PREDICTING WIKIPEDIA ARTICLE TRANSLATION COUNTS



OUTLINE

Contents [hide]

- 1 Introduction
- 2 Models
 - 2.1 Manually-curated Features
 - 2.2 Natural Language Processing (NLP)
 - 2.3 Combined Approach
- 3 Conclusions and Future Directions

Wikipedia is the 13th most visited website worldwide

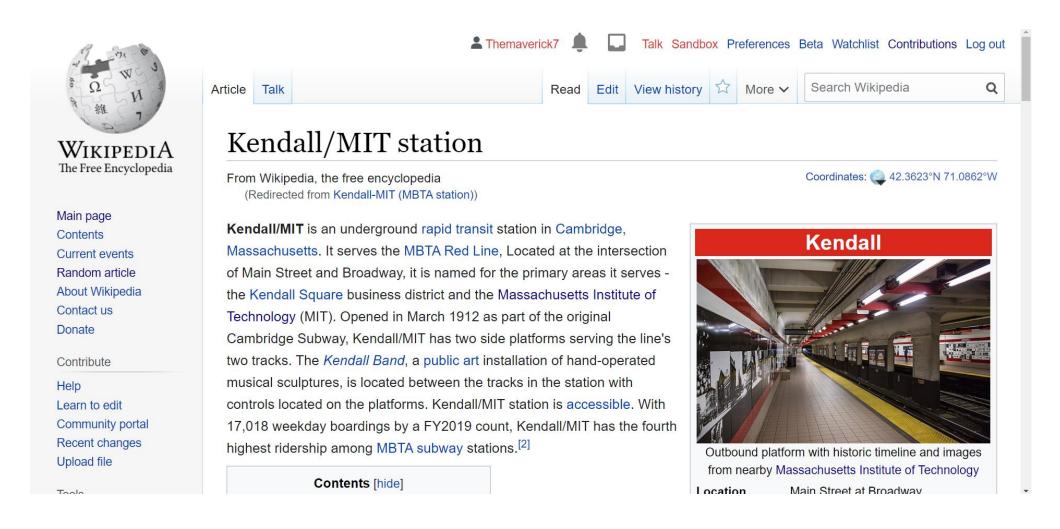


Wikipedia is the 13th most visited website worldwide

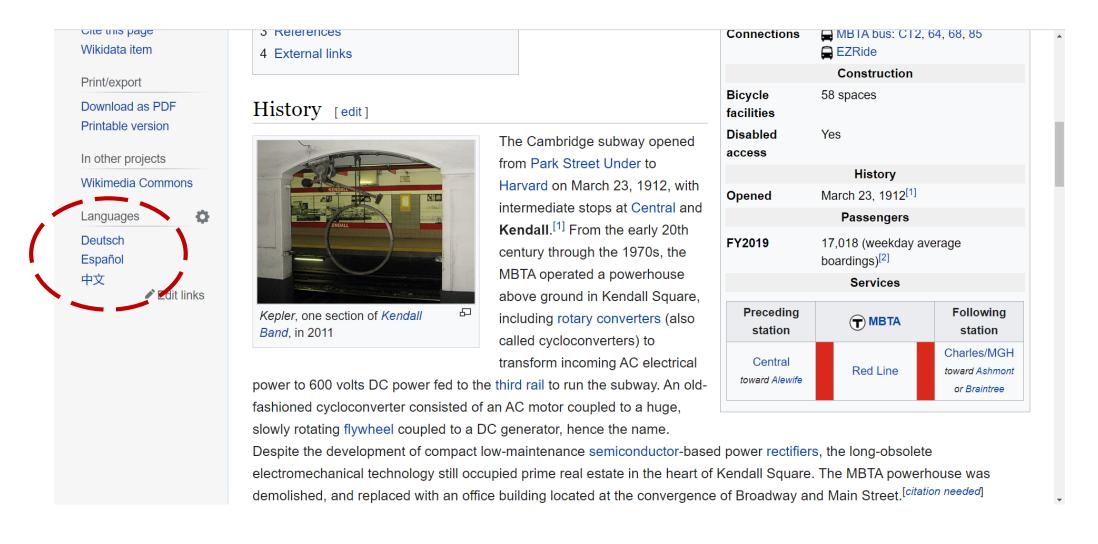
- Founded in 2001
- Available in 285 languages
- The English Language Wikipedia is the largest
- The official song of Wikipedia is "Hotel Wikipedia"



