

# Text Data in Business and Economics

Basel University – Autumn 2023

## 10. Some Extras about Recent NLP

# Outline

GPT

GPT for Research

Bias in NLP Systems

GPT = Generative Pre-Trained Transformer

# GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
  - ▶ trained on the Books corpus.
  - ▶ train on a language modeling task, as well as a multi-task that adds a supervised learning task.

# GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
  - ▶ trained on the Books corpus.
  - ▶ train on a language modeling task, as well as a multi-task that adds a supervised learning task.
- ▶ GPT-2 (2019):
  - ▶ all articles linked from Reddit with at least 3 upvotes (8 million documents, 40 GB of text)
  - ▶ dispense with supervised learning task, make some other architectural adjustments
  - ▶ make model much bigger

# GPT = Generative Pre-Trained Transformer

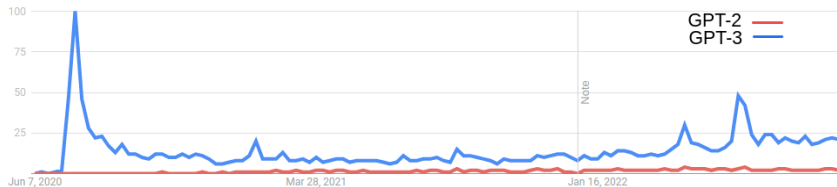
- ▶ GPT-1: the first autoregressive transformer model (2018)
  - ▶ trained on the Books corpus.
  - ▶ train on a language modeling task, as well as a multi-task that adds a supervised learning task.
- ▶ GPT-2 (2019):
  - ▶ all articles linked from Reddit with at least 3 upvotes (8 million documents, 40 GB of text)
  - ▶ dispense with supervised learning task, make some other architectural adjustments
  - ▶ make model much bigger
- ▶ GPT-3 (2020):
  - ▶ use an even bigger corpus (Common Crawl, WebText2, Books1, Books2 and Wikipedia)
  - ▶ make model much, much bigger

# OPENAI'S NEW MULTITALENTED AI WRITES, TRANSLATES, AND SLANDERS

*A step forward in AI text-generation that also spells trouble*

By James Vincent | Feb 14, 2019, 12:00pm EST

Howard, co-founder of Fast.AI agrees. "I've been trying to warn people about this for a while," he says. "We have the technology to totally fill Twitter, email, and the web up with reasonable-sounding, context-appropriate prose, which would drown out all other speech and be impossible to filter."



BUSINESS

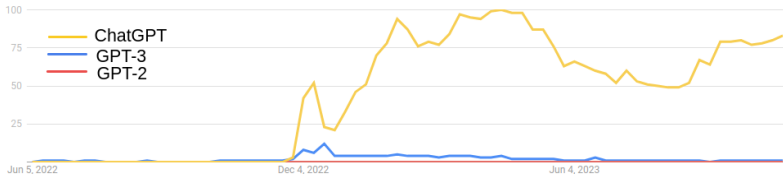
# Is ChatGPT the Start of the AI Revolution?

Analysis by The Editors | Bloomberg

December 9, 2022 at 1:39 p.m. EST

## *The New Chatbots Could Change the World. Can You Trust Them?*

Siri, Google Search, online marketing and your child's homework will never be the same. Then there's the misinformation problem.





# GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
- ▶ GPT-2 (2019)
- ▶ GPT-3 (2020)

# GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
- ▶ GPT-2 (2019)
- ▶ GPT-3 (2020)
- ▶ GPT-3.5 (2022):
  - ▶ Subclass of GPT-3 trained on data up to June 2021
  - ▶ Incorporates the base model on which ChatGPT is fine-tuned and after optimized for chat

# GPT = Generative Pre-Trained Transformer

- ▶ GPT-1: the first autoregressive transformer model (2018)
- ▶ GPT-2 (2019)
- ▶ GPT-3 (2020)
- ▶ GPT-3.5 (2022):
  - ▶ Subclass of GPT-3 trained on data up to June 2021
  - ▶ Incorporates the base model on which ChatGPT is fine-tuned and after optimized for chat
- ▶ GPT-4 (2023):
  - ▶ Multimodal model: can take also images as inputs
  - ▶ Trained in two stages:
    1. token prediction (like other GPT models)
    2. reinforcement learning with human feedback
  - ▶ Much, much, much bigger model

