



BACKGROUND

Cryptocurrency is a concept of digital currency in which encryption techniques are used to regular the generation and verify the transfer of funds based on blockchain, a distributed ledger that is inherently resistant to modification of the data.

Ever Since its birth, its nature and value has been highly debated. It is the ideal digital asset in the world of internet, that is decentralized and inherently resistant to modification of the ownership. Even though the nature and value has been highly debated, the combination of security and transparency makes it one of the most important innovation in the era of 'cloud data', where security is the last shield of privacy.

DATA SUMMARY

Data

- Bitcoin Price and Volume
 - Kraken, Coinbase, Bitstamp and Itbit; 2012/09 to present
- News Archive
 - 2017 Wall Street Journal, 43,268 valid entries
- Corpus
 - NLTK corpus
 - Wikipedia, 669,539 vocabulary with 35,556,952 documents

METHODS

Method

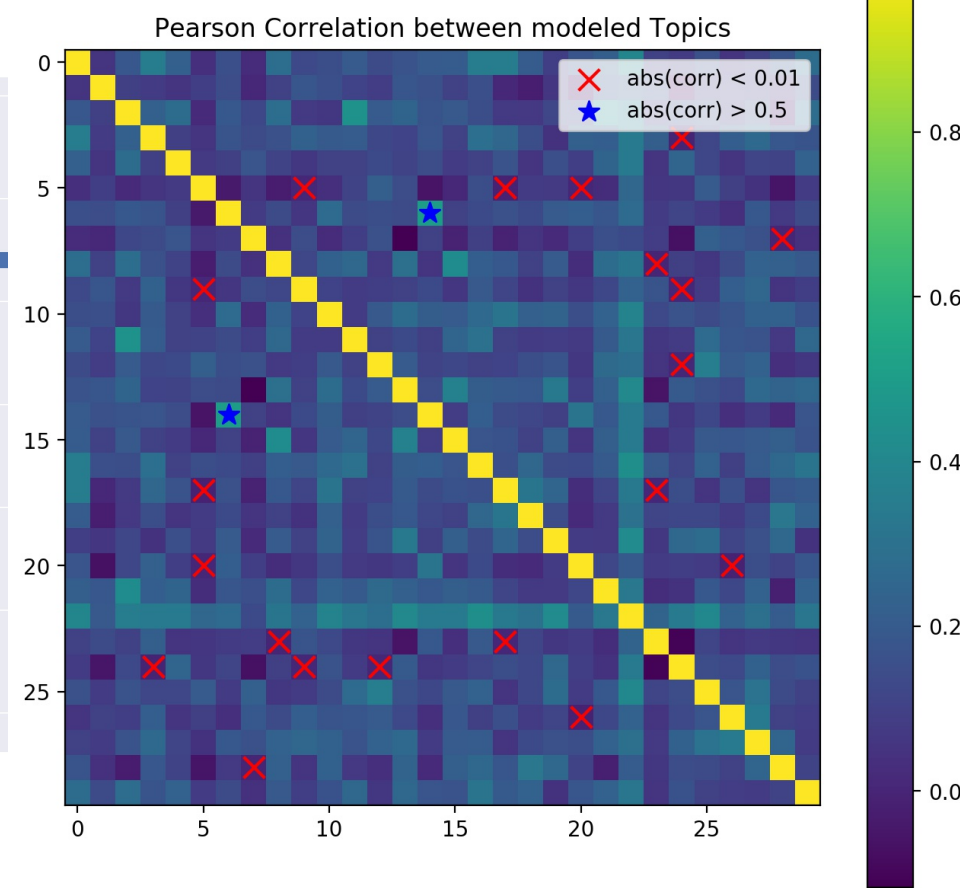
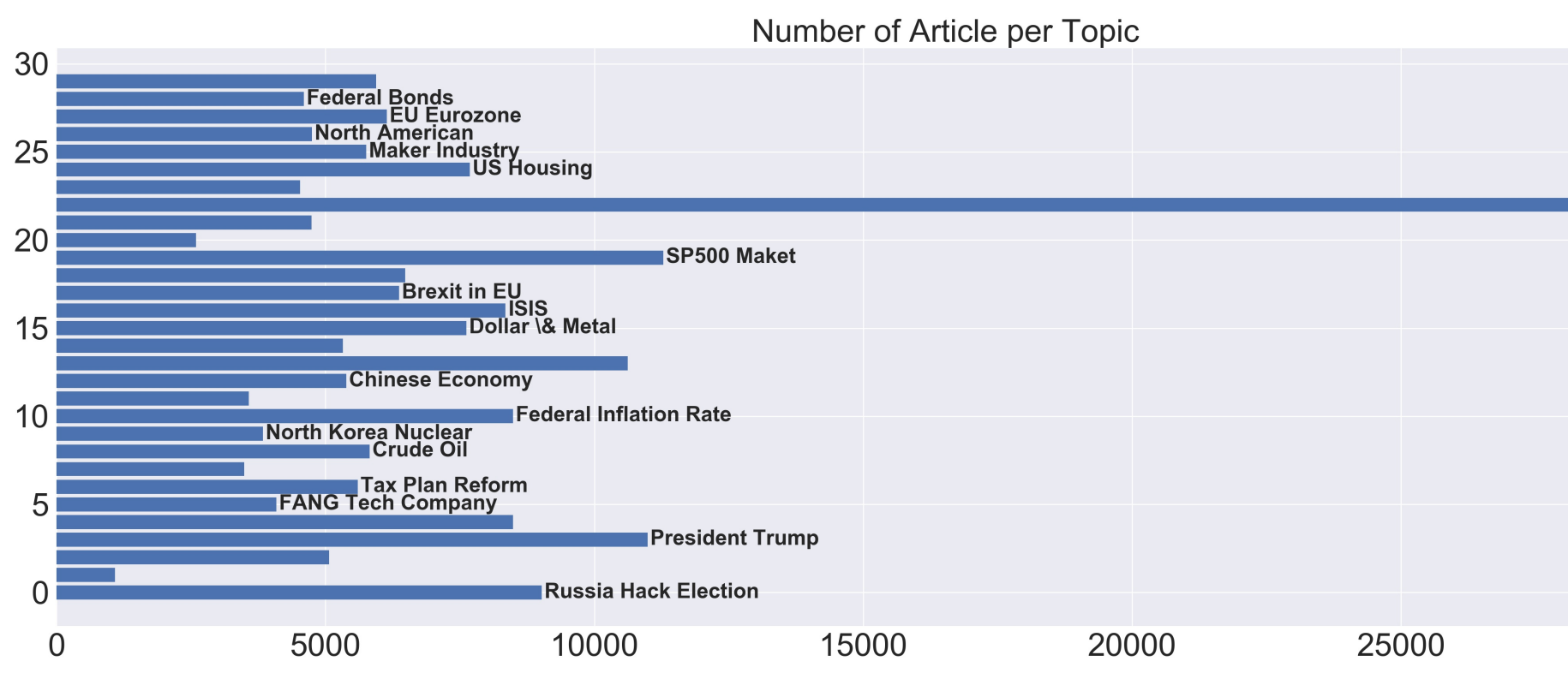
- Sentiment
 - Textblob: pre-trained on sentimental corpus
- Topic model
 - TF-IDF: term frequency and inverse document frequency
 - NMF: non-negative matrix factorization
- Semantic model
 - Doc2Vec: two layer neural network with vector representation