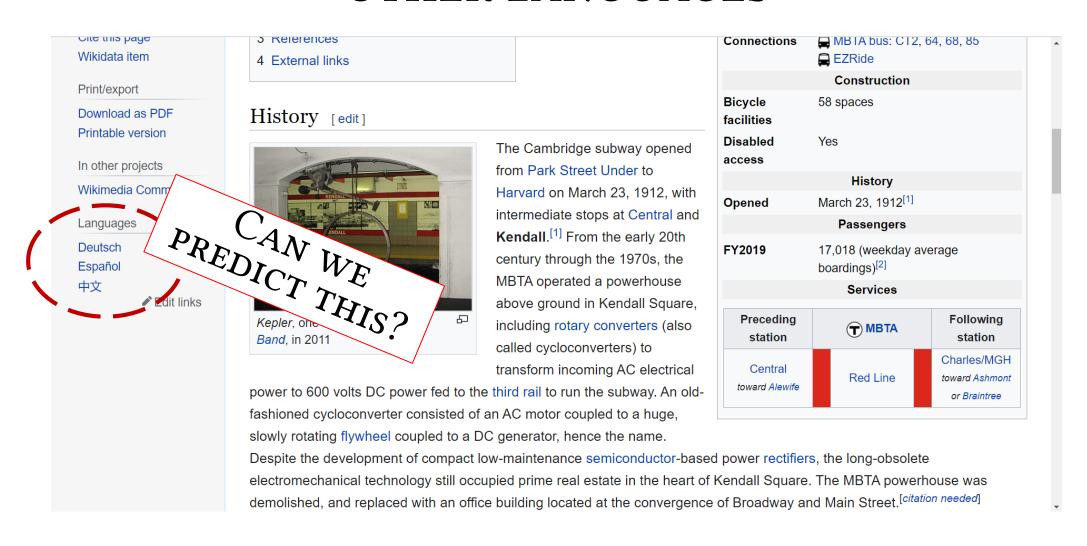
ARTICLES ARE OFTEN TRANSLATED/WRITTEN IN OTHER LANGUAGES



ARTICLES ARE OFTEN TRANSLATED/WRITTEN IN OTHER LANGUAGES



ARTICLES ARE OFTEN TRANSLATED/WRITTEN IN OTHER LANGUAGES

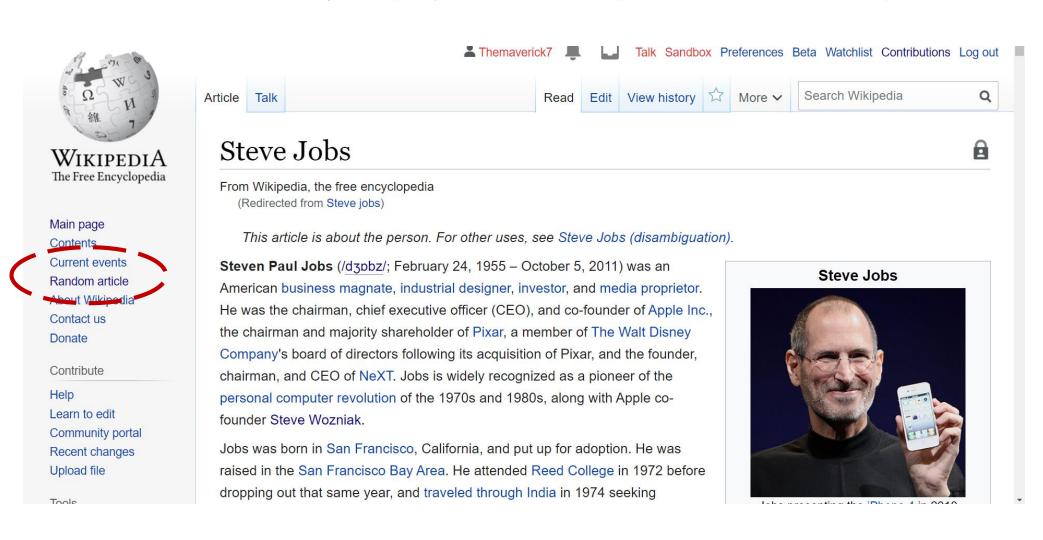


OUTLINE

Contents [hide]

- 1 Introduction
- 2 Models
 - 2.1 Manually-curated Features
 - 2.2 Natural Language Processing (NLP)
 - 2.3 Combined Approach
- 3 Conclusions and Future Directions

3655 ARTICLES WERE SCRAPED USING THE 'RANDOM ARTICLE' FEATURE



SELECTED FEATURES FOR REGRESSION

- Body text length
- Days since article creation
- # of references

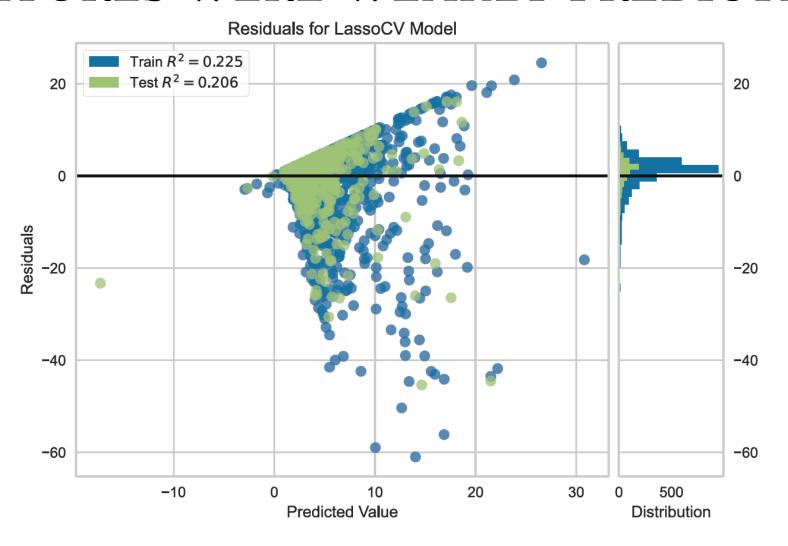


SELECTED FEATURES FOR REGRESSION

- Body text length
- Days since article creation
- # of references
- Total length of text in tables
- # of "Main article" links
- # of "See also" links
- # of "Further information" links
- # of links to other wiki articles
- Is sports?



