



JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

472/672 NLP for Computational Social Science: Introduction

1/22/2025

Today

- What is computational social science?
 - Definitions and examples
- What are the expectations and logistics for this course?

What is Computational Social Science?

“The study of social phenomena using digitized information and computational and statistical methods”
[Wallach 2018]

Social Science

- Examine to what extent recommendations affect shopping patterns vs. other factors
- Analyze the impact of gender and race on the U.S. hiring system
- When and why do senators deviate from party ideologies?

Explanation

NLP

- Recommend related products to Amazon shoppers
- Predict which candidates will be hired based on their resumes
- How many senators will vote for a proposed bill?

Prediction [Task Automation]

Example 1: Media Manipulation





What are more specific questions we might investigate?

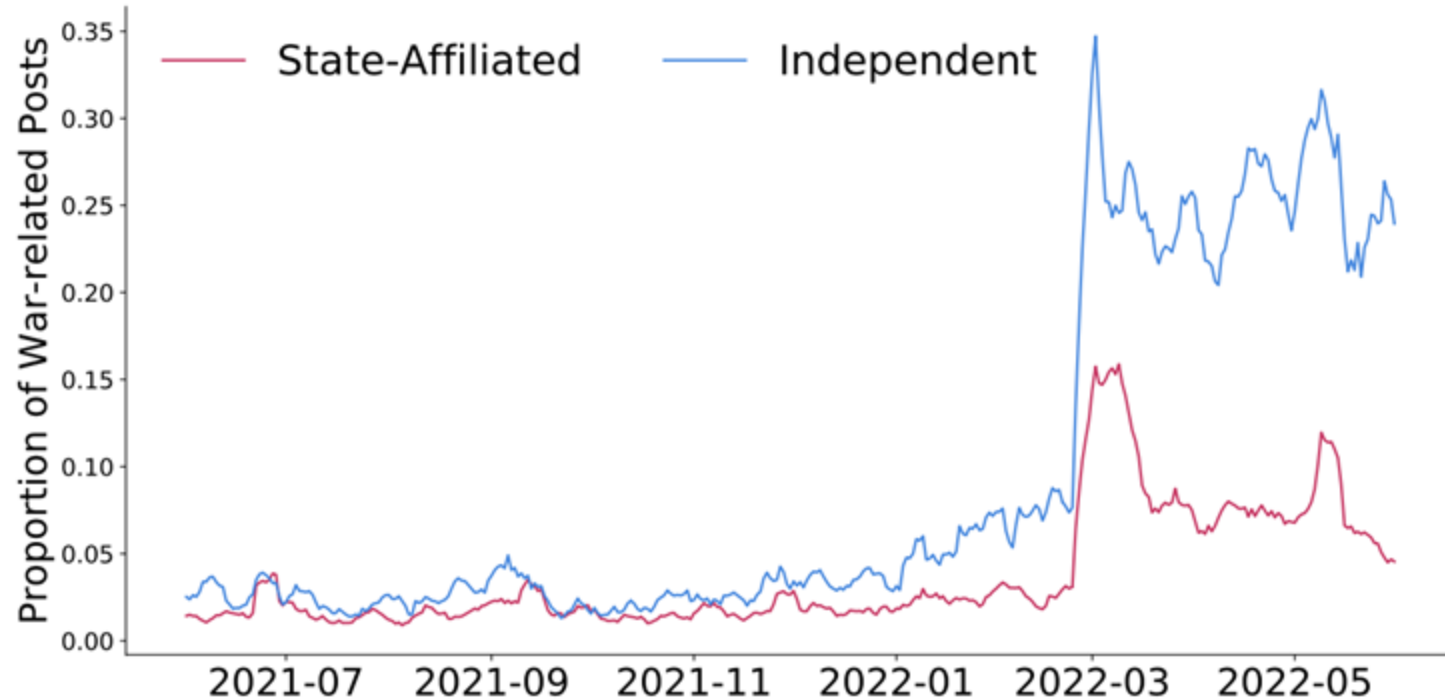
- What strategies does the Russian government use to influence public opinion?
 - Are the strategies different for internal (within Russia) vs. external audiences?
 - How do these strategies evolve over time?
- Are these strategies effective?
- What are successful ways to mitigate them?

