to find | Plenty of jobs available | | GOOD jobs are difficult to find | Plenty of GOOD jobs available |
|---|---|---|---|---|---|
| Total | 33 | 60 | Total | 45 | 48 |
| Rep/Lean Rep | 23 | 71 | Rep/Lean Rep | 36 | 58 |
| Dem/Lean Dem | 39 | 53 | Dem/Lean Dem | 55 | 39 |

Note: Respondents were randomly asked about either "jobs" or "good jobs."
Source: Survey of U.S. adults conducted Jan. 9-14, 2019.

Pew Research Center

**Neural Networks**

Classification, question answering, summation, and beyond



Transformers

PyTorch Lightning

OpenAI

# Ethics Within ML

- Data is a fundamentally social product embedded within structural inequities

- ML has great ability to reinforce bias & perpetuate harm

- Many ethical concerns within ML only engaged with after problems are retroactively discovered
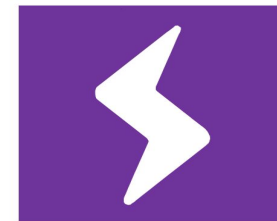
# Ethics within ML

## Ethical Questions within ML Research

- What are the potential risks associated with the model?
- What is the broadest range of potential consequences/outcomes of this research?
- How could vulnerable groups be impacted?
- What are the relationships between researcher, participants, & third parties?
- What are expectations around user autonomy, privacy, informed consent?

## Core Ethical Considerations

- Model interpretability & transparency
- Data privacy
- Participant agency & community engagement

# Key Bootcamp Takeaways

- Learned core ML concepts

- Processed and modeled structured and unstructured data

- Applied both supervised and unsupervised methods

- Thought through key decision points within ML analyses

- Compared models for different use cases

# Conclusion

Social scientists are in a unique position to address essential questions with ML given our disciplinary perspectives.

This series provided a broad overview to extend through further learning.

The world of ML can be intimidating- we truly believe that you can master these skills.
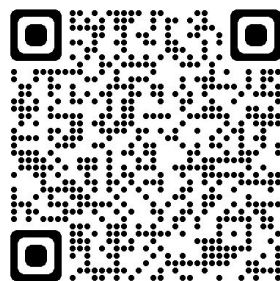
We hope this bootcamp series can serve as a source of inspiration for your own research!

CORNELL CENTER for
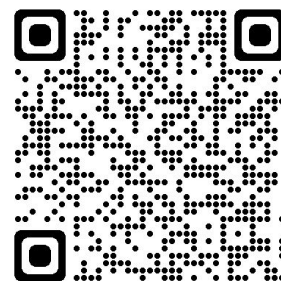SOCIAL SCIENCES

# Research Services

## Computing

- Shared & Dedicated Cloud Computing Resources
- Secure Data Enclave
- Software
  - Stata, SPSS, SAS
  - Python
  - R
  - Matlab
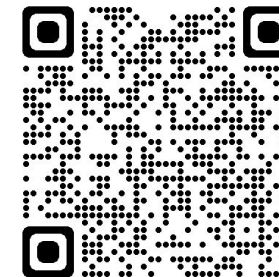  - ArcGIS
  - Atlas.ti
  - MaxQDA

## Data Services

- Data Access
- Secure Data Services
- Data Sharing & Archiving
- Results Reproduction

## Consulting

- Popular Topics:
  - Qualitative Methods
  - Data Cleaning
  - Visualizations
  - Survey Analysis
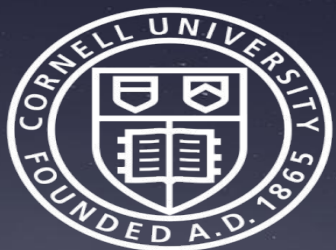  - Finding Data
  - Web Scraping
  - Text Analysis

CORNELL CENTER for
SOCIAL SCIENCES

# Please Fill Out Our Evaluation Forum!

We'd love to hear your feedback to improve this Bootcamp series.

# Thank you for joining us!

**Angel Hwang**

CCSS Senior Data Science Fellow

**Remy Stewart**

CCSS Data Science Fellow

**CCSS-ResearchSupport@cornell.edu**