Stammbach and Ash (2020),  
“e-FEVER: Explanations and Summaries for  
Automated Fact Checking”

Claim	Extractive Summary	Abstractive Summary
Marnie was directed by someone who was "The Master of Nothing".	+ Alfred Hitchcock + Sir Alfred Joseph Hitchcock ( 13 August 1899 - 29 April 1980 ) was an English film director and producer , at times referred to as " The Master of Suspense " . + Marnie (film) + Marnie is a 1964 American psychological thriller film directed by Alfred Hitchcock ."	Marnie was directed by Alfred Hitchcock, who was "The Master of Suspense".
Colin Kaepernick is a poker player.	+ Colin Kaepernick + Colin Rand Kaepernick ( born November 3 , 1987 ) is an American football quarterback who is currently a free agent . + American football + American football , referred to as football in the United States and Canada , and also known as " gridiron football " or simply " gridiron " , is a sport played by two teams of eleven players on a rectangular field with goalposts at each end .	Colin Kaepernick is a football player, not a poker player.

Stammbach, Antoniak, and Ash (2022),  
"Heroes, Villains, and Victims, and GPT-3"

Who is the villain in the following text?

*Text: In questioning the use of hydraulic fracturing in New York to help produce natural gas, you do not note that the technology has been employed and continuously improved for more than 50 years and that studies by the Environmental Protection Agency and the Ground Water Protection Council have not identified a single instance of groundwater contamination. Wells where fracturing is used are specially constructed to protect drinking water sources. Regulatory oversight is extensive. The fluids mostly water that are forced into a well to create pressure to fracture rock are pushed back out by the oil and gas flowing upward for safe processing. Protecting our water supplies is important, as are reductions in greenhouse gas emissions through use of clean-burning natural gas. Banning hydraulic fracturing would be unwarranted and shortsighted, preventing production of large amounts of natural gas that could directly benefit New York consumers for decades and create thousands of good jobs.*

**Villain: The villain in this text is the person who is questioning the use of hydraulic fracturing in New York.**

Movie	Hero	Victim	Villain
101 Dalmatians	Roger Dearly	The Dalmatian Puppies	Cruella de Vil
Aladdin	Aladdin	Aladdin	Jafar
Cinderella	Cinderella	Cinderella	Lady Tremaine
Alice In Wonderland	Alice	Alice	The Queen of Hearts
The Jungle Book	Mowgli	Mowgli	Shere Khan, a man-eating Bengal tiger
Sleeping Beauty	Prince Phillip	Aurora	Maleficent
The Lion King	Simba	Mufasa	Scar
Peter Pan	Peter Pan	Wendy, John, Michael, and the Lost Boys	Captain Hook
Mary Poppins	Mary Poppins	Mr. Banks	Mr. Dawes
The Little Mermaid	Ariel	Ariel	Ursula
Snow White	Snow White	Snow White	The Queen

Table 2: Results for Wikipedia plots of widely known Disney Movies

# Outline

GPT

GPT for Research

Bias in NLP Systems

# Generative AI for Economic Research

- ▶ Korinek (2023)<sup>1</sup>:
  - ▶ Provides an overview of how researchers can exploit these technologies

---

<sup>1</sup>The slides for this section summarize the content of the paper by Korinek (2023)

# Generative AI for Economic Research

- ▶ Korinek (2023)<sup>1</sup>:
  - ▶ Provides an overview of how researchers can exploit these technologies
  - ▶ Discusses the potential productivity-enhancing effects of cognitive automation

---

<sup>1</sup>The slides for this section summarize the content of the paper by Korinek (2023)

# Generative AI for Economic Research

- ▶ Korinek (2023)<sup>1</sup>:
  - ▶ Provides an overview of how researchers can exploit these technologies
  - ▶ Discusses the potential productivity-enhancing effects of cognitive automation
  - ▶ Discusses long-run implications of AI for economic research

---

<sup>1</sup>The slides for this section summarize the content of the paper by Korinek (2023)