$$tf(t, d) = 0.5 + 0.5 \cdot \frac{f_{t,d}}{\max\{f_{t',d} : t' \in D\}}$$

$$idf(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

$$tfidf(t, d, D) = tf(t, d) \cdot idf(t, D)$$

Topic #0	Topic #1	Topic #2	Topic #3	Topic #4	Topic #5	Topic #6	Topic #7	Topic #8	Topic #9	Topic #10	Topic #11	Topic #12	Topic #13	Topic #14	Topic #15	Topic #16	Topic #17	Topic #18	Topic #19	Topic #20	Topic #21	Topic #22	Topic #23	Topic #24	Topic #25	Topic #26	Topic #27	Topic #28	Topic #29		
election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election	election



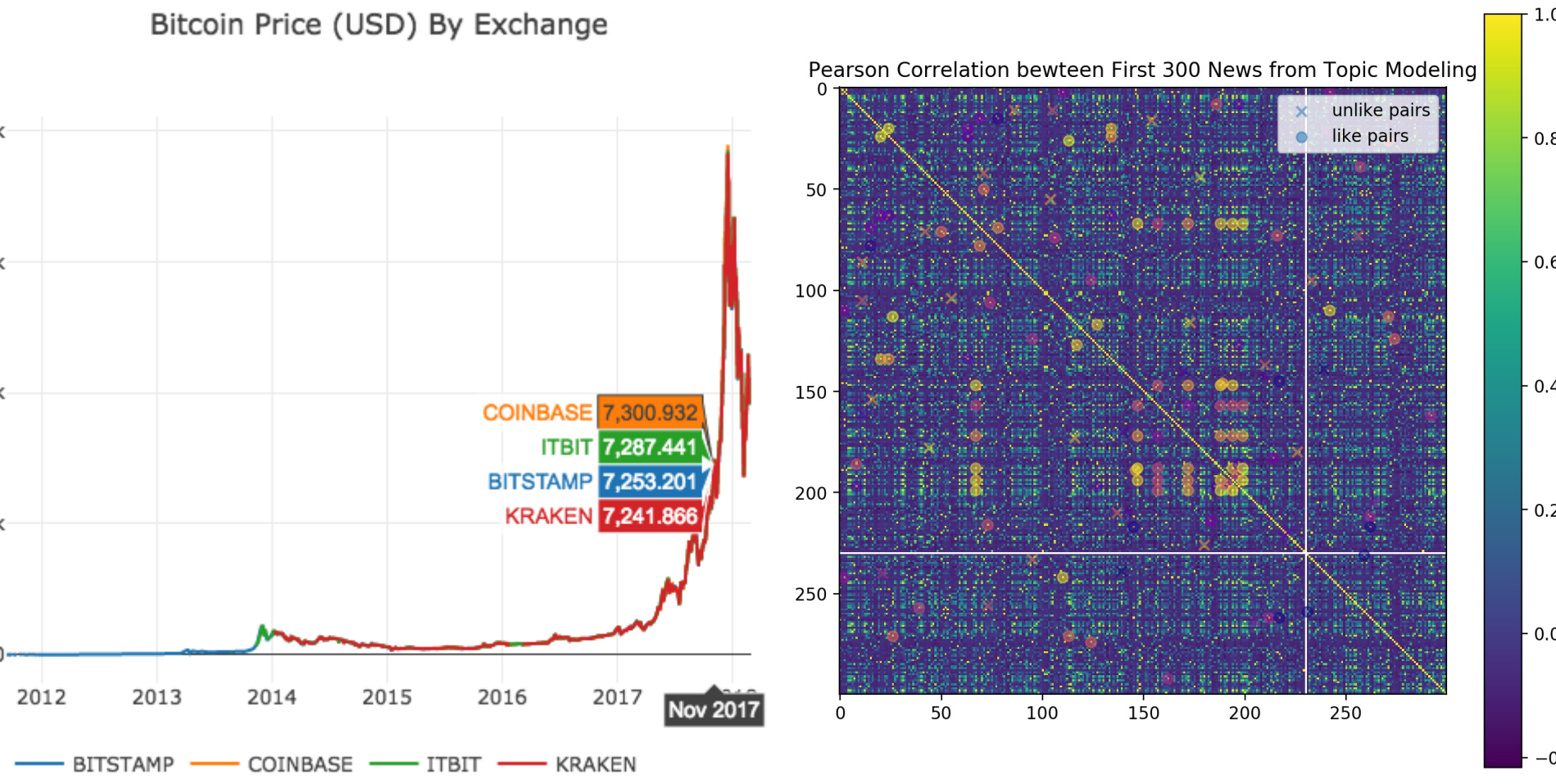
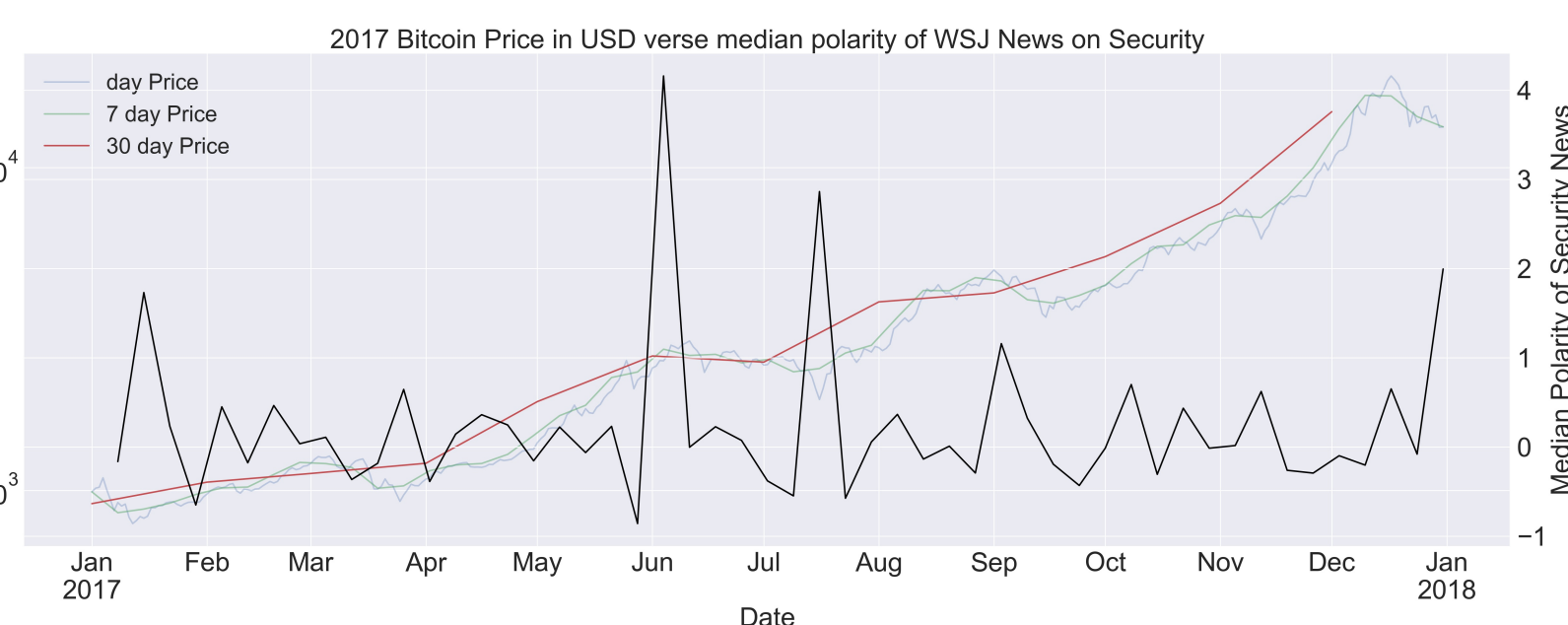
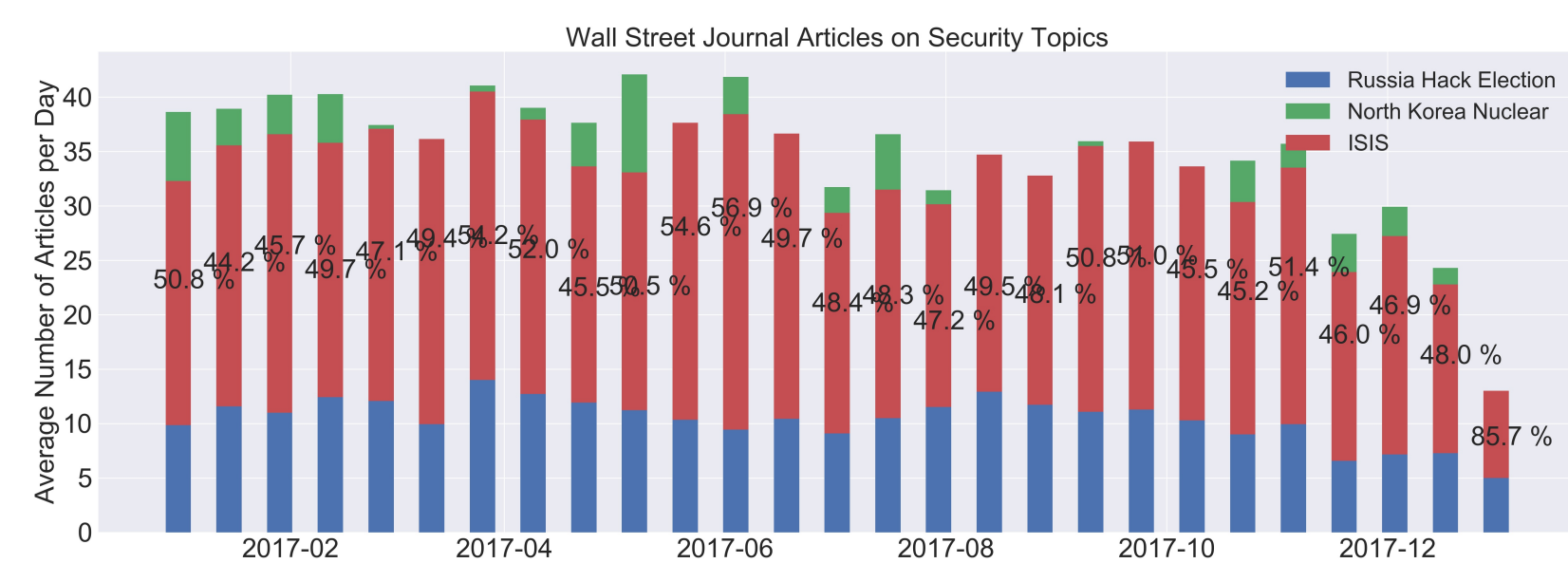
Topic Modeling:

- The popular Brown corpus from NLTK is used to exclude the stop words, in addition to the common stop words in English.
- The rest of the words are built into a vocabulary and all of the 2017 WSJ news are then factorized based on the vocabulary.
- TF-IDF is adopted to generate count statistics. Top 20 weighted key words are saved to summarize each topics.
- We use NMF to decompose news into 30 dimensions and convert the key words back. See the top 10 of each topic in table.

- A few examples on topic #9: North Korea Nuclear
- 2017 South Korea: Reacts to U.S. Missile Defense System Hunt 1 Shifted A Top South Korean national newspaper official said in a top news column about the government's commitment to a controversial U.S. missile defense system.
 - 2017 U.S. Plans Missile Defense Test Ahead Concerns Over North Korea The Pentagon is planning a missile test next week of a system designed to shoot down intercontinental ballistic missiles, U.S. officials said.
 - 2017 North Korea Nuclear: Threat From Nuclear 1 Election Focus on South Korea: Japan says it has found signs in escalating tensions with Pyongyang over its nuclear program that have changed the dynamics of...

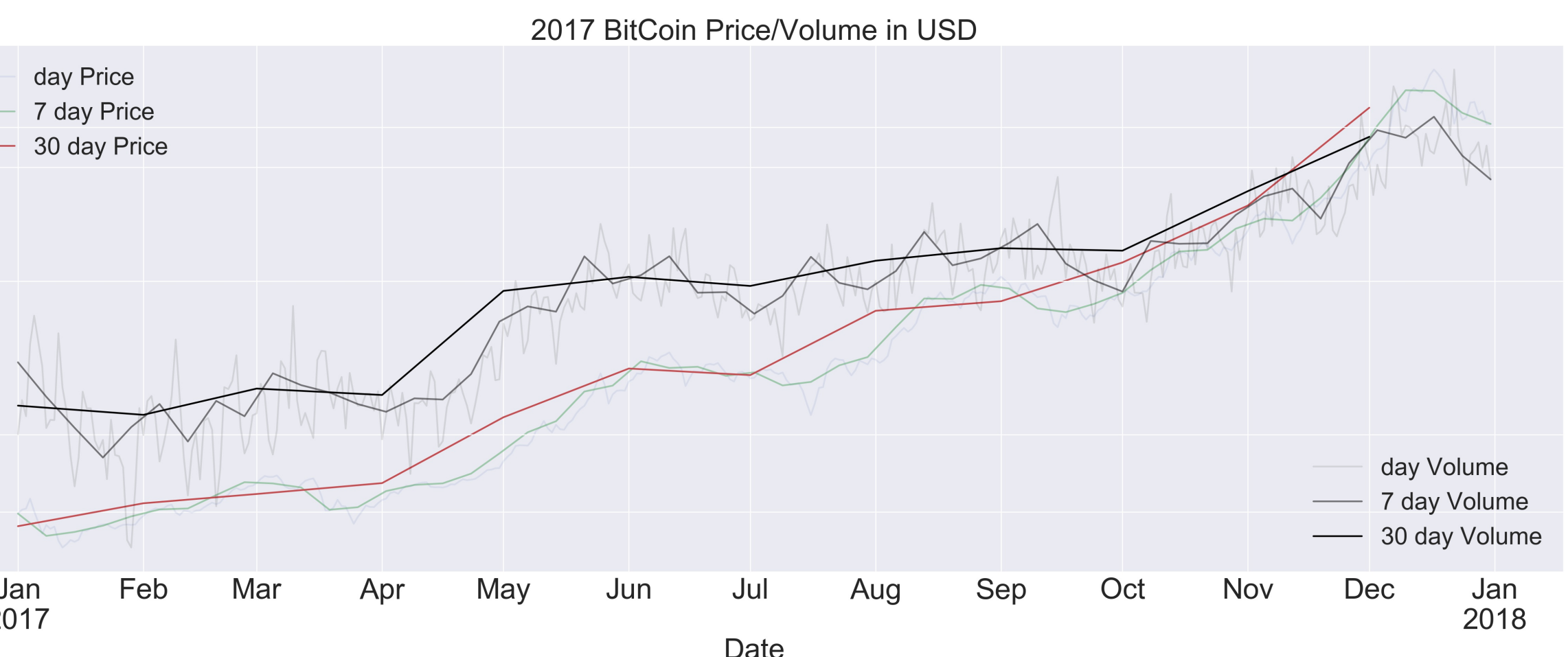
The news are grouped by the topics and then aggregated to each day of the year. The average polarity is taken to indicate the sentiment of the day. A few examples are listed:

- 2016: South Korea: Reacts to U.S. Missile Defense System Hunt 1 Shifted A Top South Korean national newspaper official said in a top news column about the government's commitment to a controversial U.S. missile defense system.
- 2017: U.S. Plans Missile Defense Test Ahead Concerns Over North Korea The Pentagon is planning a missile test next week of a system designed to shoot down intercontinental ballistic missiles, U.S. officials said.
- 2017: North Korea Nuclear: Threat From Nuclear 1 Election Focus on South Korea: Japan says it has found signs in escalating tensions with Pyongyang over its nuclear program that have changed the dynamics of...

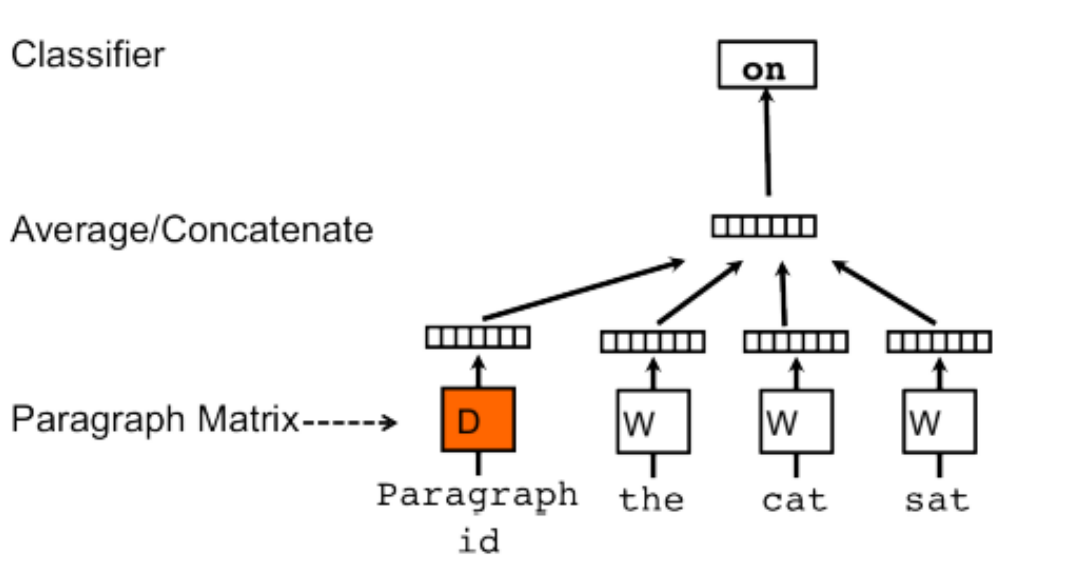


The gap in the red curve is due to the maintenance of Kraken exchange. Due to various reason, the Bitcoin data recorded from four exchanges contains faulty and missing information. Thus, the average Price and Volume are calculated based on the best available and continuous data. The changes in 2017 are primary used to be consistent with News obtained from WSJ archive.

