

# Monitoring Social Media using Machine Learning

Summer 2019

Joseph Jinn, Professor Keith VanderLinden

## Introduction

Our research is an extension of prior work by CSIRO - Commonwealth Scientific and Industrial Research Organization, Australia's national research laboratory. Our focus is on utilizing Twitter data, Tweets, as a dataset by which we measure the SLO - Social License to Operate - of various mining, gas, and oil companies. SLO is defined as the acceptability of a company's business operations by its employees, stakeholders, and the general public. The primary purpose of the summer 2019 research project is to investigate and find a methodology by which we can effectively model the topics of all the Tweets in our dataset. Topic modeling is a way of defining abstract "topics" that are prevalent in a corpus of textual documents. It is statistical in nature and is essentially unsupervised machine learning by which we attempt to cluster the Twitter data to find similarities and patterns among groups of words.

To that end, we first utilized standard data science techniques to investigate the nature of our Twitter dataset. This involves the use of the Python programming language, the Pandas data analysis library, the Matplotlib data visualization library, and other processing and visualization software. Our discoveries and results are recorded in Jupyter Notebooks – an interactive web-based application that allows researchers to easily share code, equations, visualizations, and text. We also utilize Scikit-Learn, a machine learning software suite, and Gensim, a topic modeling software suite, along with various 3rd party libraries, to implement baseline topic models from which we can begin to investigate how to best extract relevant topics from the Tweet texts.

## Objectives

1. Analyze the Twitter dataset to derive various numerical and categorical statistics in order to determine the nature of our data.
2. Preprocess, postprocess, and tokenize Tweet text using the spaCy Natural Language Processing library.
3. Evaluate baseline topic modeling algorithms using Gensim, Scikit-Learn, and other 3rd party Python packages.
4. Visualize the topic extraction results using pyldAVIS, matplotlib, and other 3rd party Python packages.
5. Understand the statistical and general mathematical construct behind each algorithm to infer the best approach for our data.
6. Create a modified baseline topic modeling algorithm that generates more coherent topics and associated words.

## Results - Topic Extraction

### Author-Topic

Label: 1 -> Words: leard go look maules work think well get have property  
Label: 2 -> Words: project road beach news win rd downer ceo 's peter  
Label: 3 -> Words: gas coal farmer seam water stop community people want protest  
Label: 4 -> Words: creek write 's wa coal | project company timor plan  
Label: 5 -> Words: gas coal forest petroleum tycoon energy price oil sale whc  
Label: 6 -> Words: tax joyce barnaby slo\_cashil pay chevron liberal rail coal origin  
Label: 7 -> Words: narrabri water gas risk basin coal national artesian > farmer  
Label: 8 -> Words: i ° ½² beach field ¼¼□□ tour \$ win home  
Label: 9 -> Words: eis land coal water inland go seam want appliance fracking  
Label: 10 -> Words: \$ wpl sto video sire share field day gain result

### Non-Negative Matrix Factorization

Topic 0: coal build away indian massive environment india mean minister power  
Topic 1: tell say tumbull way ask make question happen try issue  
Topic 2: job 000 10 create claim lie thousand tourism renewable pm  
Topic 3: australian company oil face govt financial use fail write set  
Topic 4: project gas narrabri land forest farmer seam sign community field  
Topic 5: stop tax pay profit million corporate haven office rate island  
Topic 6: time thank good come late stand start leave week long  
Topic 7: reef barrier fight help kill destroy risk protect let save  
Topic 8: labor support green vote lnp win election shorten alp party  
Topic 9: ¼i beach ¼i read love ½² slo\_hash letter oh ¾i  
Topic 10: want people know billion future dollar listen industry mega video  
Topic 11: day action protest court group right join wrong native meet

## Processes

### Twitter API

- Tweet Attributes
- JSON file format
- UTF-8 Encoding

### Raw Tweet Dataset

- Determine company
- Determine language
- Drop no-Company Tweets
- Drop non-English Tweets
- CSV file format