# Generative AI for Economic Research

- ▶ Korinek (2023)<sup>1</sup>:
  - ▶ Provides an overview of how researchers can exploit these technologies
  - ▶ Discusses the potential productivity-enhancing effects of cognitive automation
  - ▶ Discusses long-run implications of AI for economic research
- ▶ Main takeaway on LLMs for research:

“Treat LLMs as if it was an internet who is: 1) smart and incredibly motivated and eager to help 2) but completely lacks the context of what you are doing”

---

<sup>1</sup>The slides for this section summarize the content of the paper by Korinek (2023)

# Generative AI for Economic Research

- ▶ Korinek (2023)<sup>1</sup>:
  - ▶ Provides an overview of how researchers can exploit these technologies
  - ▶ Discusses the potential productivity-enhancing effects of cognitive automation
  - ▶ Discusses long-run implications of AI for economic research
- ▶ Main takeaway on LLMs for research:

“Treat LLMs as if it was an internet who is: 1) smart and incredibly motivated and eager to help 2) but completely lacks the context of what you are doing”
- ▶ Good practices:
  - ▶ Provide context
  - ▶ Iterate
  - ▶ Be patient

---

<sup>1</sup>The slides for this section summarize the content of the paper by Korinek (2023)

Category	Task	Usefulness
Ideation and Feedback	Brainstorming	●
	Feedback	◐
	Providing counterarguments	◐
Writing	Synthesizing text	●
	Editing text	●
	Evaluating text	●
	Generating catchy titles & headlines	●
	Generating tweets to promote a paper	●
Background Research	Summarizing Text	●
	Literature Research	○
	Formatting References	●
	Translating Text	●
	Explaining Concepts	◐
Coding	Writing code	◐
	Explaining code	◐
	Translating code	●
	Debugging code	◐
Data Analysis	Creating figures	◐
	Extracting data from text	●
	Reformatting data	●
	Classifying and scoring text	◐
	Extracting sentiment	◐
	Simulating human subjects	◐
Math	Setting up models	◐
	Deriving equations	○
	Explaining models	◐

## Ideation and Feedback

Category	Task	Usefulness
Ideation and Feedback	Brainstorming	●
	Feedback	◐
	Providing counterarguments	◐

# Ideation and Feedback: Feedback

## SYSTEM

You are a referee report for an economics journal who needs to evaluate a paper by looking at the abstract.

Highlight 3 strengths and 3 limitations of the paper which could undermine their conclusions, writing one sentence for each point.

## ASSISTANT

Strengths:

1. The paper utilizes a large dataset of over 250 million job vacancy postings across five English-speaking countries, providing a comprehensive analysis of the shift to remote work.
2. The authors employ a state-of-the-art language processing framework that achieves 99% accuracy in identifying job postings that advertise remote work, surpassing other methods.
3. The study covers a significant time period from 2019 to early 2023, allowing for an examination of the long-term trends and changes in the share of job postings offering remote work.

Limitations:

1. The paper focuses solely on English-speaking countries, which may limit the generalizability of the findings to other regions and cultures.
2. While the authors achieve high accuracy in identifying remote work postings, there may still be some misclassification or ambiguity in the language used, leading to potential measurement errors.
3. The study highlights non-uniformity in the adoption of remote work across cities, industries, occupations, and companies, but it does not delve into the underlying reasons for these differences, which could be crucial for understanding the broader implications of remote work.

# Writing

---

Writing	Synthesizing text	●
	Editing text	●
	Evaluating text	●
	Generating catchy titles & headlines	●
	Generating tweets to promote a paper	●

---

# Writing: Synthesizing Text

**USER**

Write a paragraph that integrates the following arguments and goes beyond. Include a topic sentence.

- \* large language models automate cognitive tasks
- \* the work of economists is largely cognitive [add examples]
- \* we economists should therefore worry about recent progress in language models
- \* in the short run, these models can help us be more productive
- \* in the long run, there is a good likelihood that they will displace human economists
- \* we should think how to prepare [offer solutions]

**ASSISTANT**

The rapid advancement of large language models, which are capable of automating cognitive tasks, has significant implications for the field of economics, where the work is predominantly cognitive in nature, such as analyzing data, forecasting trends, and developing policy recommendations. As economists, it is crucial that we acknowledge and address the potential impact of these language models on our profession. In the short term, these models can undoubtedly enhance our productivity by streamlining research, automating data analysis, and even generating insightful reports. However, in the long run, there is a considerable likelihood that these advanced technologies will displace human economists, as they become increasingly adept at performing complex tasks. To prepare for this eventuality, we must consider strategies such as continuous upskilling, focusing on areas where human intuition and creativity are irreplaceable, and advocating for responsible development and deployment of these technologies to ensure that they complement, rather than replace, our expertise.