Dataset Collection

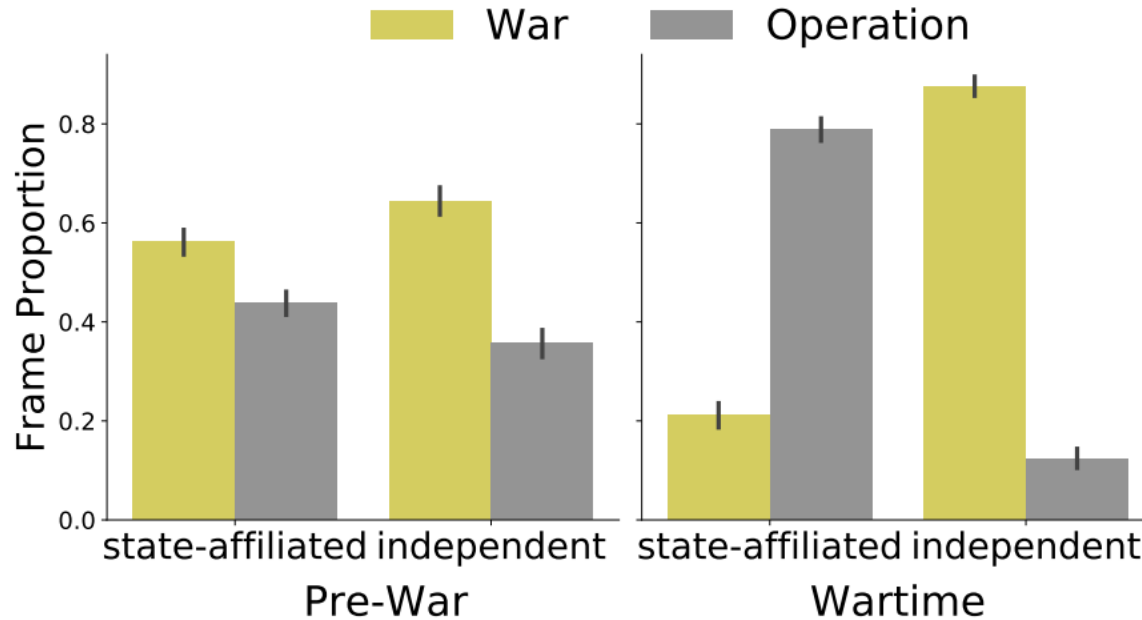
- Jan 01 2021 ~ May 31 2022
- Three dimensions
 - **Time:** pre-war, during-war
 - **Platform:** Twitter, VKontakte (VK)
 - **Media ownership:** state-affiliated, independent

23 State-affiliated outlets		20 Independent outlets	
RT_com	rbc	tvrain	snob_project
life	ria	Forbes	golosameriki
tassagency	gazeta	novgaz	svobodaradio
tv5	vesti	meduzaproject	BBC
rgru	Ukraine RU	rtvi	The insiders

