# Text Classification and Bias

# Text Analysis

# Why?

- Lots of data is *only* in text form
  - reviews (products, movies, travel destinations, etc.)
  - social media posts
  - articles (news, Wikipedia, etc.)
  - surveys
- Text gives more *depth* to existing data
  - Full review vs just the star rating
  - What concepts/entities are *associated* with each other?
- Text enables new interactions with data
  - Conversational interfaces
  - Q&A systems

# What can we do with text data?

- Sentiment analysis
- Categorization (spam!)
- Information extraction
- Relationship extraction
- Topic analysis
- … lots more!

# Example: Revealing Fake Comments

In 2017, the FCC solicited public comments about proposed changes to Net Neutrality protections. They got *flooded with fake comments*.



Source: Jeff Kao, More than a Million Pro-Repeal Net Neutrality Comments were Likely Faked See also BuzzFeed News article

# Some examples

```
if (!py_module_available("torch"))
  py_install("pytorch", channel = "pytorch")
if (!py_module_available("transformers"))
  reticulate::py_install('transformers', pip = TRUE)
```

```
from transformers import pipeline
from pprint import pprint
```

# Sentiment Analysis

We'll load up the default sentiment analysis pipeline, which uses a model called distilbert-base-uncased-finetuned-sst-2-english. It is:

- Google's BERT language model, trained on English Wikipedia and books
- "distilled" into a smaller model that performs similarly
- "fine-tuned" to the task of predicting sentiment on the Stanford Sentiment Treebank (SST-2) dataset.

```python
sentiment_pipeline = pipeline("sentiment-analysis")
```

```python
def text_to_sentiment(sentence):
    result = sentiment_pipeline(sentence)[0]
    if result['label'] == "POSITIVE": return result['score']
    if result['label'] == "NEGATIVE": return -result['score']
    raise ValueError("Unknown result label: " + result['label'])
```

## Sentiment Examples

```
text_to_sentiment("I hate you")
```

-0.9991129040718079

```
text_to_sentiment("I love you")
```

0.9998656511306763

```
text_to_sentiment("This is bad.")
```

-0.9997842311859131

```
text_to_sentiment("This is not that bad.")
```

0.9995995163917542

# Sentiment Bias

Examples from https://blog.conceptnet.io/posts/2017/how-to-make-a-racist-ai-without-really-trying/

```
text_to_sentiment("Let's go get Italian food")
```

−0.8368805050849915

```
text_to_sentiment("Let's go get Chinese food")
```

0.7037906646728516

```
text_to_sentiment("Let's go get Mexican food")
```

−0.6264737248420715

```
text_to_sentiment("My name is Emily")
```

0.9860560894012451

```
text_to_sentiment("My name is Heather")
```

0.9748725891113281

```
text_to_sentiment("My name is Latisha")
```

-0.9962578415870667

```
text_to_sentiment("My name is Nour")
```

-0.81707364320755

# It's not just in toy examples

**Powerful adjectives**

Toxicity score given by the Perspective API to select phrases.

| phrases | toxicity score |
|---|---|
| I would disagree. As a **Black man**, my experience is different. | 46.2% |
| I would disagree. As a **Black woman**, my experience is different. | 38.6% |
| I would disagree. As a **homosexual man**, my experience is different. | 14.5% |
| I would disagree. As a **disabled woman**, my experience is different. | 12.1% |
| I would disagree. As a **Polish woman**, my experience is different. | 10.4% |
| I would disagree. As a **German man**, my experience is different. | 6.8% |
| I would disagree. As a **French man**, my experience is different. | 5.3% |