# Background Research

---

Background Research	Summarizing Text	●
	Literature Research	○
	Formatting References	●
	Translating Text	●
	Explaining Concepts	◐

---



# Background Research: Literature Review

C

What does the economic literature say about gender differences in media representation of gender-based violence? Please provide me with the most important references from the economics literature (no need for a summary)



I can provide you with some important references from the economics literature on gender differences in media representation of gender-based violence. Please note that the availability of these references may vary, and it's always a good idea to access academic databases or libraries for the most up-to-date research. Here are some relevant references:

1. Anderson, D. A., & Dill, K. E. (2000). Video games and aggressive thoughts, feelings, and behavior in the laboratory and in life. *Journal of Personality and Social Psychology*, 78(4), 772-790.
2. Bandura, A., Ross, D., & Ross, S. A. (1963). Vicarious reinforcement and imitative learning. *The Journal of Abnormal and Social Psychology*, 67(6), 601-607.
3. Barroso, A. P., & Machin, S. (2011). Crime and the depenalization of cannabis possession: Evidence from a policing experiment. *The Economic Journal*, 121(552), 130-165.
4. Gentzkow, M., Shapiro, J. M., & Stone, D. F. (2015). Media slant against foreign owners: Downsizing. NBER Working Paper No. 20957.

# Background Research: Literature Review

C

What does the economic literature say about gender differences in media representation of gender-based violence? Please provide me with the most important references from the economics literature (no need for a summary)



I can provide you with some important references from the economics literature on gender differences in media representation of gender-based violence. Please note that the availability of these references may vary, and it's always a good idea to access academic databases or libraries for the most up-to-date research. Here are some relevant references:

1. Anderson, D. A., & Dill, K. E. (2000). Video games and aggressive thoughts, feelings, and behavior in the laboratory and in life. *Journal of Personality and Social Psychology*, 78(4), 772-790.
2. Bandura, A., Ross, D., & Ross, S. A. (1963). Vicarious reinforcement and imitative learning. *The Journal of Abnormal and Social Psychology*, 67(6), 601-607.
3. Barroso, A. P., & Machin, S. (2011). Crime and the depenalization of cannabis possession: Evidence from a policing experiment. *The Economic Journal*, 121(552), 130-165.
4. Gentzkow, M., Shapiro, J. M., & Stone, D. F. (2015). Media slant against foreign owners: Downsizing. NBER Working Paper No. 20957.

- ▶ ChatGPT invents papers
- ▶ Elicit does a much better job with this

# Coding

- ▶ ChatGPT does a decent job with writing code from scratch and debugging:
  - ▶ Good substitute to StackOverflow for some tasks
  - ▶ For more complicated problems needs a lot of back and forth (real life example)

# Coding

- ▶ ChatGPT does a decent job with writing code from scratch and debugging:
  - ▶ Good substitute to StackOverflow for some tasks
  - ▶ For more complicated problems needs a lot of back and forth (real life example)
  - ▶ Github Copilot is a **great** coding assistant

```
#write a function to load data from pickle
#and csv files depending on the ending of the file name
def load_data(file):
    if file.endswith('.pkl'):
        df = pd.read_pickle(file)
    elif file.endswith('.csv'):
        df = pd.read_csv(file)
    return df
```

# Coding

- ▶ ChatGPT does a decent job with writing code from scratch and debugging:
  - ▶ Good substitute to StackOverflow for some tasks
  - ▶ For more complicated problems needs a lot of back and forth (real life example)
  - ▶ Github Copilot is a **great** coding assistant

```
#write a function to load data from pickle
#and csv files depending on the ending of the file name
def load_data(file):
    if file.endswith('.pkl'):
        df = pd.read_pickle(file)
    elif file.endswith('.csv'):
        df = pd.read_csv(file)
    return df
```

- ▶ It does a great job translating code across languages

# Data Analysis

---

Data Analysis	Creating figures	○
	Extracting data from text	●
	Reformatting data	●
	Classifying and scoring text	○
	Extracting sentiment	○
	Simulating human subjects	○

---

# Data Analysis: Extract Data

## SYSTEM

You are an investigator assistant, who helps to extract information about crimes from Italian newspaper articles. In particular, I want you to answer the following questions: Does the article mention the murder of a human being that happened in real life in Italy (i.e., not in fiction)? Who is the victim of the murder? Who is the perpetrator or the accused perpetrator? When did the murder happen? What is the perpetrator to the victim (e.g., partner, ex-partner, relative, acquaintance, unknown)?