REGRESSION USING MANUALLY-CURATED FEATURES WERE WEAKLY PREDICTIVE



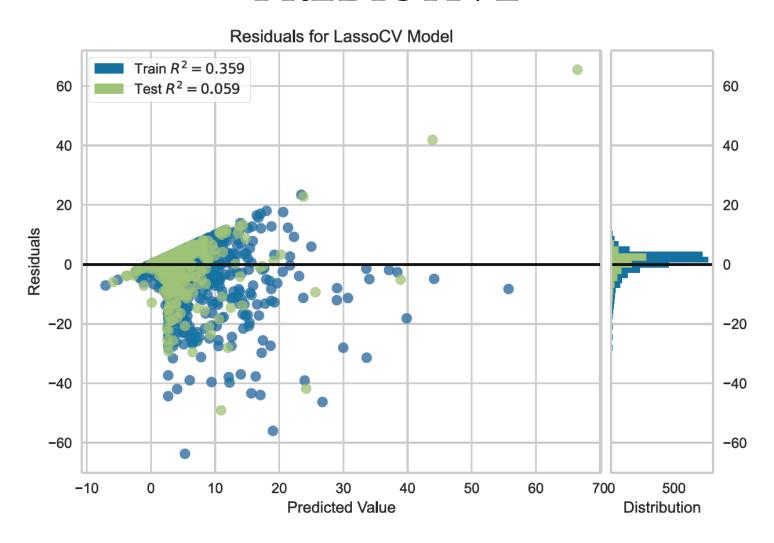
NLP METHOD: BAG OF WORDS

This is a article about the future of our country. It looks quite grim. However, we believe we can make it through all the hardships like we have before. Our country is strong. The author of this article approves this message.

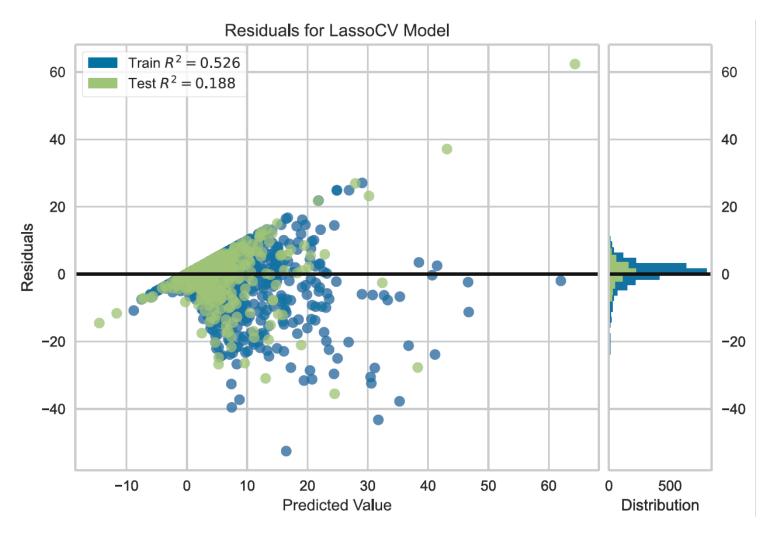


'article'	'future'	'country'	'comfort'	'author'	'message'
2	1	2	0	1	1

REGRESSION USING BASIC NLP WAS LESS PREDICTIVE



Combination of the two approaches did not yield an noticeable improvement



OUTLINE

Contents [hide]

- 1 Introduction
- 2 Models
 - 2.1 Manually-curated Features
 - 2.2 Natural Language Processing (NLP)
 - 2.3 Combined Approach
- 3 Conclusions and Future Directions

CONCLUSIONS AND FUTURE DIRECTIONS

- # of languages can be weakly predicted
- The most under-translated articles tend to be regional articles in the United States, often featuring sports
- Future improvements:
 - More refined NLP methods, or neural networks
 - Try incorporating PageRank
 - Frequency of editing per year





THANK YOU WIKIPEDIA!

APPENDIX: LIST OF COEFFICIENTS FOR THE MOST IMPACTFUL WORDS

```
population
            1.212857
            0.814047
   county
            0.779921
     roman
            0.703669
      also
            0.666219
       . . .
   united -0.249402
     well -0.271809
      show -0.279516
      date -0.377511
 economic -0.519919
```

Appendix: Histogram of Language distributions