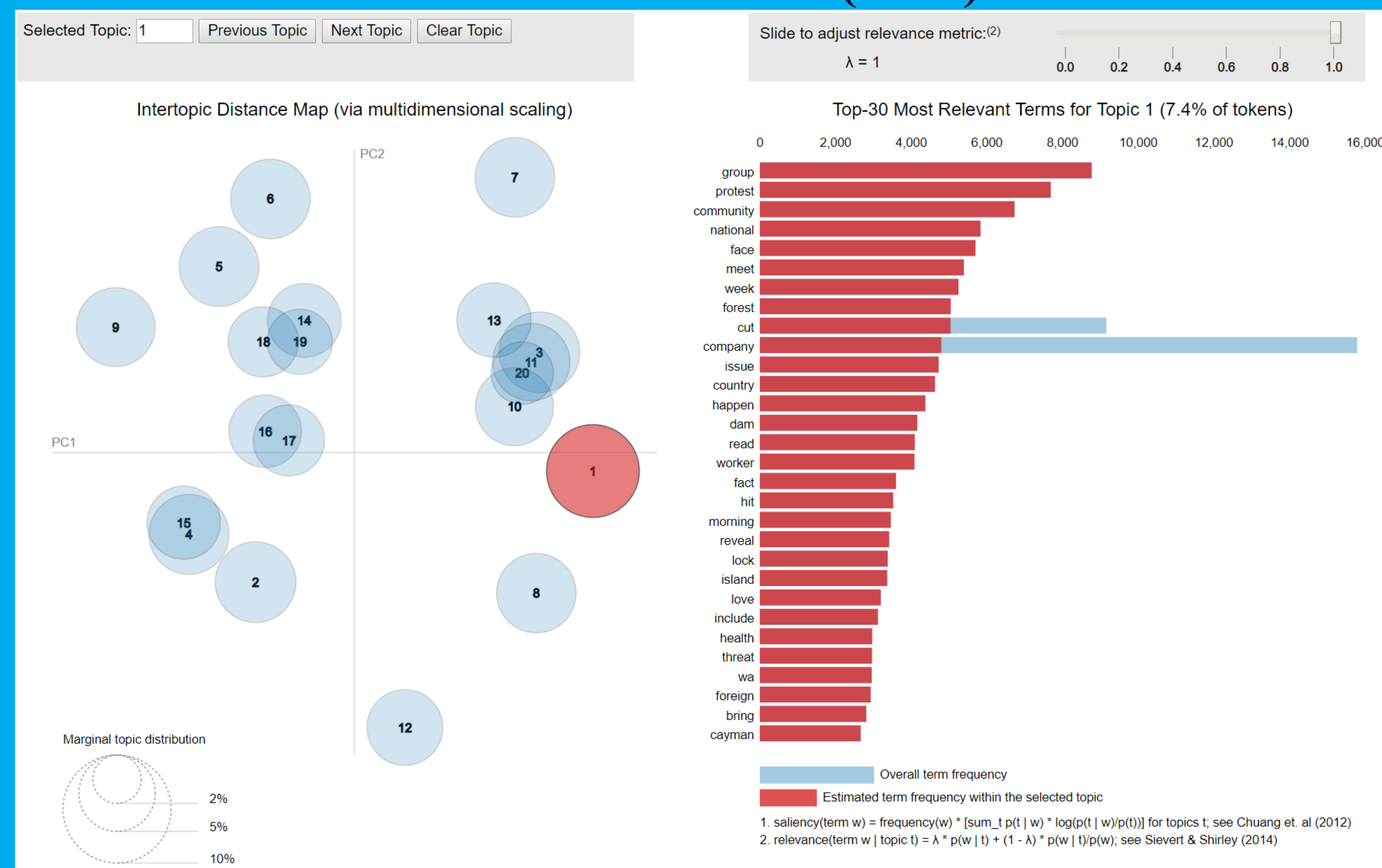
### Dataset

- Preprocessing
- Postprocessing
- Tokenization
- Topic Extraction

### Latent Dirichlet Allocation

Topic 0: work court risk world group port native live title coal  
Topic 1: want good coalmine claim profit finance clear gov voter poll  
Topic 2: gas time coal know oppose thank industry narrabri think ask  
Topic 3: project fund reef climate coal change wo barrier kill watch  
Topic 4: fortescue protest cost local pm forest fine coal financial hand  
Topic 5: support fight right join wrong campaign demand planet country happen  
Topic 6: need farmer deal cut rule coal premier run royalty investment  
Topic 7: water australian future land question open coal pollution morning concern  
Topic 8: labor coal rail energy slo\_cashn line win shorten election stop  
Topic 9: stop billion public action destroy dollar naif sign national party  
Topic 10: plan say business big council use end law year production  
Topic 11: new loan govt coal bank price giant make support long  
Topic 12: slo\_mention tell help leave start environment talk massive away turn  
Topic 13: job government tumbull lnp look 000 year way lose create  
Topic 14: people queensland power let coal alp community thing stand close  
Topic 15: coal ½i build ½² ¼i galilee ahead basin approve china  
Topic 16: company state iron indian ore news high lie face carbon  
Topic 17: tax pay break approval federal ceo subsidy issue slo\_cash miner  
Topic 18: india money report come taxpayer point day barnaby joyce week  
Topic 19: green vote minister share sell canavan matt promise buy read

## Data Visualization (LDA)



## Conclusion

Baseline topic modeling algorithm libraries do not work well on our dataset. Latent Dirichlet Allocation and all derivatives of this algorithm work best on corpora containing documents that are long and written in a formal grammatical style. Tweets suffer from limited character length, grammatical inconsistency, and Twitter-specific linguistic elements, which makes it difficult to extract coherent topics.

