

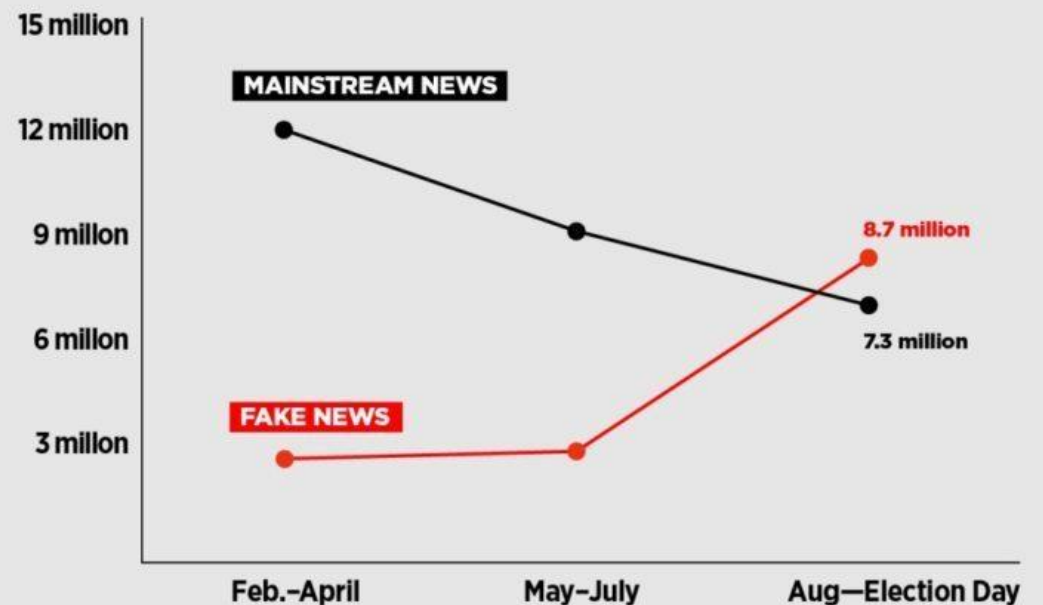
# Fake News Detection With Machine Learning

Vladyslav Maksyk

# Motivation

- INEFFICIENCY OF BASIC COUNTERMEASURES
- EXCESSIVE AMOUNT OF FAKE NEWS IN THE MEDIA
- CURRENT DEVELOPMENT IN THIS AREA

## Total Facebook Engagements for Top 20 Election Stories



ENGAGEMENT REFERS TO THE TOTAL NUMBER OF SHARES, REACTIONS, AND COMMENTS FOR A PIECE OF CONTENT ON FACEBOOK SOURCE: FACEBOOK DATA VIA BUZZSUMO

# Problem Definition

- Develop a machine learning program to identify the credibility of an article bases on the content



# Data

	Credibility	Description	SearchResults	Speaker	Statement	Subjects
0	Half-True	Rudy Giuliani has repeatedly said that he cut ...	[[{"url": "http://www.nytimes.com/2007/01/17/ny..."}]]	Rudy Giuliani	"I cut taxes 23 times when I was mayor of New ...	Taxes
1	Half-True	In a TV ad airing in Iowa and New Hampshire, S...	[[{"url": "http://www.washingtonpost.com/wp-dyn..."}]]	Hillary Clinton	"Hillary stood up for universal health care wh...	Health Care
2	Mostly True	Among the field of Democratic candidates, Sen....	[[{"url": "https://www.washingtonpost.com/wp-dy..."}]]	Joe Biden	"First, he was in favor of my plan, now he's a...	Iraq
3	Mostly True	On ethanol, McCain has maintained a long-stand...	[[{"url": "http://grist.org/article/mccain_fact..."}]]	Mitt Romney	"(McCain) was opposed to ethanol. Now he's for...	Energy
4	True	Romney is right that McCain switched on the ta...	[[{"url": "https://www.senate.gov/?congress=107..."}]]	Mitt Romney	"Senator McCain voted against the Bush tax cut...	Taxes



# WorkFlow



# Before

# After

```
count_vectorizer.get_feature_names()

['00',
 '000',
 '0000',
 '00000031',
 '000035',
 '00006',
 '0001',
 '0001pt',
 '000ft',
 '000km']
```

```
tfidf_vectorizer.get_feature_names()

['aam',
 'aba',
 'aback',
 'abandon',
 'abandoned',
 'abandonment',
 'abas',
 'abbas',
 'abbey',
 'abbreviate']
```

## Text Preprocessing

- ▶ Perform text cleaning operations to remove noise from vectors
- ▶ For Doc2Vec, convert to comma separated word format.