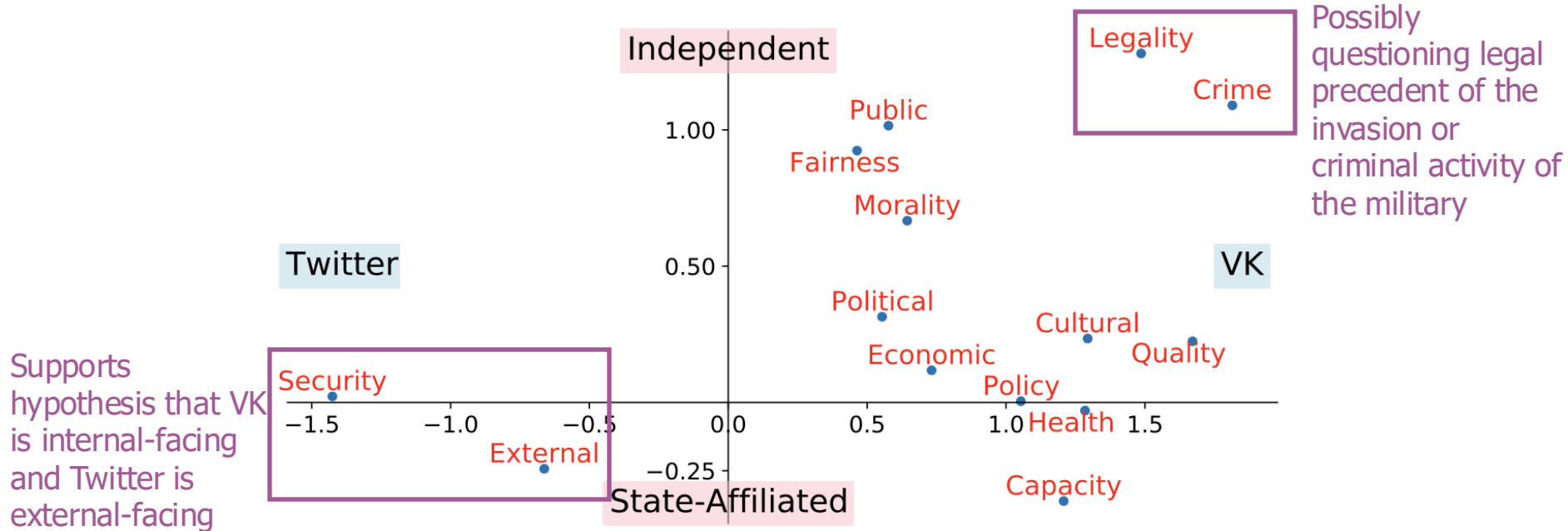
Starting Point: Count how often outlets mention the war (just count words)



Starting Point: Count how often outlets use “war” vs. “operation”



Deeper analysis: issue-generic frames



Example Social Science Domains

- Political science
 - What strategies do authoritative governments use to control public opinion?
- Linguistics
 - How dialects of English differ by geographic region?

Example 2: Changes in Word Meaning

Under what conditions do words change meaning?

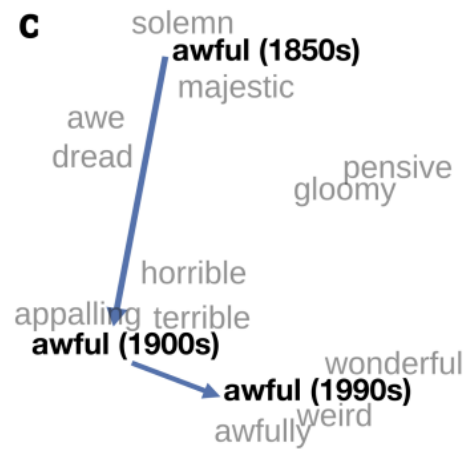
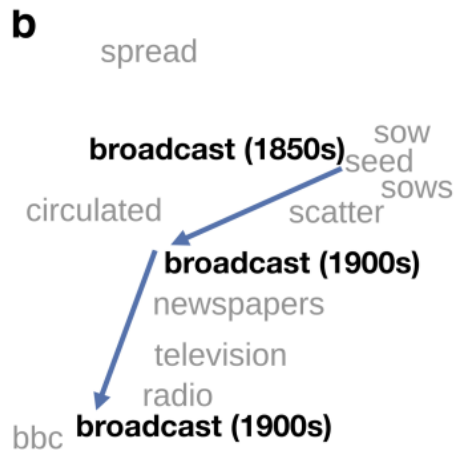
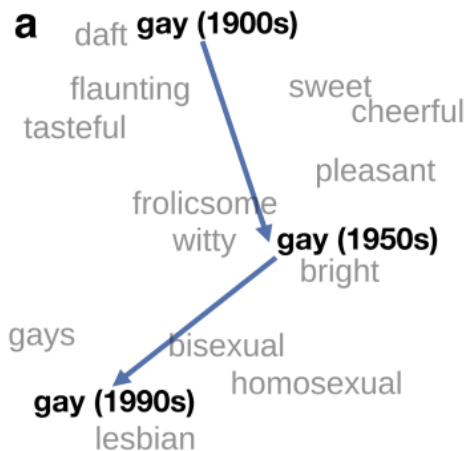
- Easy to observe that words change meaning over time
- “Awful”
 - 1850s
 - “Full of awe”, majestic, inspiring wonder (or fear)
 - 1900s
 - Terrible, dreadful
- Lots of examples, case studies conducted by hand, what are more systematic trends?

