



What Makes the News?

A Topic Modeling Project by Max Lan

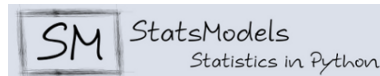


The State of the Media

- US Newsroom employment declined 25% from 2008 to 2018
- 36% of young adults use social media for news; only 27% use traditional news websites
- What do news publications currently write about and where can they expand?

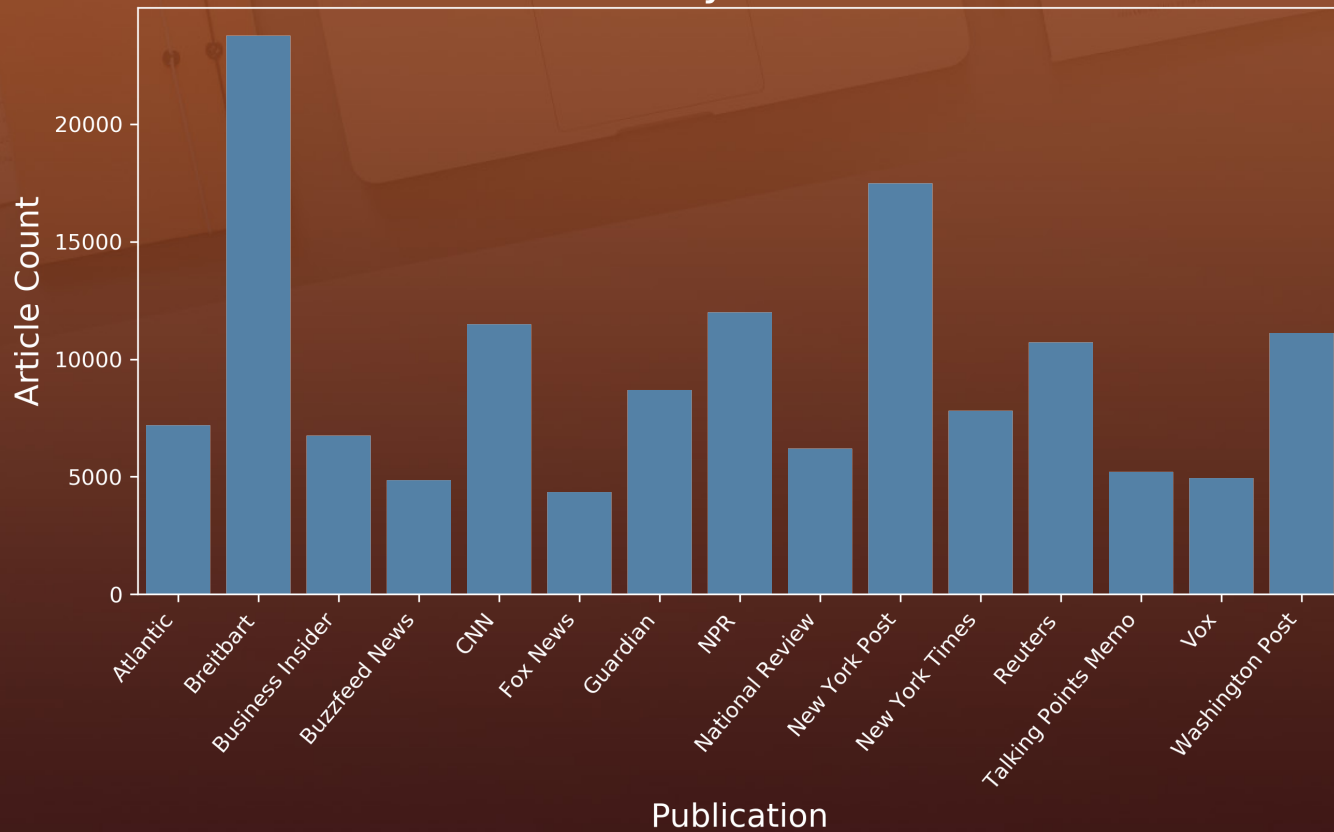
Methodology

- Data: 143,000 articles from 15 news publications on Kaggle
- Topic Modeling: CountVectorizer, LDA, and A/B Testing
- Readability Analysis with spaCy readability
- Primary Packages used:





