

## Motivation

- Stock movement prediction has long attracted both investors and researchers. It is a challenging problem: the market is highly volatile and it is affected by numerous factors.
- Discover whether and how news articles and social media would affect the stocks market trend.

## Goals

- Use NLP methods to do feature extraction on news and twitter dataset to predict stock price.
- Optimize the deep learning model with time-series dataset to predict the stock price trend or stock price
- Build an interactive platform to acquire latest news article, tweets, and stock price for several tickers and Nasdaq Index and predict with real-time processing.

## Pipeline

### Data Collection and Integration

**Storage:** Json and pkl files

**Data Sources:** New York Times, Bloomberg, CNN, The Washington News, Reuters, The Guardian, BBC

**APIs:** Quandl, Alpha Vintage, NYT, The Guardian, Twitter

**Time Frame:** 2014-2019

Quandl BBC

The Washington Post

REUTERS®

### Data Preprocessing Exploratory Data Analysis

#### Text Preprocessing:

- Tokenize, remove stop words, lemmatize

#### EDA:

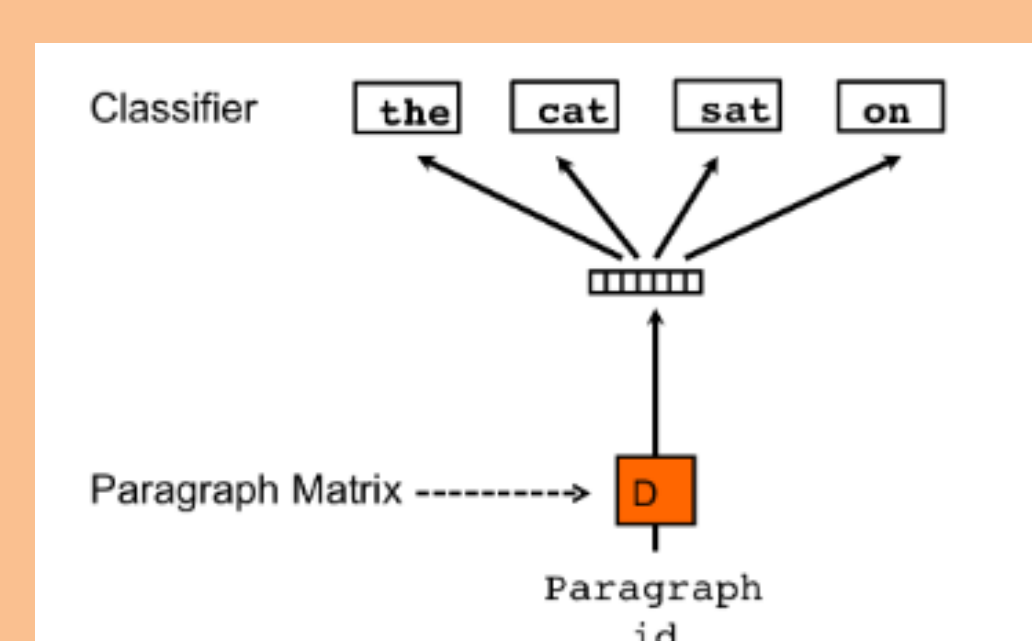
- Plot the stock trend with the positive and negative news count
- Use LDA to do topic modeling on news content
- Find the news distribution in different news categories



### Feature Extraction and Engineering

#### Methods:

- Bag of words into N-gram model
- Words embedding into Glove model on news headline
- Words embedding with Doc2Vec on news content
- Use normalized methods to transform stock price including MinMaxScaler