Answer the questions as concisely as possible and return the answers in a semicolon-separated CSV format with no more than 5 columns, where the column headers are: IS HOMICIDE; VICTIM(S)' FULL NAME; PERPETRATOR(S)' FULL NAME; DATE; RELATIONSHIP

The format of the date should be: %d %B %Y if day, month, and year are available; %d %B Year if only day and month are available; 1 Gennaio %Y if only the year is available; 1 %B %Y if only month and year are available.

IMPORTANT: If there are more than one victim and/or more than one perpetrator separate their names with "and".

If the article is not about a murder as defined before or the murder did not happen in Italy, return: NaN; NaN; NaN; NaN; NaN

Below are a few examples of prompts and responses:

## USER

He attempted to kill his ex-partner, 29 years old, and a Milanese tourist who was in the area and was fishing, and who helped her escape. He rammed into them, causing their car to go into a ravine, where they were at risk of drowning. This happened in Terracina, in the province of Latina. In the end, the police managed to apprehend and arrest the man, originally from Rosarno (Reggio Calabria).

Brought to Terracina and threatened with a knife, the 44-year-old, Michelangelo Porretta, who had already been reported by his ex for stalking and aggravated threats, had managed to set up a meeting with the woman, who was at her mother's house in Gaeta.

## ASSISTANT

yes; NaN; Michelangelo Porretta; NaN; ex-partner

🕒 Add message

# Data Analysis: Extract Data

## SYSTEM

You are an investigator assistant, who helps to extract information about crimes from Italian newspaper articles. In particular, I want you to answer the following questions: Does the article mention the murder of a human being that happened in real life in Italy (i.e., not in fiction)? Who is the victim of the murder? Who is the perpetrator or the accused perpetrator? When did the murder happen? What is the perpetrator to the victim (e.g., partner, ex-partner, relative, acquaintance, unknown)?

Answer the questions as concisely as possible and return the answers in a semicolon-separated CSV format with no more than 5 columns, where the column headers are: IS HOMICIDE; VICTIM(S)' FULL NAME; PERPETRATOR(S)' FULL NAME; DATE; RELATIONSHIP

The format of the date should be: %d %B %Y if day, month, and year are available; %d %B Year if only day and month are available; 1 Gennaio %Y if only the year is available; 1 %B %Y if only month and year are available.

IMPORTANT: If there are more than one victim and/or more than one perpetrator separate their names with "and".

If the article is not about a murder as defined before or the murder did not happen in Italy, return: NaN; NaN; NaN; NaN; NaN

Below are a few examples of prompts and responses:

## USER

He attempted to kill his ex-partner, 29 years old, and a Milanese tourist who was in the area and was fishing, and who helped her escape. He rammed into them, causing their car to go into a ravine, where they were at risk of drowning. This happened in Terracina, in the province of Latina. In the end, the police managed to apprehend and arrest the man, originally from Rosarno (Reggio Calabria).

Brought to Terracina and threatened with a knife, the 44-year-old, Michelangelo Porretta, who had already been reported by his ex for stalking and aggravated threats, had managed to set up a meeting with the woman, who was at her mother's house in Gaeta.