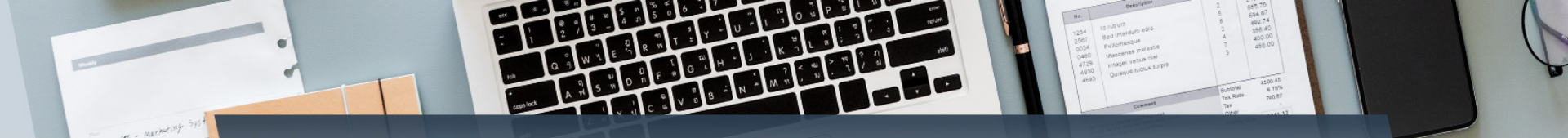
Articles by Publication





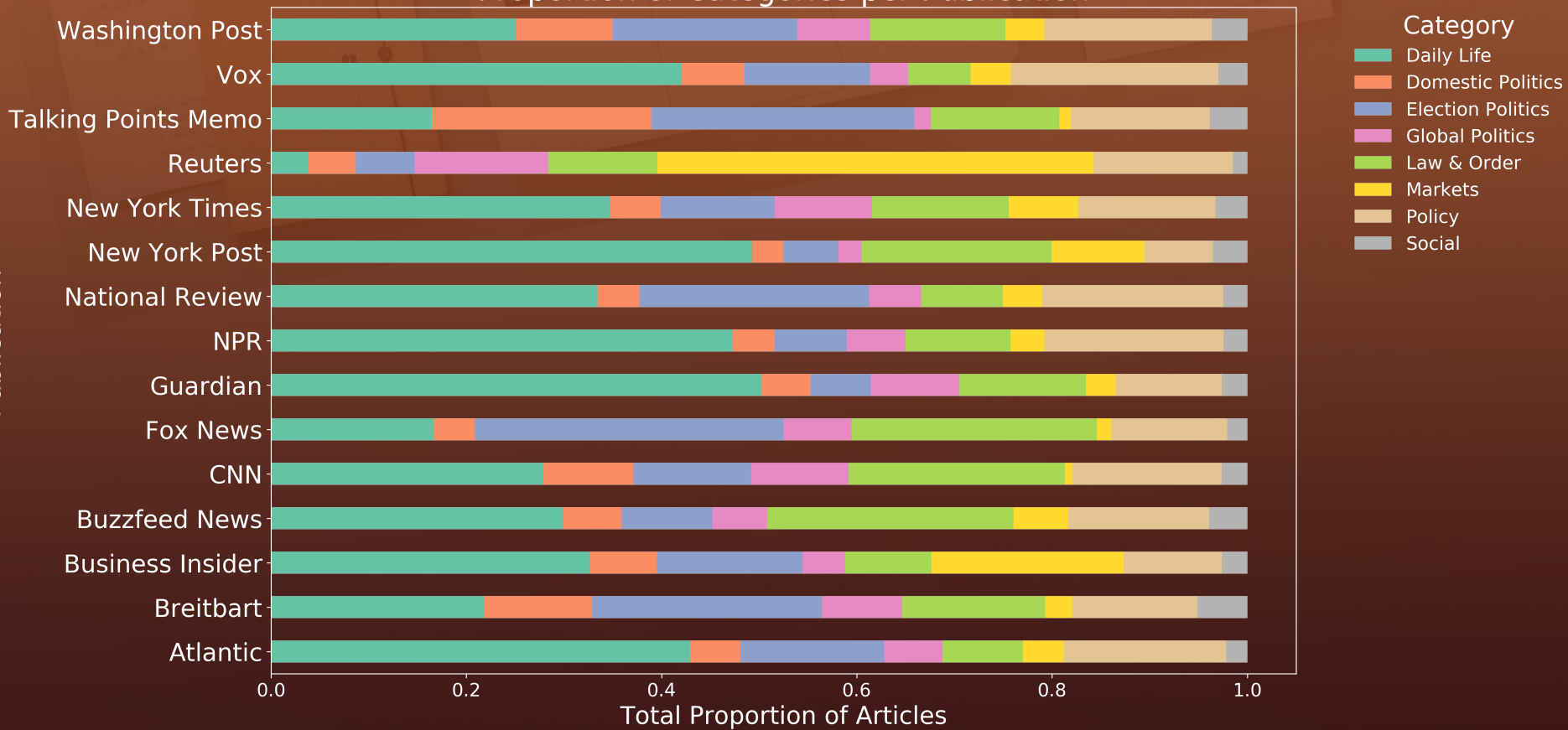
Topics and Categories

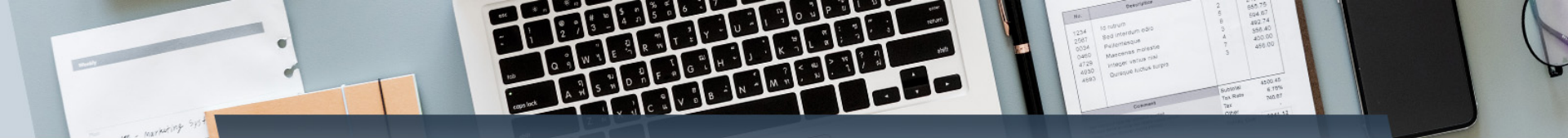
- 24 topics from LDA
 - Most relevant words used
 - pyLDavis
- 8 primary categories from 24 topics
 - Cosine similarity
- Examples:
 - Topic:
 - Law: court, case, law, justice
 - Police: police officer, shooting
 - Categories:
 - Daily Life: sports, lifestyle, social media