Examine changes in word meanings over decades

- Use NLP (word embeddings) to quantify how words change meaning over time

Name	Language	Description	Tokens	Years	POS Source
ENGALL	English	Google books (all genres)	8.5×10^{11}	1800-1999	(Davies, 2010)
ENGFIC	English	Fiction from Google books	7.5×10^{10}	1800-1999	(Davies, 2010)
COHA	English	Genre-balanced sample	4.1×10^8	1810-2009	(Davies, 2010)
FREALL	French	Google books (all genres)	1.9×10^{11}	1800-1999	(Sagot et al., 2006)
GERALL	German	Google books (all genres)	4.3×10^{10}	1800-1999	(Schneider and Volk, 1998)
CHIALl	Chinese	Google books (all genres)	6.0×10^{10}	1950-1999	(Xue et al., 2005)

Validate methodology with known examples



Proposed Laws of Semantic Change

- Remember, starting question: Under what conditions do words change meaning?
- **The law of conformity:** “Rates of semantic change scale with a negative power of word frequency.”
 - [Frequent words change more slowly]
- **The law of innovation:** After controlling for frequency, polysemous words have significantly higher rates of semantic change.
 - [Polysemy: words with multiple meanings]

Example Social Science Domains

- Political science
 - What strategies do authoritative governments use to control public opinion?
- Linguistics
 - How dialects of English differ by geographic region?
- Psychology
 - What types of language do readers of online mental health support forums perceive as empathetic?
- Sociology:
 - How do social media users engage in collective action?

Example 3: Emotions in Social Movements

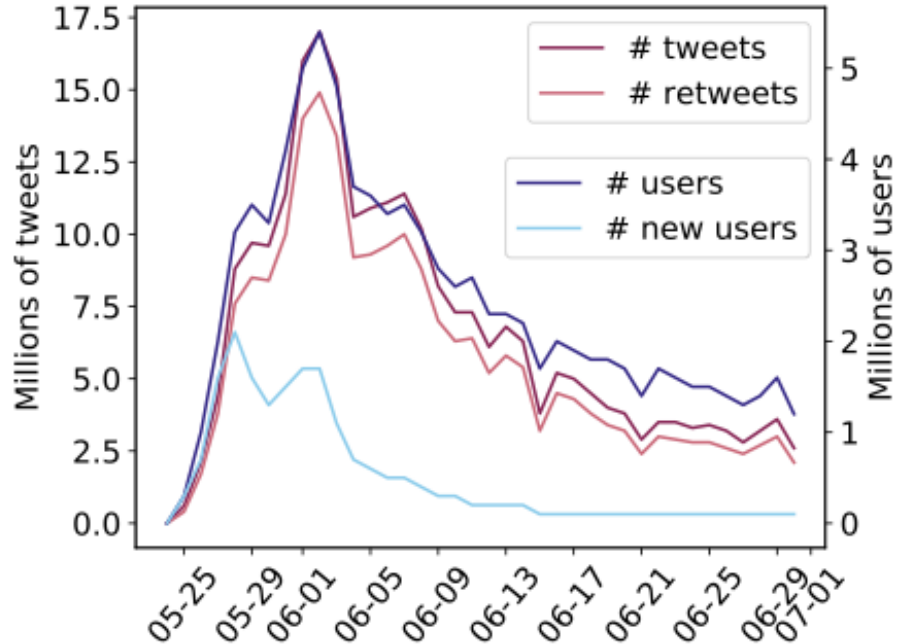
Sociology and Social Psychology models of emotions