## ASSISTANT

yes; NaN; Michelangelo Porretta; NaN; ex-partner

🕒 Add message



# Math

Math	Setting up models	●
	Deriving equations	○
	Explaining models	●

- Integration with Wolfram Mathematica works well

## Some Limitations...

- ▶ **Weaker in analytic concepts** due to their *nature* as LLM

## Some Limitations...

- ▶ **Weaker in analytic concepts** due to their *nature* as LLM
- ▶ **Limit on training data** up to 2021 and **fixed input** length

## Some Limitations...

- ▶ **Weaker in analytic concepts** due to their *nature* as LLM
- ▶ **Limit on training data** up to 2021 and **fixed input** length
- ▶ **Limits to reproducibility**
- ▶ **Hallucinations**: LLMs can easily make things up, which limits how much we can leverage their knowledge base
- ▶ **Bias**

## ... but a Lot of Potential Applications!

- ▶ Fixing OCR errors in digitized text
- ▶ Extract information about homicides based on newspaper articles
- ▶ Tag articles as believing/skeptical of climate change
- ▶ Identifying “heroes”, “villans”, and “victims”
- ▶ Interpreting old proverbs

## ... but a Lot of Potential Applications!

- ▶ Fixing OCR errors in digitized text
- ▶ Extract information about homicides based on newspaper articles
- ▶ Tag articles as believing/skeptical of climate change
- ▶ Identifying “heroes”, “villans”, and “victims”
- ▶ Interpreting old proverbs
- ▶ Other ideas...?

# Outline

GPT

GPT for Research

Bias in NLP Systems

# Bias in NLP Systems

## Sentiment Analysis

```
text_to_sentiment("Let's go get Italian food")  
2.0429166109  
text_to_sentiment("Let's go get Chinese food")  
1.4094033658  
text_to_sentiment("Let's go get Mexican food")  
0.3880198556
```

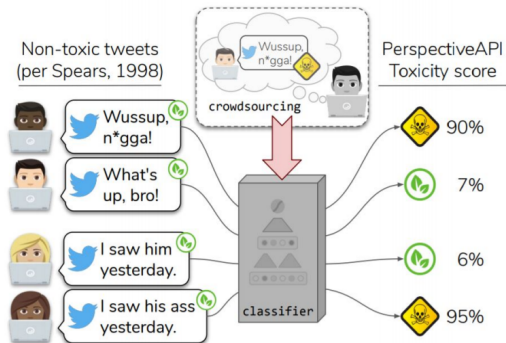
```
text_to_sentiment("My name is Emily")  
2.2286179365  
text_to_sentiment("My name is Heather")  
1.3976291151  
text_to_sentiment("My name is Yvette")  
0.9846380213  
text_to_sentiment("My name is Shaniqua")  
-0.4704813178
```

**Is this sentiment model racist?**



# Bias in NLP Systems

## Toxicity Detection



Within dataset proportions

DWMW17	% false identification				
	Group	Acc.	None	Offensive	Hate
	AAE	94.3	1.1	<b>46.3</b>	0.8
	White	87.5	<b>7.9</b>	9.0	<b>3.8</b>
	Overall	91.4	2.9	17.9	2.3

			% false identification		
FDCL18	Group	Acc.	None	Abusive	Hateful
	AAE	81.4	4.2	<b>26.0</b>	<b>1.7</b>
	White	82.7	<b>30.5</b>	4.5	0.8
	Overall	81.4	20.9	6.6	0.8

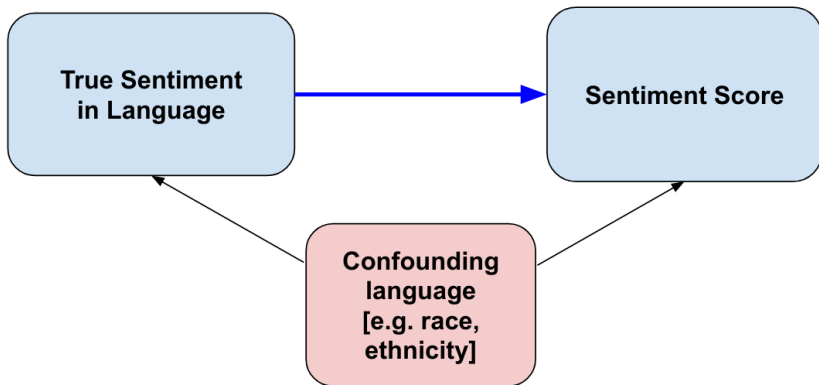
Is this toxicity detection model racist?

## NLP “Bias” is statistical bias

- ▶ Sentiment scores that are trained on annotated datasets also learn from the correlated non-sentiment information.