Proportion of Categories per Publication

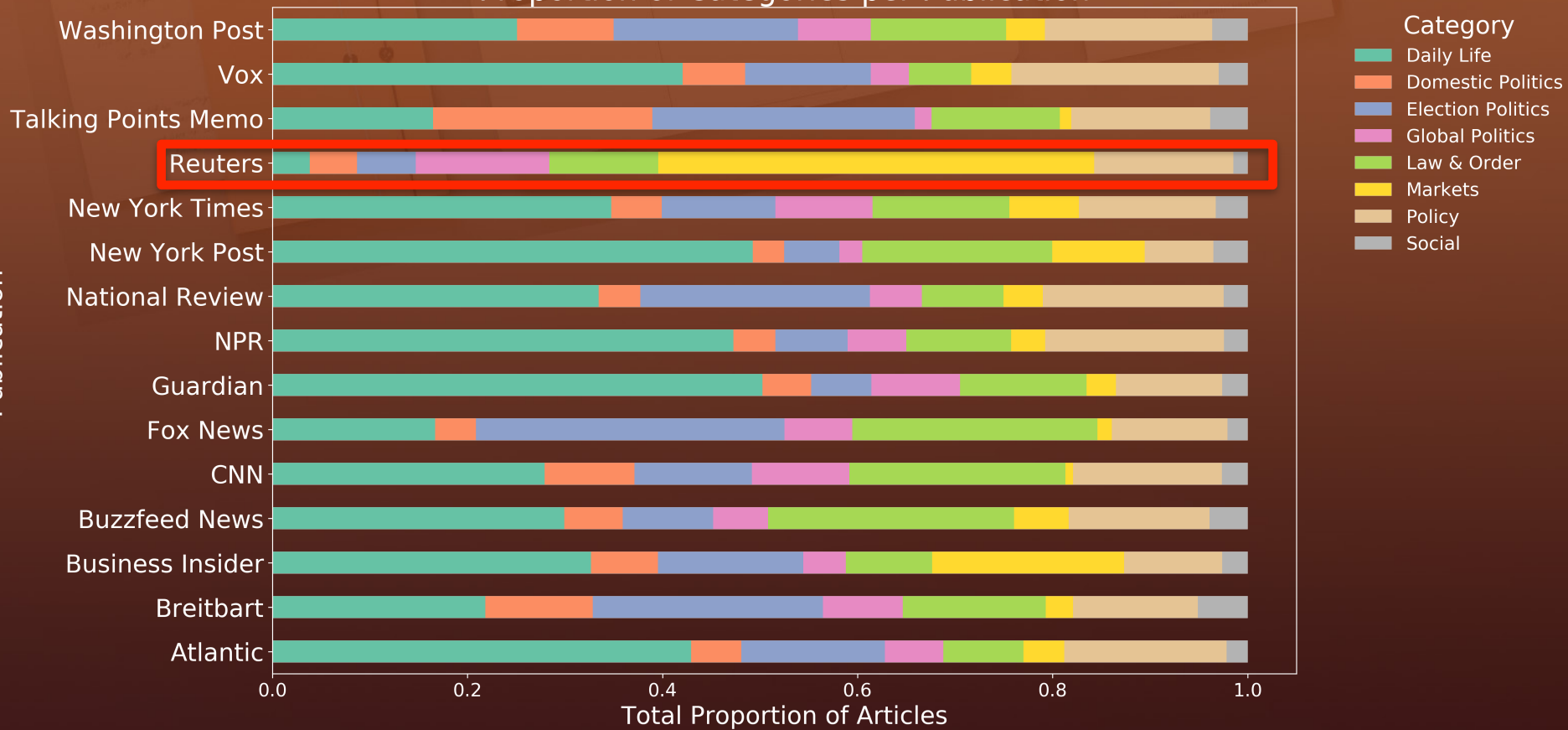
Publication





Proportion of Categories per Publication

Publication



Categories - A Statistical Approach