- What motivates social movements?
- “Moral shocks” can cause people to join social movements, but sense of camaraderie, optimism, and hope for change are necessary for sustained involvement
 - Qualitative interview studies: we can ask people why they joined protests
 - Potentially limited in scope
 - Is this still true in the era of online social movements?

Analysis Data: 34M tweets about the #BlackLivesMatter Movement

The term #BlackLivesMatter originated in posts made by activists Alicia Garza and Patrisse Cullors in 2013

#BlackLivesMatter
#JusticeForGeorgeFloyd
#ICantBreathe



Supervised machine learning pipeline

Annotate small
amount of data



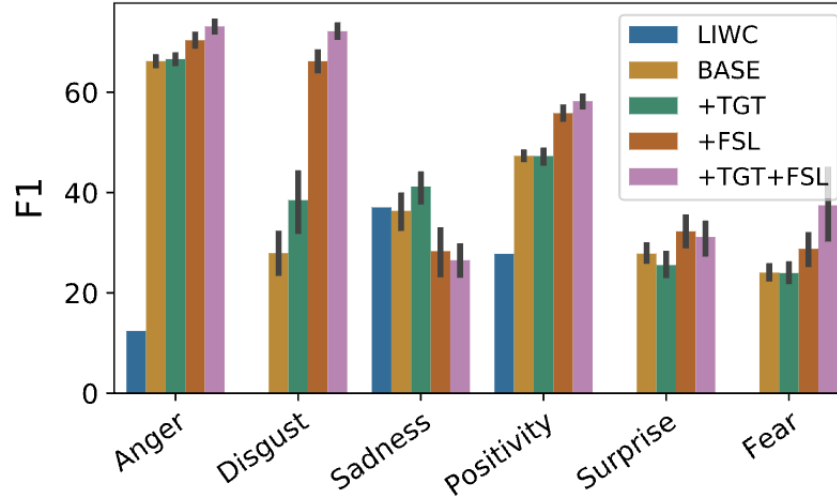
Train ML
model



Evaluate ML