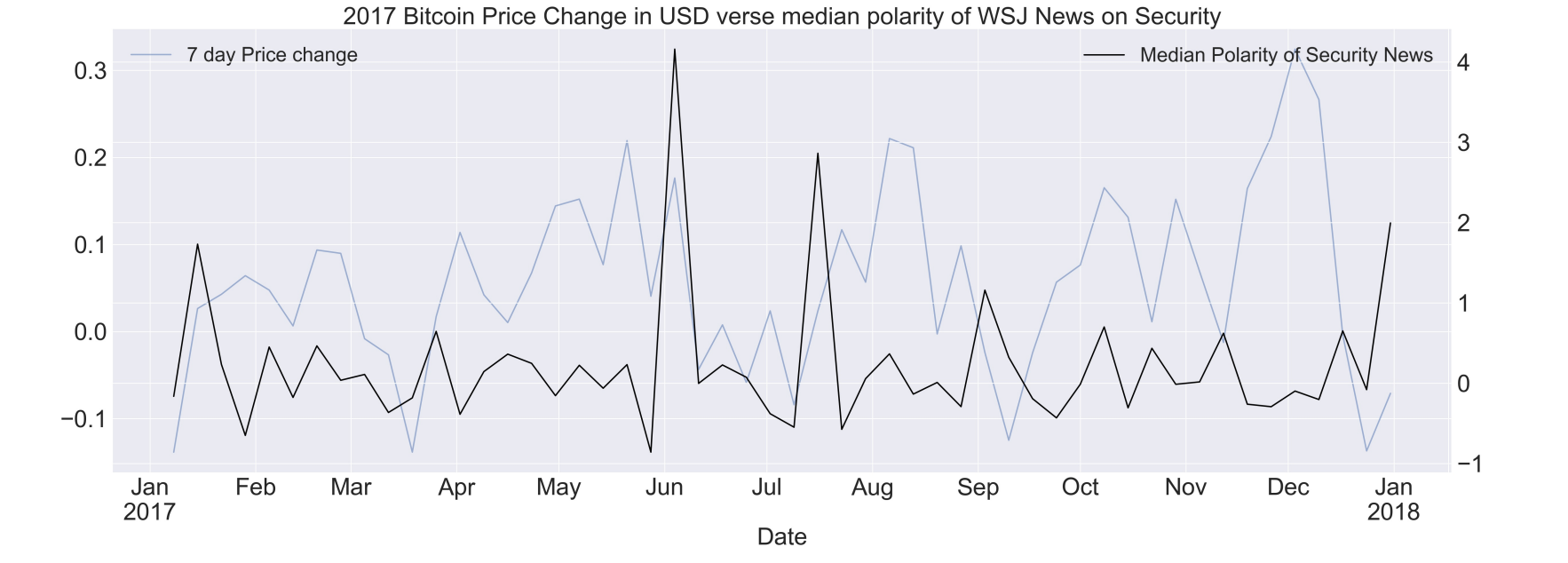
Exchange	Online Days	2017
BITSTAMP	2339	902
COINBASE	902	1538
ITBIT	1538	1465
KRAKEN	1465	0