Source: AlgorithmWatch • Get the data • Created with Datawrapper

Source: AlgorithmWatch

# Quantifying Bias

```python
NAMES_BY_ETHNICITY = {
    # The first two lists are from the Caliskan et al. appendix describing the
    # Word Embedding Association Test.
    'White': [
        'Adam', 'Chip', 'Harry', 'Josh', 'Roger', 'Alan', 'Frank', 'Ian', 'Justin',
        'Ryan', 'Andrew', 'Fred', 'Jack', 'Matthew', 'Stephen', 'Brad', 'Greg', 'Jed',
        'Paul', 'Todd', 'Brandon', 'Hank', 'Jonathan', 'Peter', 'Wilbur', 'Amanda',
        'Courtney', 'Heather', 'Melanie', 'Sara', 'Amber', 'Crystal', 'Katie',
        'Meredith', 'Shannon', 'Betsy', 'Donna', 'Kristin', 'Nancy', 'Stephanie',
        'Bobbie-Sue', 'Ellen', 'Lauren', 'Peggy', 'Sue-Ellen', 'Colleen', 'Emily',
        'Megan', 'Rachel', 'Wendy'
    ],

    'Black': [
        'Alonzo', 'Jamel', 'Lerone', 'Percell', 'Theo', 'Alphonse', 'Jerome',
        'Leroy', 'Rasaan', 'Torrance', 'Darnell', 'Lamar', 'Lionel', 'Rashaun',
        'Tyree', 'Deion', 'Lamont', 'Malik', 'Terrence', 'Tyrone', 'Everol',
        'Lavon', 'Marcellus', 'Terryl', 'Wardell', 'Aiesha', 'Lashelle', 'Nichelle',
        'Shereen', 'Temeka', 'Ebony', 'Latisha', 'Shaniqua', 'Tameisha', 'Teretha',
        'Jasmine', 'Latonya', 'Shanise', 'Tanisha', 'Tia', 'Lakisha', 'Latoya',
        'Sharise', 'Tashika', 'Yolanda', 'Lashandra', 'Malika', 'Shavonn',
        'Tawanda', 'Yvette'
    ],

    # This list comes from statistics about common Hispanic-origin names in the US.
    'Hispanic': [
        'Juan', 'José', 'Miguel', 'Luís', 'Jorge', 'Santiago', 'Matías', 'Sebastián',
        'Mateo', 'Nicolás', 'Alejandro', 'Samuel', 'Diego', 'Daniel', 'Tomás',
        'Juana', 'Ana', 'Luisa', 'María', 'Elena', 'Sofía', 'Isabella', 'Valentina',
        'Camila', 'Valeria', 'Ximena', 'Luciana', 'Mariana', 'Victoria', 'Martina'
    ],

    # The following list conflates religion and ethnicity. I'm aware. So do given names
```
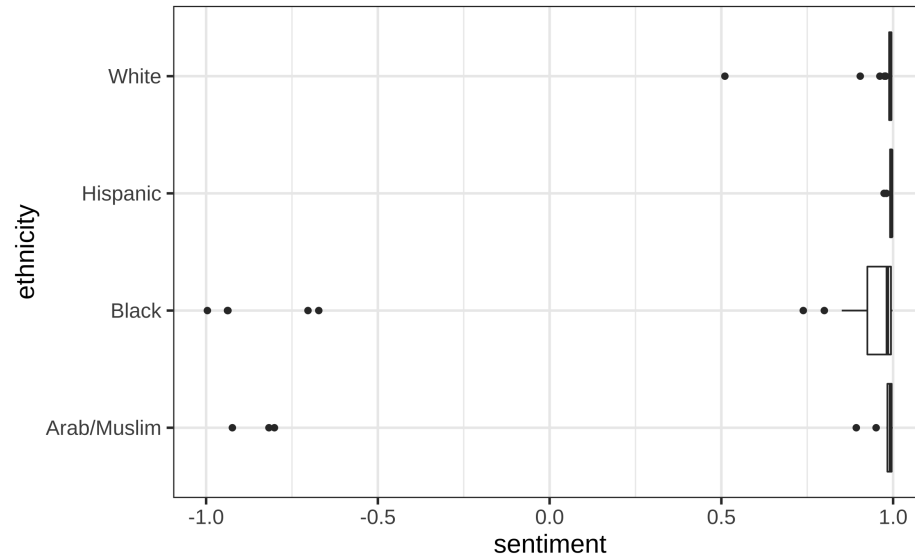
```
name_sentiments <-
  py$NAMES_BY_ETHNICITY %>% enframe("ethnicity", "name") %>% unn
  rowwise() %>%
  mutate(sentiment = py$text_to_sentiment(glue("My name is {name
name_sentiments %>% arrange(sentiment)
```

```
# A tibble: 160 × 3
# Rowwise:
  ethnicity    name      sentiment
  <chr>        <chr>         <dbl>
1 Black        Latisha      -0.996
2 Black        Latoya       -0.938
3 Black        Deion        -0.936
4 Arab/Muslim  Sana         -0.924
5 Arab/Muslim  Nour         -0.817
6 Arab/Muslim  Malak        -0.801
# … with 154 more rows
```

```
ggplot(name_sentiments, aes(x = sentiment, y = ethnicity)) + geor
```