model



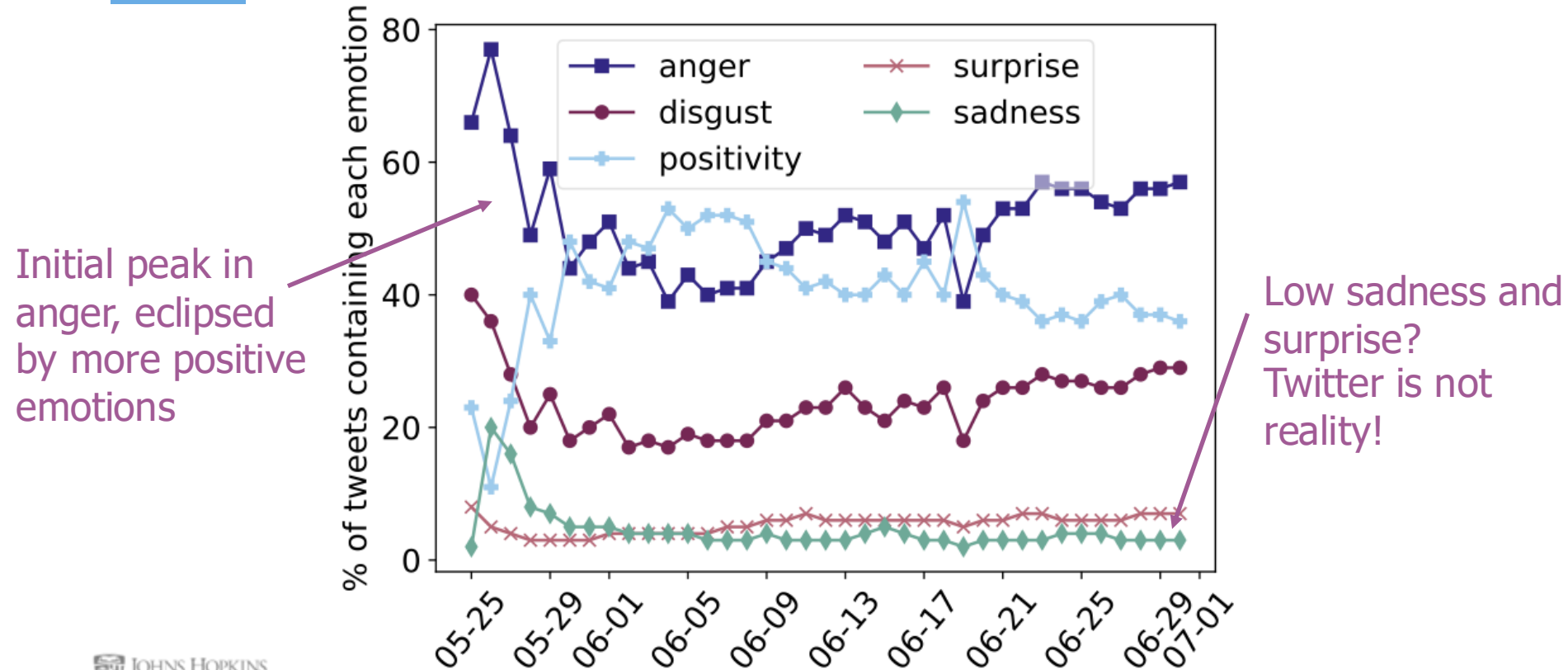
Analyze Full
data set



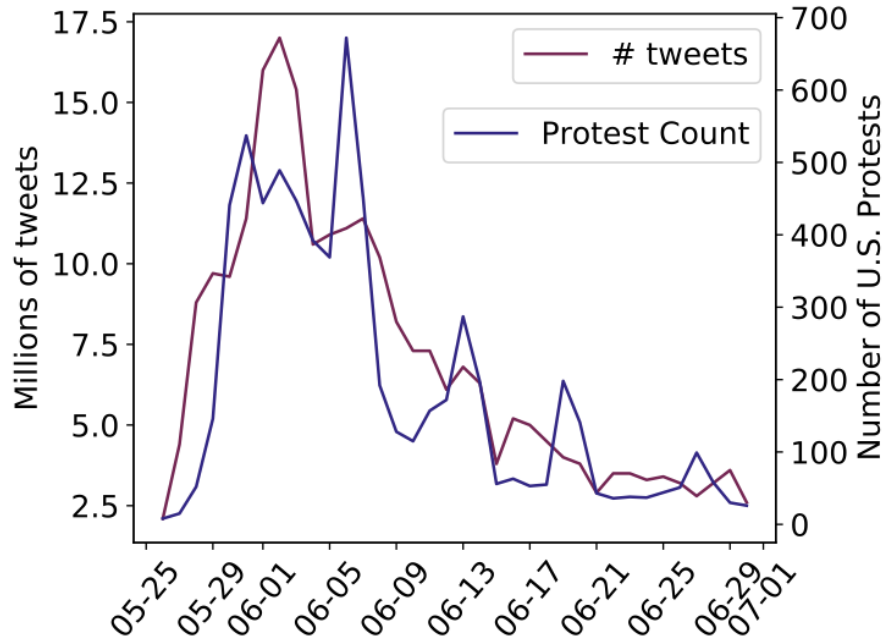
Limitations and Ethical Considerations

- Sample of tweets may not be representative
- Measuring emotions perceived in tweets
 - Cannot draw conclusions about what emotions people actually experienced
- Privacy and consent
 - Not showing any specific examples or usernames from the data

Emotions over time in tweets with pro-BLM hashtags



Positivity is correlated with in-person protests



	Correlation with protest across states	Correlation with protests across cities
Anger	-0.43*	-0.16*
Disgust	-0.24	-0.21*
Positivity	0.48*	0.12*
Sadness	-0.38*	0.06
Surprise	-0.25	0.09

Why NLP for Computational Social Science?

