

Topic Modeling and Sentiment Analysis on News Articles and Comments

SFU

WORK
INTEGRATED
LEARNING

Simon Fraser University, Big Data Programming II
Chithra Bhat, Ruoting Liang, Tianpei Shen

SOCC: The SFU Opinion and Comments Corpus



Opinion articles and the comments published in the Canadian newspaper *The Globe and Mail* in the five year period between 2012 and 2016.

- 10,339 articles + 1,280,454 comments
- 6,895,696 words(articles) + 77,238,179 words (comments)

Motivation

Create a succinct summary of opinions on articles, issues, or policies among Canadian citizens to help organizations make better decisions.

- ❖ Observe the [change in hot topic](#) for these 5 years (2012 - 2016).
- ❖ Discover the [surrounding topics](#) on comments under a given article.
- ❖ Discover the [constructiveness/toxicity/sentiment](#) for the comments under a given topic.

Case Study:

IAN BURUMA What drives anti-immigrant sentiment?

IAN BURUMA Special to The Globe and Mail Published Saturday, Jan. 03, 2015 3:00AM EST Last updated Saturday, Jan. 03, 2015 3:00AM EST

Comments

LOL! Good one.

Mass immigration is an environmental and social rolling disaster.

He should go to hell!

An impartial analysis would do much to dispell the misinformation and propaganda as well as establish the facts...

These are folks whose costs of education and healthcare while they were growing up were borne by other countries but whose brains were drained to Canada for the benefit of our country. Sometimes we should ...