Hyperparameter tuning may improve results to some extent but would be a time-consuming and exhaustive process. Biterm execution runs take almost half a day per. Hierarchical LDA suffers from RAM overflow issues due to its recursive nature. Utilization of Calvin College's Borg Supercomputer could expedite matters but would require parallelization of our codebase and the construction of a Singularity container.

## References

- David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. J. Mach. Learn. Res. 3 (March 2003), 993-1022.
- David M. Blei, Michael I. Jordan, Thomas L. Griffiths, and Joshua B. Tenenbaum. 2003. Hierarchical topic models and the nested chinese restaurant process. In Proceedings of the 16th International Conference on Neural Information Processing Systems (NIPS'03), S. Thrun, L. K. Saul, and B. Schölkopf (Eds.). MIT Press, Cambridge, MA, USA, 17-24.
- Michal Rosen-Zvi, Thomas Griffiths, Mark Steyvers, and Padhraic Smyth. 2004. The author-topic model for authors and documents. In Proceedings of the 20th conference on Uncertainty in artificial intelligence (UAI '04). AUAI Press, Arlington, Virginia, United States, 487-494.
- R. Zhao and V. Y. F. Tan, "Online Nonnegative Matrix Factorization With Outliers," in IEEE Transactions on Signal Processing, vol. 65, no. 3, pp. 555-570, 1 Feb.1, 2017.  
doi: 10.1109/TSP.2016.2620967
- Xiaohui Yan, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng. 2013. A biterm topic model for short texts. In Proceedings of the 22nd international conference on World Wide Web (WWW '13). ACM, New York, NY, USA, 1445-1456. DOI: <https://doi.org/10.1145/2488388.2488514>
- Yee Whye Teh, Michael I Jordan, Matthew J Beal & David M Blei (2006) Hierarchical Dirichlet Processes, Journal of the American Statistical Association, 101:476, 1566-1581, DOI: 10.1198/016214506000000302

## Future Work

Our plans are to continue Objectives 5 and 6. We do not have a full grasp of the statistical and mathematical construct behind each topic modeling algorithm. This understanding will be essential to any attempt to create a modified baseline topic modeling algorithm that will hopefully improve topic extraction results on our Twitter dataset. It is our hope that we can minimize model perplexity while maximizing topic coherence metric values.

We also hope to obtain an updated Twitter dataset with more recent Tweets from CSIRO.