- Many social issues manifest in text and are intricately tied with language
 - Information manipulation
 - Social movements
 - Toxicity, stereotypes, and prejudice
- Cross-cutting methodology
 - Natural language processing
 - Methods for *analysis* and *explanation*, not only *prediction*
 - Combining data types and metadata
- *Deductive* and *Inductive* approaches

Deductive vs. Inductive Reasoning Approach

- Deductive:
 - Establish theory
 - Infer predictions
 - Gather data to test predictions
- Inductive
 - Only after examining and understanding data can we narrow to research questions

This Course

Learning Objectives

- Quantitative analysis of social phenomena
- Methods for text analysis
 - [Course is organized by methodologies]
- Example applications to social science fields, such as political science, sociolinguistics, sociology, and economics

Course Topics

- Unsupervised (off-the-shelf) approaches
 - Word statistics, topic modeling, word embeddings, lexicons
- Supervised approaches
 - Data annotating, classification models, interpreting model outputs
- Incorporating meta data
 - Time series analysis, network analysis, causal inference
- Current state-of-the art methods
 - Language models
- Draws from *NLP (ML)*, *statistics*, and *social science*

Logistics

- 4 HW assignments (40%)
 - (one for each meta-topic)
 - HW 1, 3, and 4 are primarily coding assignments (Python). HW 2 focuses on data annotation design
- 5 late days
 - Intended to be used for any circumstances that you may have (conflicting deadlines, illness, etc.)
 - Generally other extensions will not be granted unless there are extreme extenuating circumstances
 - If you use a late day on a group assignment it will count as a late day for everyone in the group
 - Cannot be used for final project

Logistics

- Project (25%)
 - Part 1: Proposal
 - Part 2: Results
- 1 in-class midterm (25%)
- Participation / course surveys (10%)
 - In-class iClicker quizzes (dropping lowest 1/3), course goals, midterm survey

Academic Integrity

- This class abides by JHU academic policies
- You are encouraged to discuss assignments with your classmates; however, what you hand in should be your own work
- Generally okay to use open-source software with proper acknowledgements
- Copying/reusing code is not allowed; strict action will be taken if similarities are found
- Copying content from other published work (without citing it) is also not allowed, and is considered plagiarism

Generative AI Policy

- **What you hand in should be your own work**
- Allowed uses of AI:
 - Use cases that aid in your understanding of course content (e.g. querying about content from lecture)
 - Translation or grammar correction
- Disallowed uses of AI:
 - Fully generating code or written text for assignments
- If you use AI for a HW assignment:
 - Your written report must contain a brief paragraph about how you used AI and your starting prompts

Course Staff



Miriam Wanner
(Teaching Assistant)

Kuleen Sasse
(Course Assistant)

Rohan Allen
(Course Assistant)

This Course

- Course website: <http://nlp-css-601-672.cs.jhu.edu/sp2025/>
- Canvas/Gradescope
 - Assignment submissions
 - Lecture recordings (Panopto)
- Piazza: <https://piazza.com/jhu/spring2025/en601472672>

Four principles of quantitative text analysis

1. All quantitative models of language are wrong—but some are useful
2. Quantitative methods for text amplify resources and augment humans
3. There is no globally best method for automated text analysis
4. Validate, Validate, Validate.

Why Computational Social science?

“Despite all the hype, machine learning is not a be-all and end-all solution. We still need social scientists if we are going to use machine learning to study social phenomena in a responsible and ethical manner.” [Wallach 2018]

Tells about your goals for this class

- <https://forms.office.com/r/BYyDj66xjM>