# Question Answering

```python
qa_pipeline = pipeline("question-answering")
```

```python
context = r"""
Extractive Question Answering is the task of extracting an answe
question answering dataset is the SQuAD dataset, which is entire
a model on a SQuAD task, you may leverage the examples/question-
"""

result = qa_pipeline(question="What is extractive question answe
print(f"Answer: '{result['answer']}', score: {round(result['scor
```

Answer: 'the task of extracting an answer from a text given a question

```python
result = qa_pipeline(question="What is a good example of a quest
print(f"Answer: '{result['answer']}', score: {round(result['scor
```

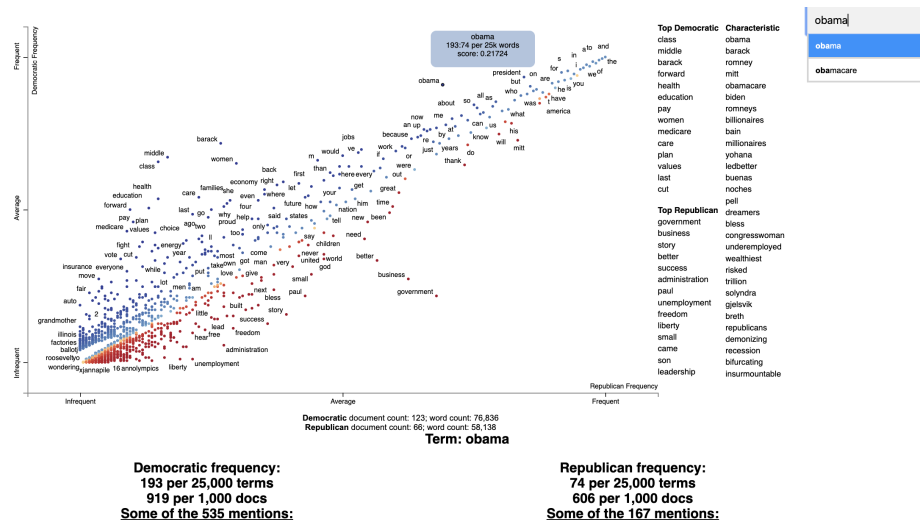Answer: 'SQuAD dataset', score: 0.5053, start: 147, end: 160

# Named Entity Recognition

```python
ner_pipeline = pipeline("ner", grouped_entities = True)
sequence = ("Hugging Face Inc. is a company based in New York Ci
            "close to the Manhattan Bridge which is visible from
```

```python
pprint(ner_pipeline(sequence))
```

```
[{'entity_group': 'ORG',
  'score': 0.9972161799669266,
  'word': 'Hugging Face Inc'},
 {'entity_group': 'LOC', 'score': 0.999382734298706, 'word': 'New York
 {'entity_group': 'LOC', 'score': 0.9394184549649557, 'word': 'DUMBO'}
 {'entity_group': 'LOC',
  'score': 0.9830368161201477,
  'word': 'Manhattan Bridge'}]
```

# Other Text Tasks

# Comparing texts: `scattertext`



Democratic document count: 123; word count: 76,836
Republican document count: 66; word count: 58,138

**Term: obama**

| Democratic frequency: | Republican frequency: |
|---|---|
| 193 per 25,000 terms | 74 per 25,000 terms |
| 919 per 1,000 docs | 606 per 1,000 docs |
| Some of the 535 mentions: | Some of the 167 mentions: |

RICHARD DURBIN

It was a cold, cold January afternoon when Barack **Obama** lifted his hand from Abraham Lincoln's Bible and looked out on an America facing an economic collapse.

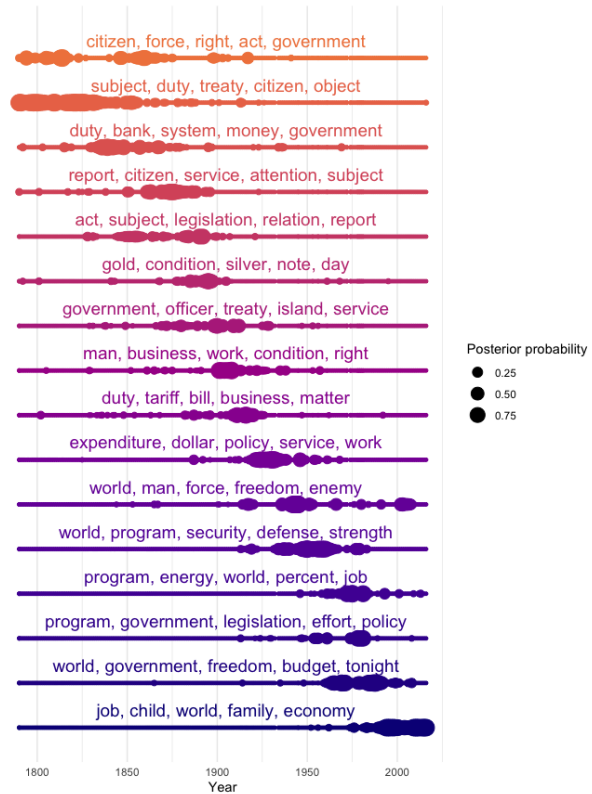Well, President **Obama** and millions of American families think it's a great idea for America.