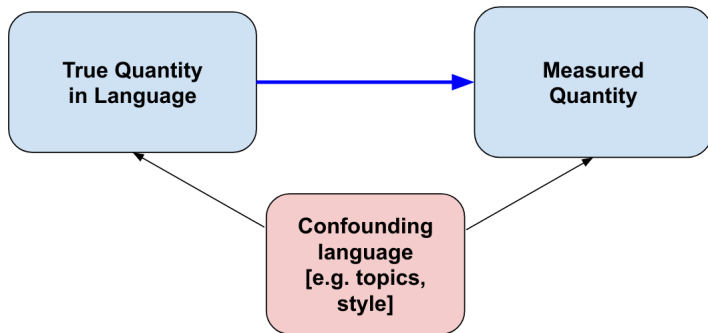
## NLP “Bias” is statistical bias

- ▶ Sentiment scores that are trained on annotated datasets also learn from the correlated non-sentiment information.



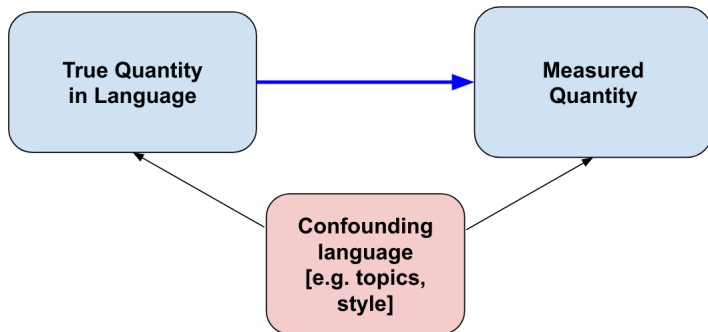
- ▶ Supervised sentiment models are confounded by correlated language factors.
  - ▶ e.g., in the training set maybe people complain about Mexican food more often than Italian food.

## This is a universal problem



- ▶ supervised models (classifiers, regressors) learn features that are correlated with the label being annotated.
- ▶ unsupervised models (topic models, word embeddings) learn correlations between topics / contexts.

## This is a universal problem



- ▶ supervised models (classifiers, regressors) learn features that are correlated with the label being annotated.
- ▶ unsupervised models (topic models, word embeddings) learn correlations between topics / contexts.
- ▶ An important exception: dictionary methods (perhaps explaining why they are often used by economists). But they have other serious limitations.

## Examples: Confounders in Measurement from Text

**What quantity do we care about? vs. What do we measure?**

## Examples: Confounders in Measurement from Text

What quantity do we care about? vs. What do we measure?

- ▶ Positive/negative sentiment → Count positive/negative words, or predict text annotations.
- ▶ Toxicity → Count toxic words, or predict text annotations.

**confounders?**

## Examples: Confounders in Measurement from Text

What quantity do we care about? vs. What do we measure?

- ▶ Positive/negative sentiment → Count positive/negative words, or predict text annotations.
- ▶ Toxicity → Count toxic words, or predict text annotations.

**confounders?**

- ▶ Student performance → predicted essay grade based on labeled essay documents.
- ▶ Credit worthiness → predicted probability of default based on loan application documents.

**confounders?**



## Examples: Confounders in Measurement from Text

What quantity do we care about? vs. What do we measure?

- ▶ Positive/negative sentiment → Count positive/negative words, or predict text annotations.
- ▶ Toxicity → Count toxic words, or predict text annotations.

**confounders?**

- ▶ Student performance → predicted essay grade based on labeled essay documents.
- ▶ Credit worthiness → predicted probability of default based on loan application documents.

**confounders?**

- ▶ Political partisanship → predicted probability being Democrat/Republican based on speeches.

**confounders?**

## Examples: Confounders in Measurement from Text

**What quantity do we care about? vs. What do we measure?**

- ▶ Positive/negative sentiment → Count positive/negative words, or predict text annotations.
- ▶ Toxicity → Count toxic words, or predict text annotations.

**confounders?**

- ▶ Student performance → predicted essay grade based on labeled essay documents.
- ▶ Credit worthiness → predicted probability of default based on loan application documents.

**confounders?**

- ▶ Political partisanship → predicted probability being Democrat/Republican based on speeches.