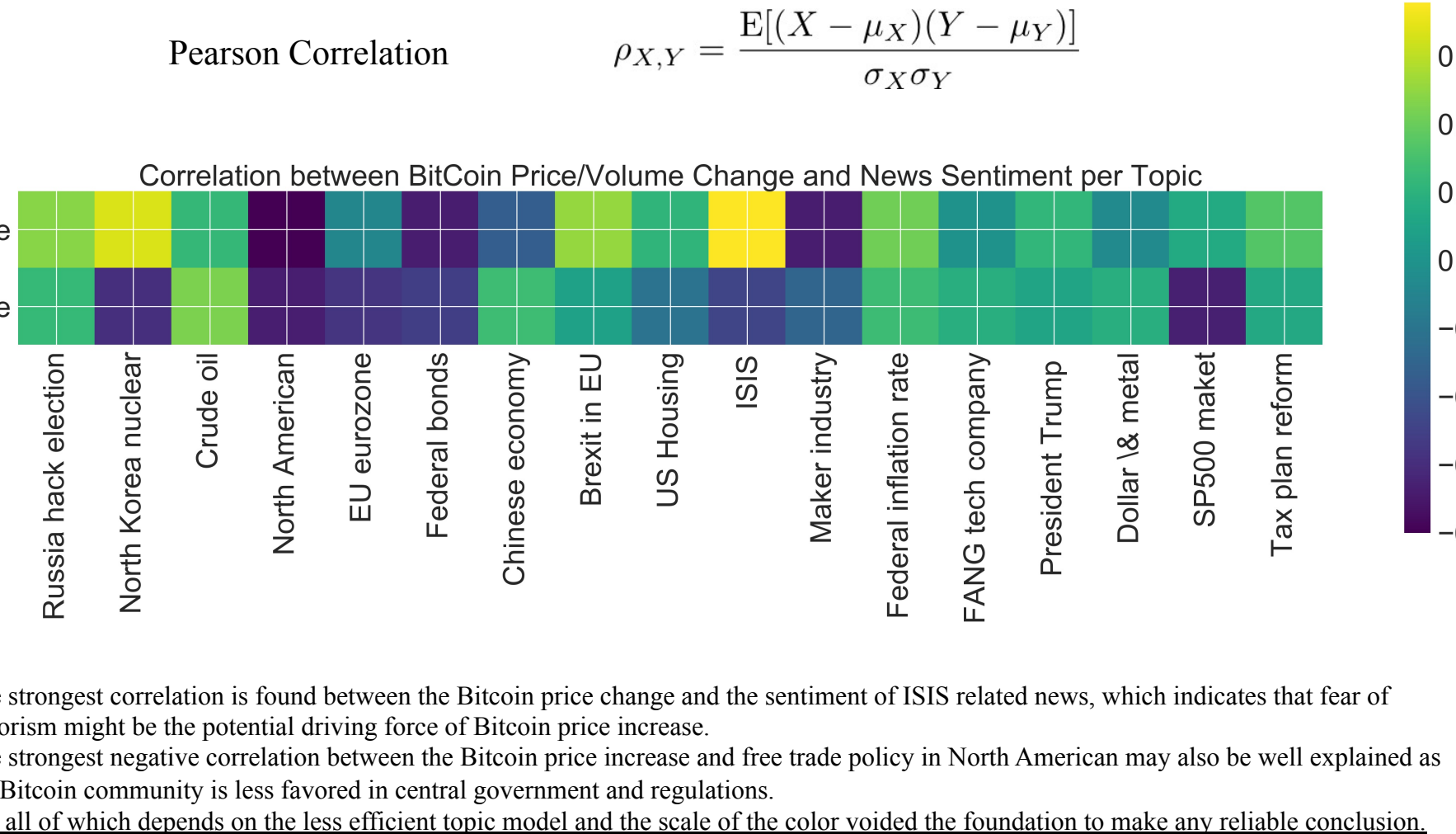
- ### Doc2Vec
- Representation
 - Vector dimension: 300 float
 - Architectures
 - Distributed bag-of-words: window size 15
 - Softmax pulling
 - Negative sampling: 5 noise words will be drawn
 - Learning algorithm
 - Gradient descent