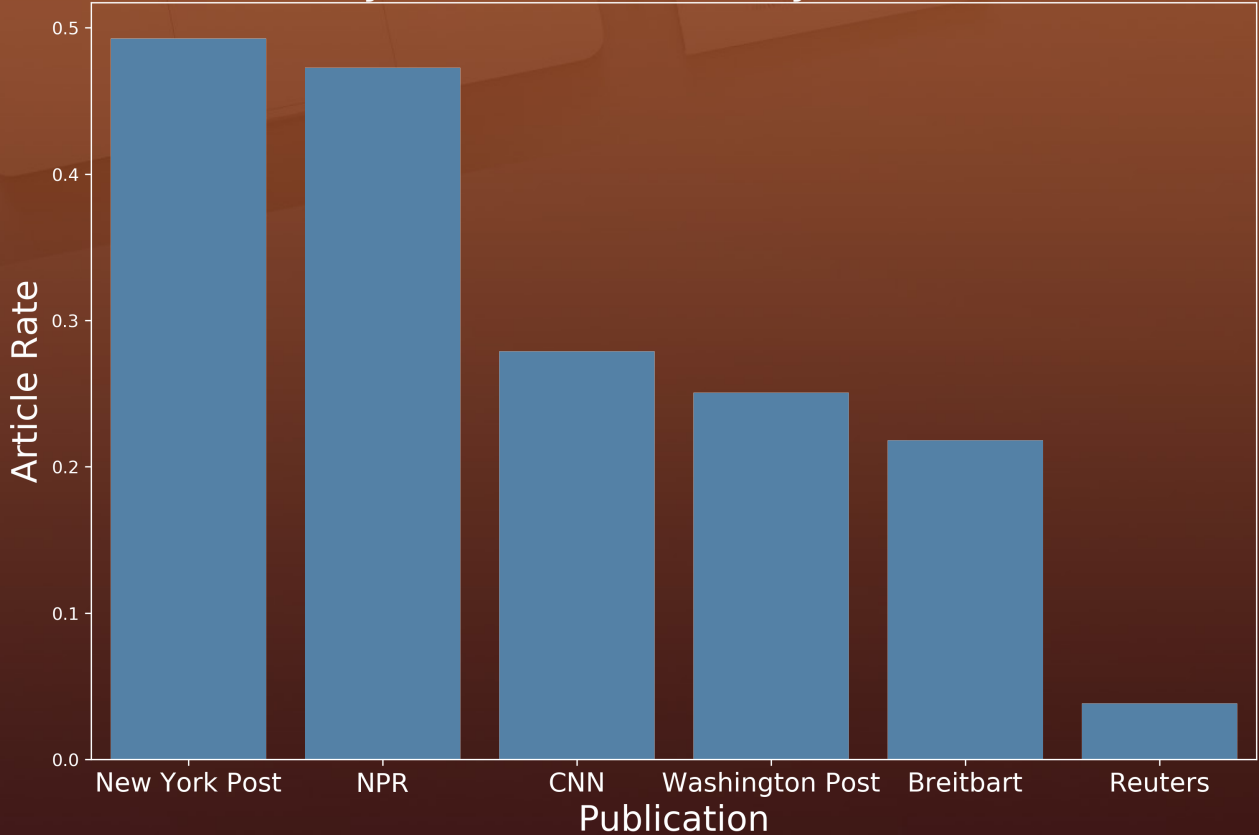
- Multiple hypothesis A/B testing by category
 - Top 6 publications by article count
 - Holm-Bonferroni Correction
- Determine statistically significant differences in article proportion