**confounders?**

- ▶ Policy priorities → predicted probability of speeches/laws being about a particular policy topic.

**confounders?**

## When is measurement confounding important?

- ▶ By itself, producing measurements that are biased by confounders might not be a problem.
- ▶ e.g.:
  - ▶ an NLP-based credit score that learns confounders → not a problem unless debtors learn about it and strategically alter their documents.
  - ▶ similarly with automated essay grading

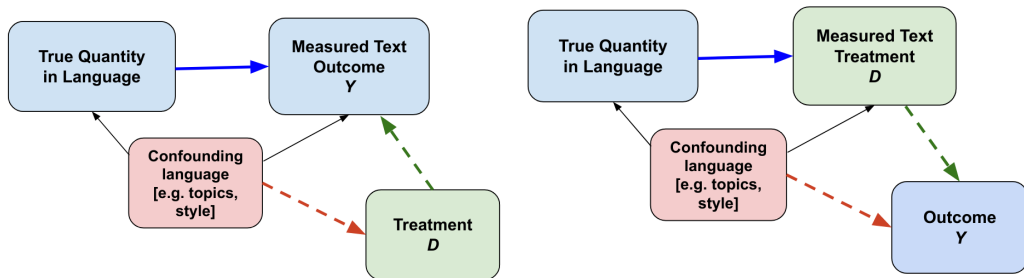
## When is measurement confounding important?

- ▶ By itself, producing measurements that are biased by confounders might not be a problem.
- ▶ e.g.:
  - ▶ an NLP-based credit score that learns confounders → not a problem unless debtors learn about it and strategically alter their documents.
    - ▶ similarly with automated essay grading
- ▶ for measuring political divisiveness or policy priorities
  - ▶ probably won't matter for in-domain summary statistics
  - ▶ but would matter a lot for summary statistics in a new domain

## When is measurement confounding important?

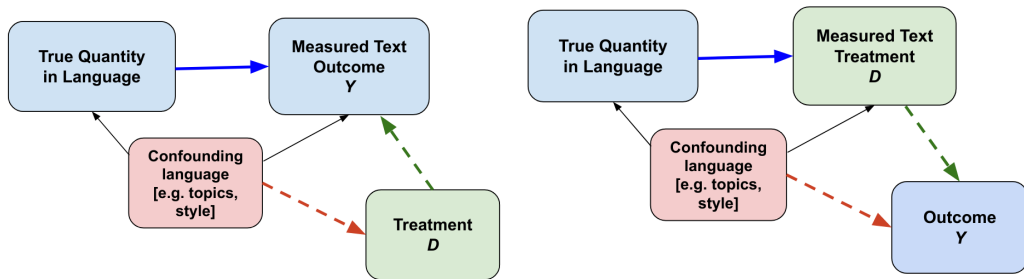
- ▶ By itself, producing measurements that are biased by confounders might not be a problem.
- ▶ e.g.:
  - ▶ an NLP-based credit score that learns confounders → not a problem unless debtors learn about it and strategically alter their documents.
  - ▶ similarly with automated essay grading
- ▶ for measuring political divisiveness or policy priorities
  - ▶ probably won't matter for in-domain summary statistics
  - ▶ but would matter a lot for summary statistics in a new domain
- ▶ even in domain, will matter for assessing the causal effect of a treatment, e.g. the electoral cycle:
  - ▶ elections might cause politicians to focus on social issues rather than economic issues,
  - ▶ if social/economic issues are confounded with partisanship, the resulting estimates are biased.

## When is measurement confounding important?



- ▶ When text is outcome, the confounders cannot be correlated with the treatment.
- ▶ When text is treatment, the confounders cannot be correlated with the outcome.

## When is measurement confounding important?



- ▶ When text is outcome, the confounders cannot be correlated with the treatment.
- ▶ When text is treatment, the confounders cannot be correlated with the outcome.
  - ▶ e.g.: estimating the effect of politician speech sentiment on his/her reelection chances?

## Steps for de-biasing

- ▶ Language features that are often confounded with the quantity of interest:
  - ▶ stopwords
  - ▶ named entities: person/organization/place names
- ▶ These can be dropped during pre-processing to reduce the influence of confounders in subsequent measurements.
- ▶ Can control for topic or style features or other potential confounders in regressions.