Topic: anti-immigration

Positive

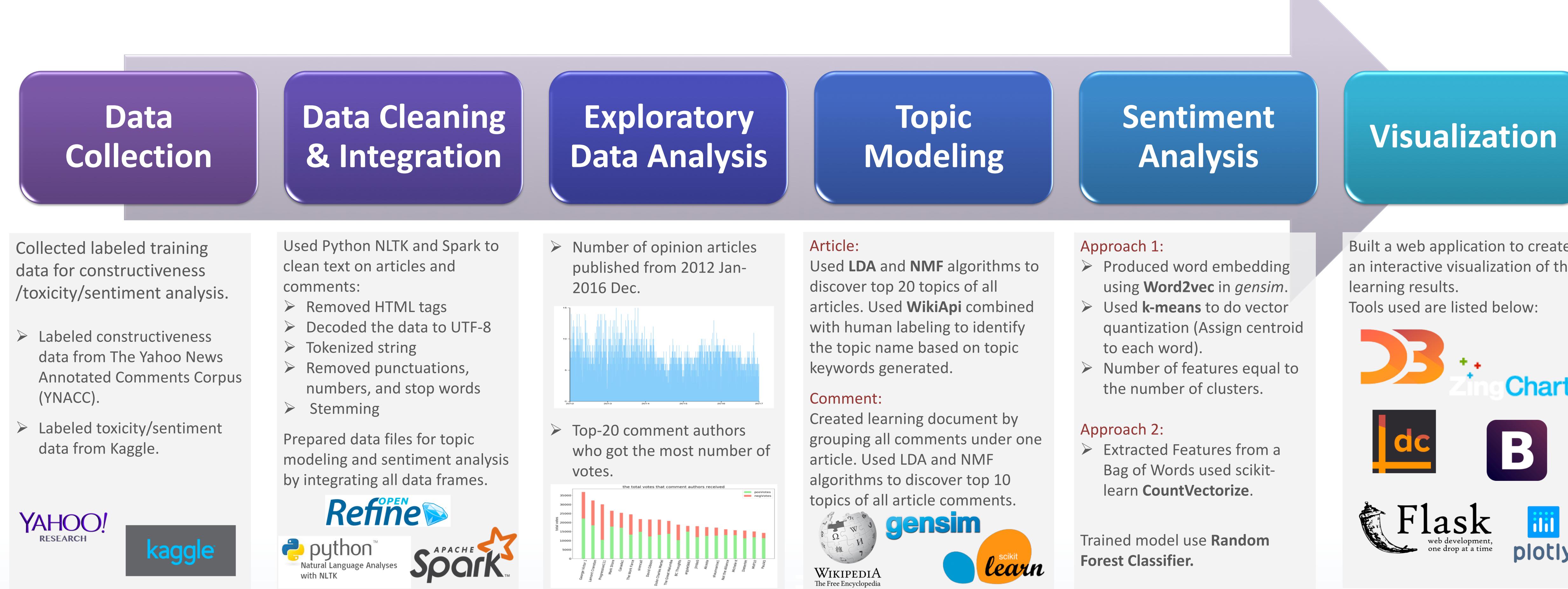
Negative

Toxic

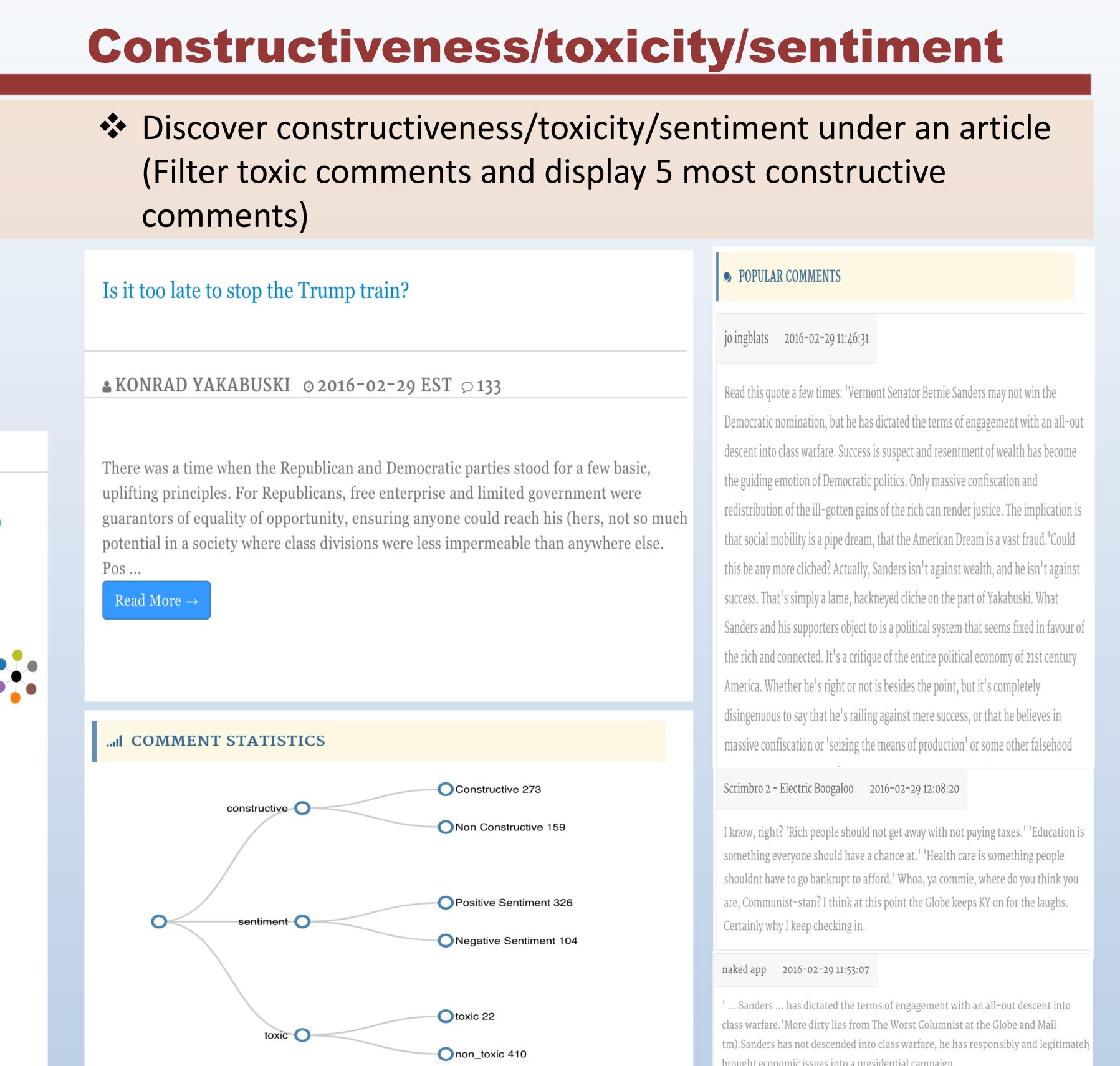
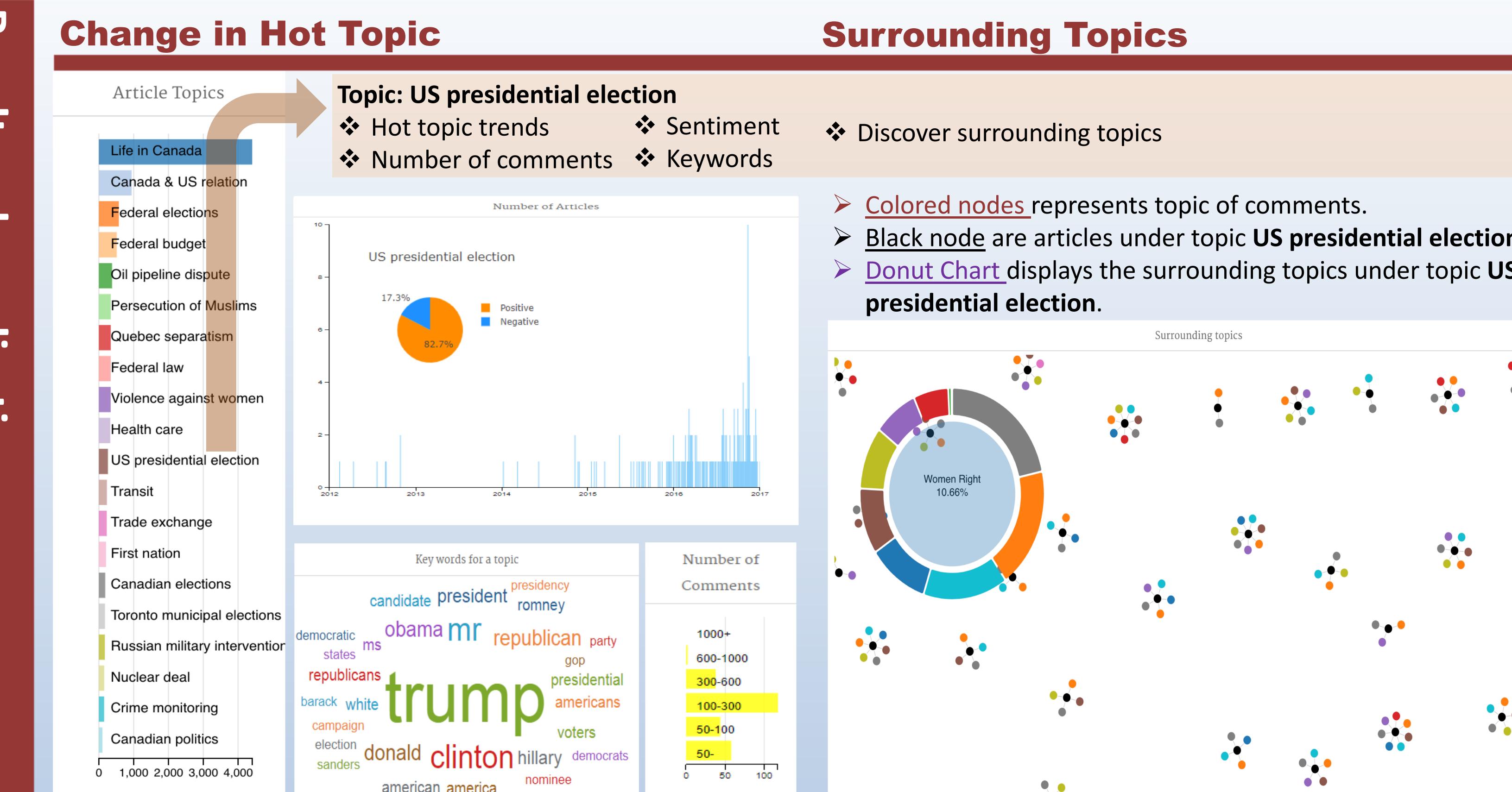
Constructive