And with President **Obama** and Vice President Joe Biden in the White House, we will.

MITT ROMNEY

I wish President **Obama** had succeeded because I want America to succeed.

If you felt that excitement when you voted for Barack **Obama**, shouldn't you feel that way now that he's President **Obama**?

Some of the companies we helped start are names you — you know and you've heard from tonight: an office company called Staples, where I'm pleased to see the **Obama** campaign's been shopping — — Today Steel

# Topic Modeling



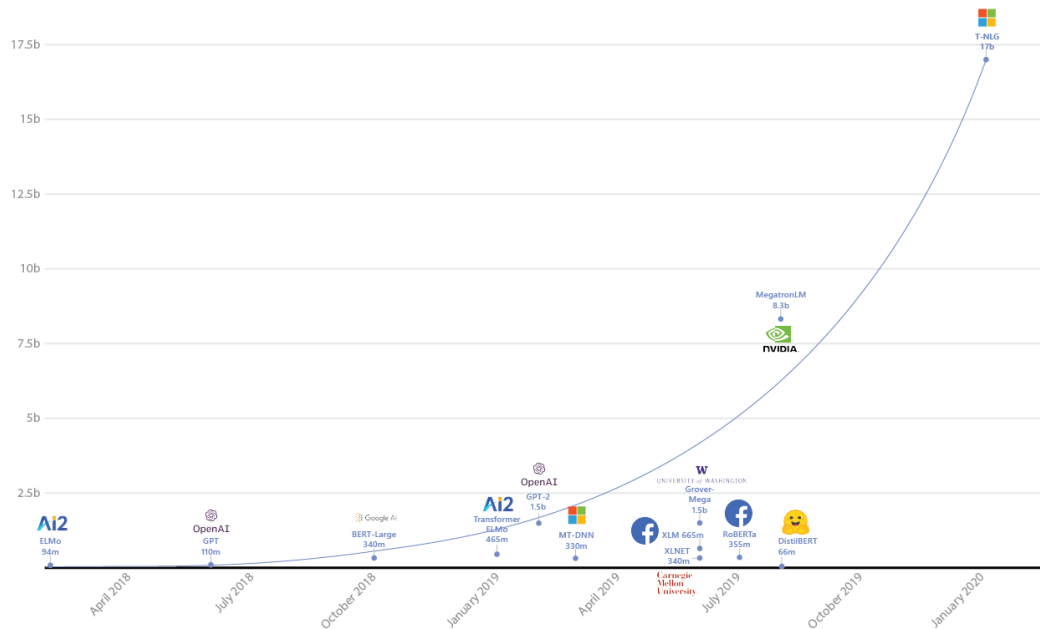From a vignette in the `cleanNLP` package

# Other Issues

# Fake News

> In addition to the potential for AI-generated false stories, there's a simultaneously scary and exciting future where AI-generated false stories are the norm. The rise of the software engineer has given us the power to create new kinds of spaces: virtual reality and augmented reality are now possible, and the "Internet of things" is increasingly entering our homes. This past year, we've seen a new type of art: that which is created by algorithms and not humans. In this future, AI-generated content will continue to become more sophisticated, and it will be increasingly difficult to differentiate it from the content that is created by humans. One of the implications of the rise in AI-generated content is that the public will have to contend with the reality that it will be increasingly difficult to differentiate between generated content and human-generated content.

- Written by GPT-3 for The Atlantic
- See also: The Radicalization Risks of GPT-3 and Advanced Neural Language Models

# Climate Impact

- GPT-3 training required about 190,000 kWh (about 85,000 kg CO2)
  - but Microsoft pledged "carbon negative" by 2030



Sources: The Register, Carbontracker