Category	Test Result
Daily Life	Pass
Domestic Politics	Fail
Election Politics	Fail
Global Politics	Fail
Law & Order	Fail
Markets	Fail
Policy	Fail
Social	Fail

Daily Life Article Rate by Publication

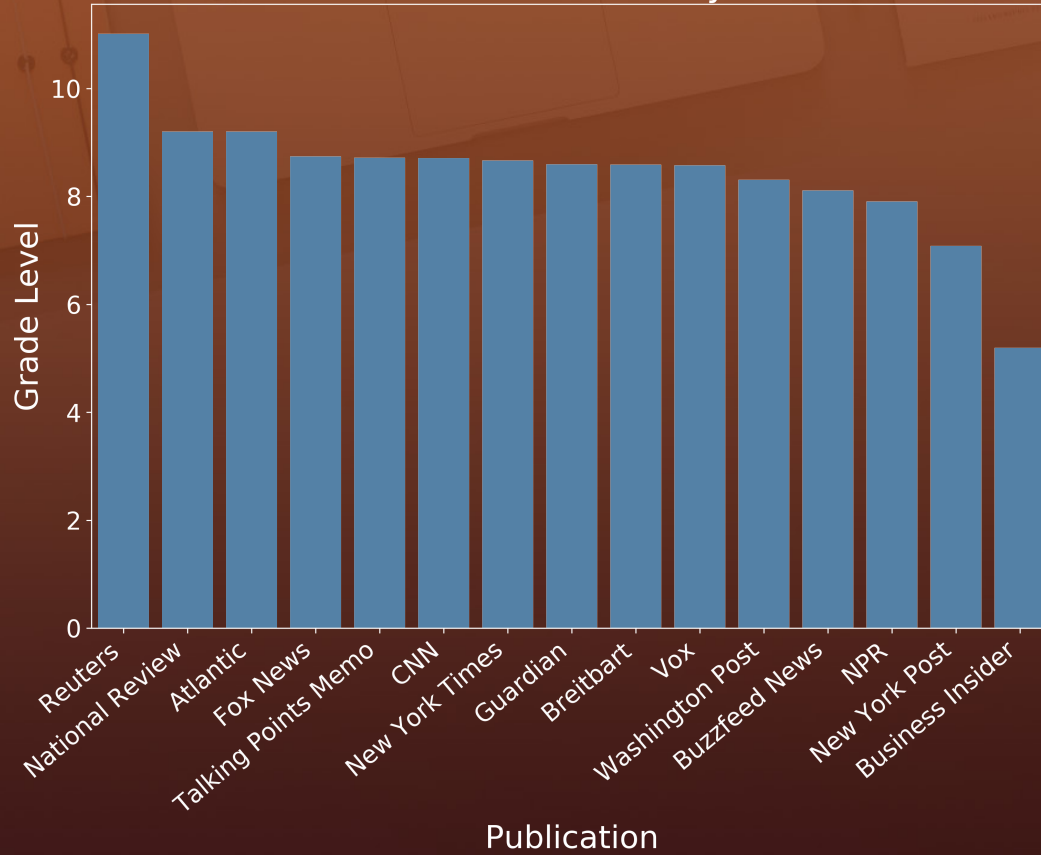


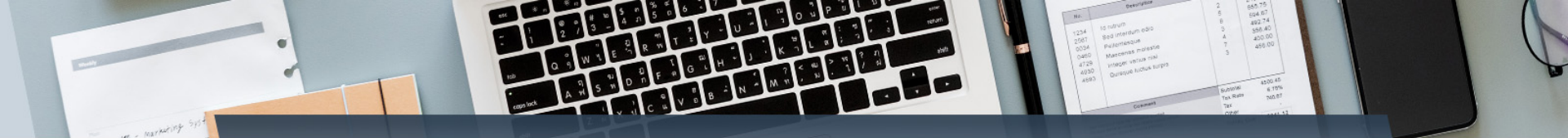


Readability Analysis

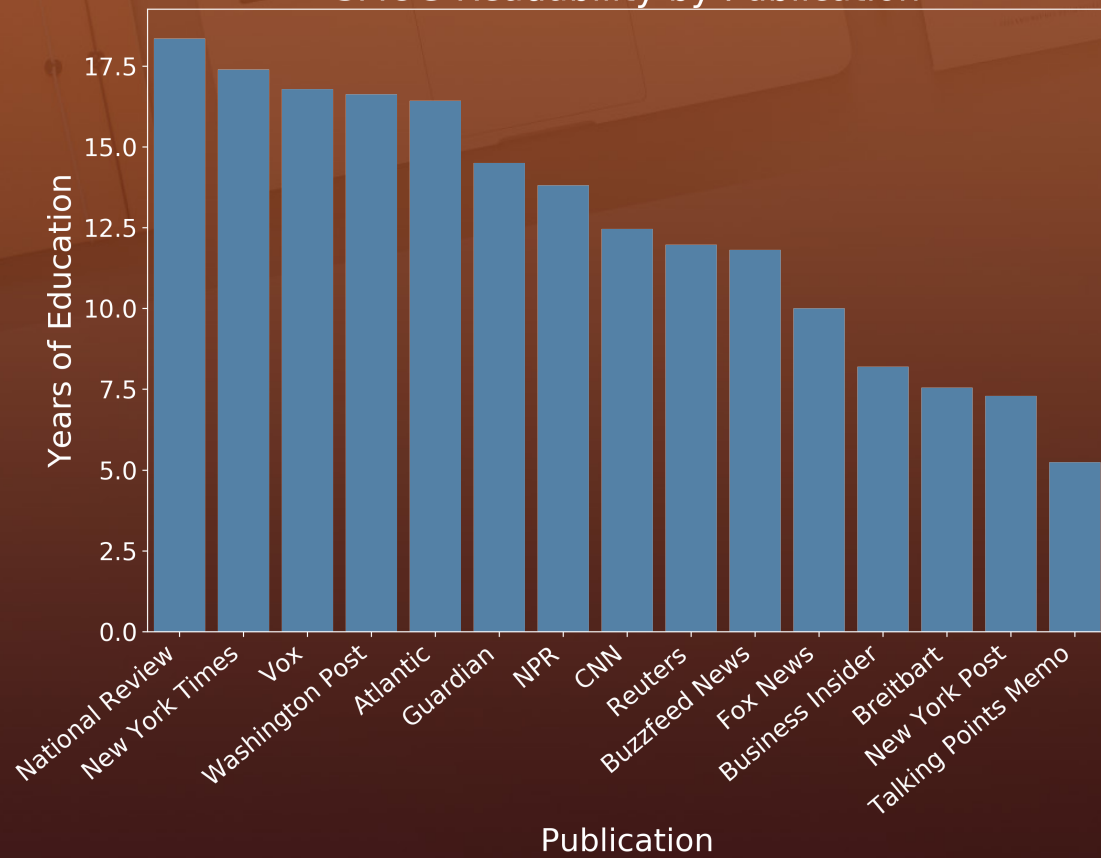
- Examine publications by multiple readability metrics
- Determine if publications focus on different writing levels

Flesh-Kincaid Grade Level by Publication





SMOG Readability by Publication



A photograph of a desk setup. In the center is a laptop keyboard. To the left is a pen and a document with a table. To the right is a smartphone. The background is a solid brown color.

Recommendations

- Daily life is a differentiator:
 - Reuters: more reader-friendly publication with daily life articles
 - New York Post, NPR: consider separate focus on other topics
- Reading level: cater to your audience
 - Breitbart and Talking Points Memo: consider longer/more sophisticated articles or pieces to expand



Additional Topics

- Obtaining more articles across more publications
- Creating a similar article recommender
- Identifying plagiarism between news publications



Thank You!

- Github: https://github.com/mlan93/project_04_NLP
- LinkedIn: <https://www.linkedin.com/in/maxlan93/>
-