Education, Economic, Society

Data-Science Pipeline



Result-web application



Discussion

Stemming is not always good: The algorithm results in a list of word-vector, which is not human-readable. When performing the human labeling we noticed that stemming sometimes will give us strange word such as 'parti', 'conser'. Therefore we ignore the stemming step during cleaning data.

How to deal with large dataset: The comment dataset is 350MB. So we used Spark to process the data and created User Defined Function for code implementation.

How to do topic modeling on short-text: LDA is one of the most popular topic models and it models a document as a mixture of topics. However, it is not the case for short text, such as comments. Therefore we grouped all comments of one article and treated them as one document in order to perform the topic modeling.

The dataset doesn't contain label for constructiveness and toxicity: In order to analyze the constructiveness and toxicity of comments, we need to train the model with a labeled dataset. But the original data doesn't contain these labels, therefore, we use other sources to train the model and test on the original dataset.

Future Work

Crowdsourcing on article comments

As discussed above, the SOCC datasets do not contain labeled data. In order to build training dataset that better represent SOCC comment data, one approach is to use crowdsourcing and get labeled data with constructiveness, toxicity, and sentiment.

Make real time prediction

In the future, we can add more feature to our web application. For example, predict a topic in real time by giving an article URL or article text.

Reference:

- C., J., E., B., & A. (2017). *Finding Good Conversations Online: The Yahoo News Annotated Comments Corpus*. Valencia, Spain: Association for Computational Linguistics. Page 13-23
- Blei, D.M., Ng, A.Y., Jordan, M.I., 2003. Latent dirichlet allocation. J. Machine Learn.
- Cheng, X., Yan, X., Lan, Y., & Guo, J. (2014). BTM: Topic Modeling over Short Texts. IEEE Transactions On Knowledge And Data Engineering, 26(12), 2928-2941.