### Predictive Modeling

#### Base Model (~52% Acc):

- Naïve Bayes, random forest, linear regression, and logistic regression

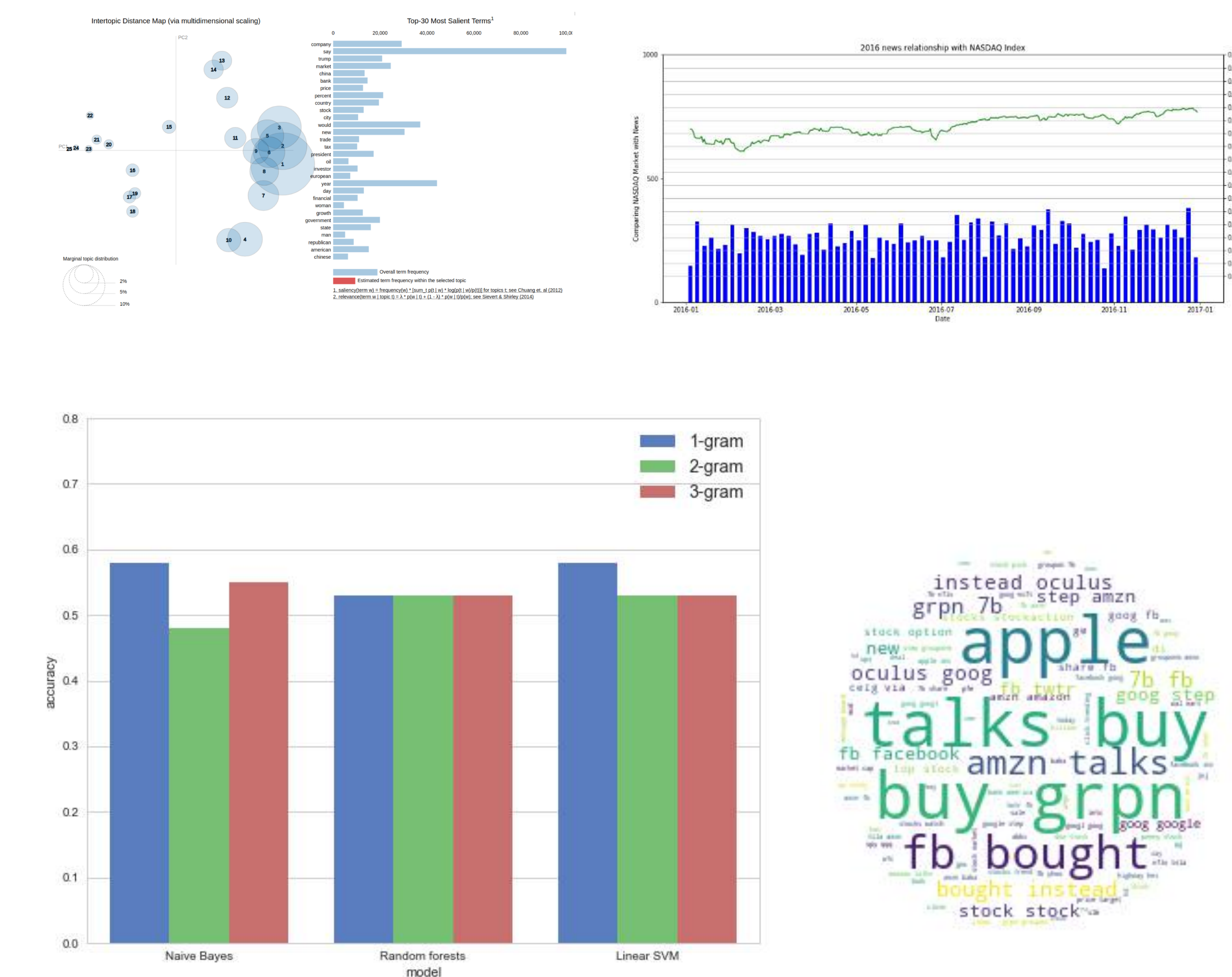
#### LSTM Classification Model (54.5% Acc):

- Trained a recurrent neural network (LSTM) to predict whether stock price rises or drops in next following days

#### LSTM Regression Model:

- Perform regression to predict the stock price in next few days given the sentiment on news headlines and minmax featured stock price

## Exploratory Data Analysis



## Learning & Future Work

### Learning:

- Data mining is open ended and LDA helps understand the latent topics of corpus.
- NLP methods like Doc2Vec, Glove to do feature engineering on word embedding can boost model performance.
- But difficult to measure the performance of embedding.

### Future Work:

- Visualize high dimensional embedding vectors for Doc2Vec.
- Combine topic models for downstream work with unsupervised method like K-means clustering.
- Fine-tune and build more complex model to fit the embedding text data.

## Model Plot and Web Interface