RESULTS



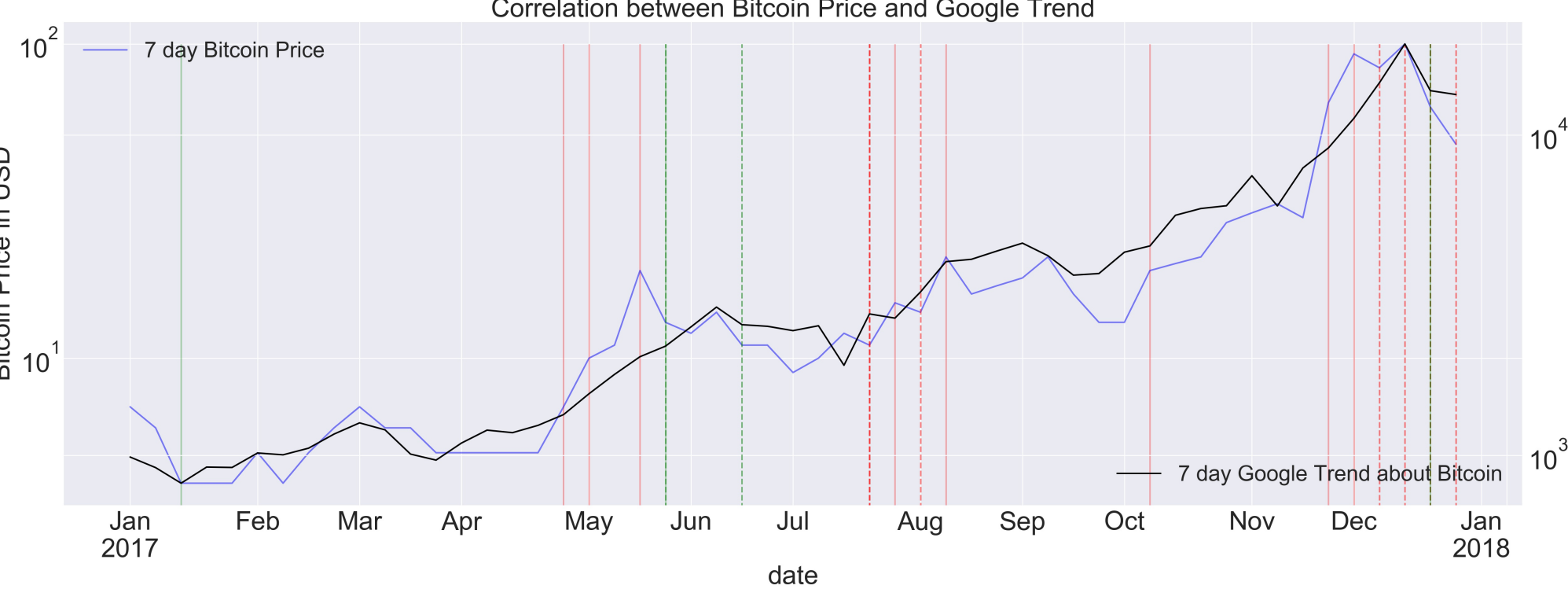
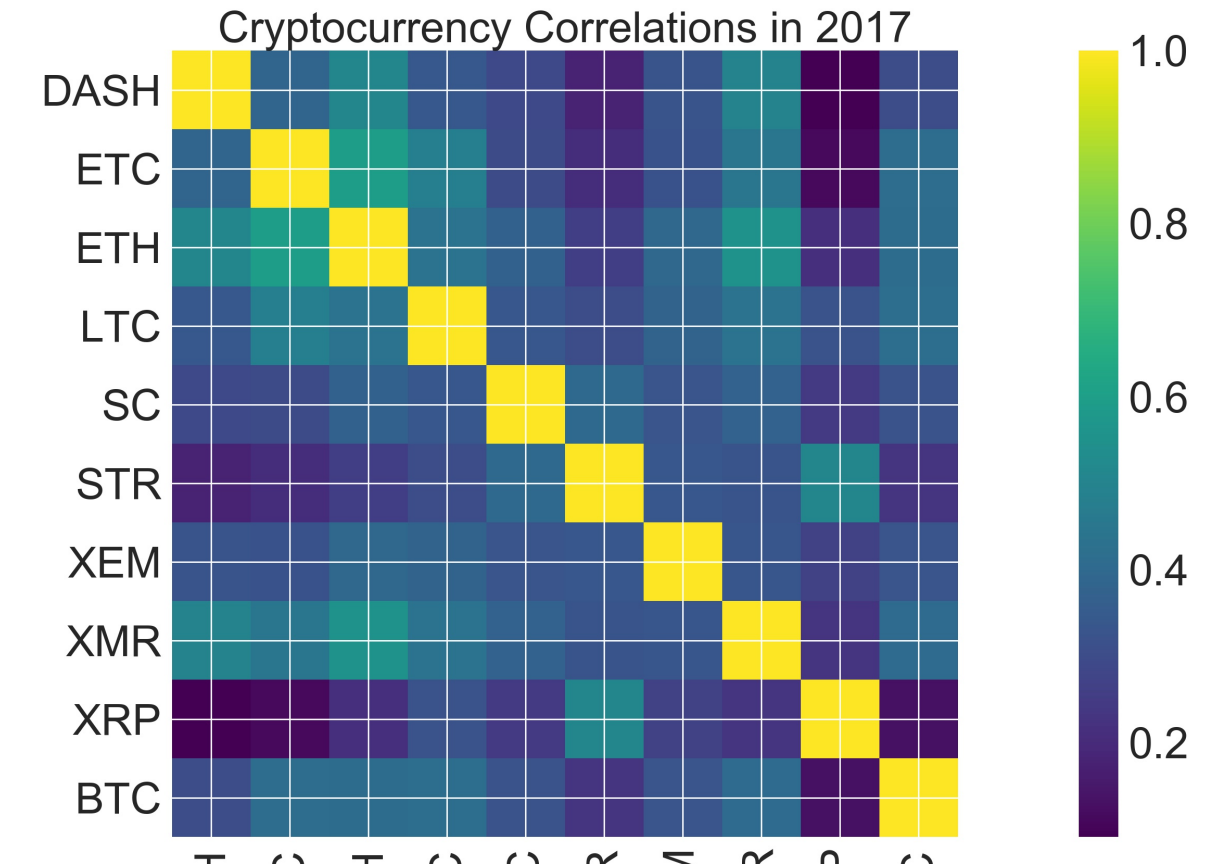
As we can tell from the graph, it is subjective to draw any conclusion about correlation between the 7 day price change and the 7 day average of the daily median polarity for news in security topics. To quantitatively describe the correlation between Bitcoin and WSJ news in 2017, the total sentiment of a news topic is defined to be the median value of the normalized polarity to minimize the bias we find in WSJ news. Then we used Pearson's equation to calculate the correlation coefficient between the Price/Volume change and the daily news sentiment in each topic.



The strongest correlation is found between the Bitcoin price change and the sentiment of ISIS related news, which indicates that fear of terrorism might be the potential driving force of Bitcoin price increase. The strongest negative correlation between the Bitcoin price increase and free trade policy in North American may also be well explained as the Bitcoin community is less favored in central government and regulations. But all of which depends on the less efficient topic model and the scale of the color voided the foundation to make any reliable conclusion.

NEW HOPE

- ### Conclusions
- As we demonstrated, Doc2Vec can model news semantic more precise and accurate than TFIDF.
 - WSJ News are biased on sentiment and topics. Google Trend and Twitter will be better.
 - It is too early and vague to connect Bitcoin with unlawful crime and security exposure.
 - Correlation between different cryptocurrency exists and worth exploring.



Ramage, Daniel, Susan T. Dumais, and Daniel J. Liebling. "Characterizing microblogs with topic models." ICWSM 10, no. 1 (2010): 16. I.e. Quoc, and Tomas Mikolov. "Distributed representations of sentences and documents." In International Conference on Machine Learning, pp. 1188-1196, 2014.