# Doc2Vec model

- ▶ Provides additional information about the entire document
- ▶ Based on Word2Vec

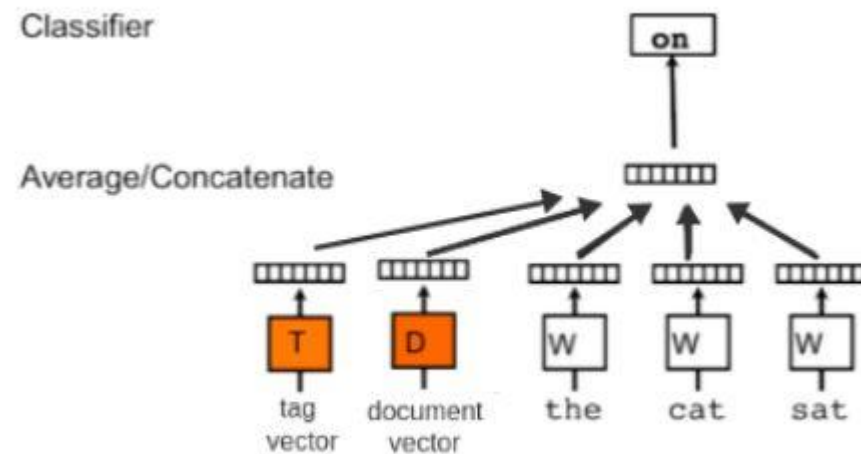


fig 5—doc2vec model with tag vector

# Models Used and Best Results obtained for PoliFact Dataset

Model	Feature Generator	True Accuracy	False Accuracy	F1 Score AVG
Naïve Bayes	N-grams	0.67	0.6	0.64
Support Vector Machine	Tf-idf	0.72	0.68	0.7
Stochastic Gradient Descent	N-grams	0.71	0.65	0.68
Random Forrest	Tf-idf	0.76	0.71	0.71

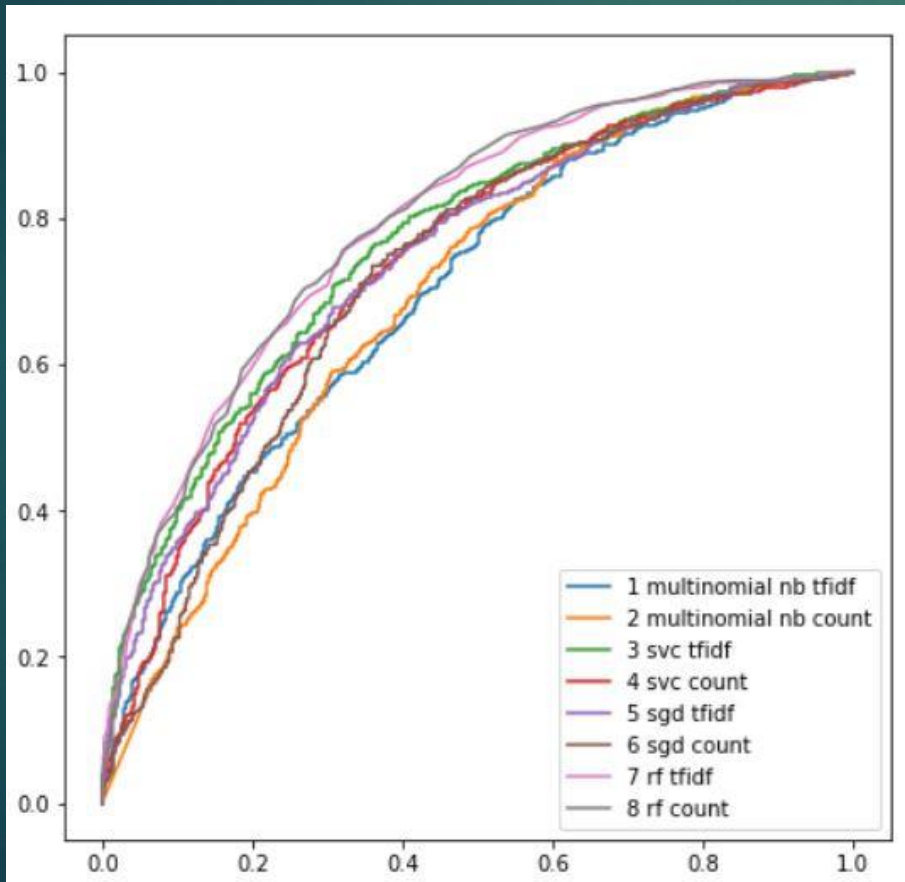


# Models Used and Best Results obtained for Snopes Dataset

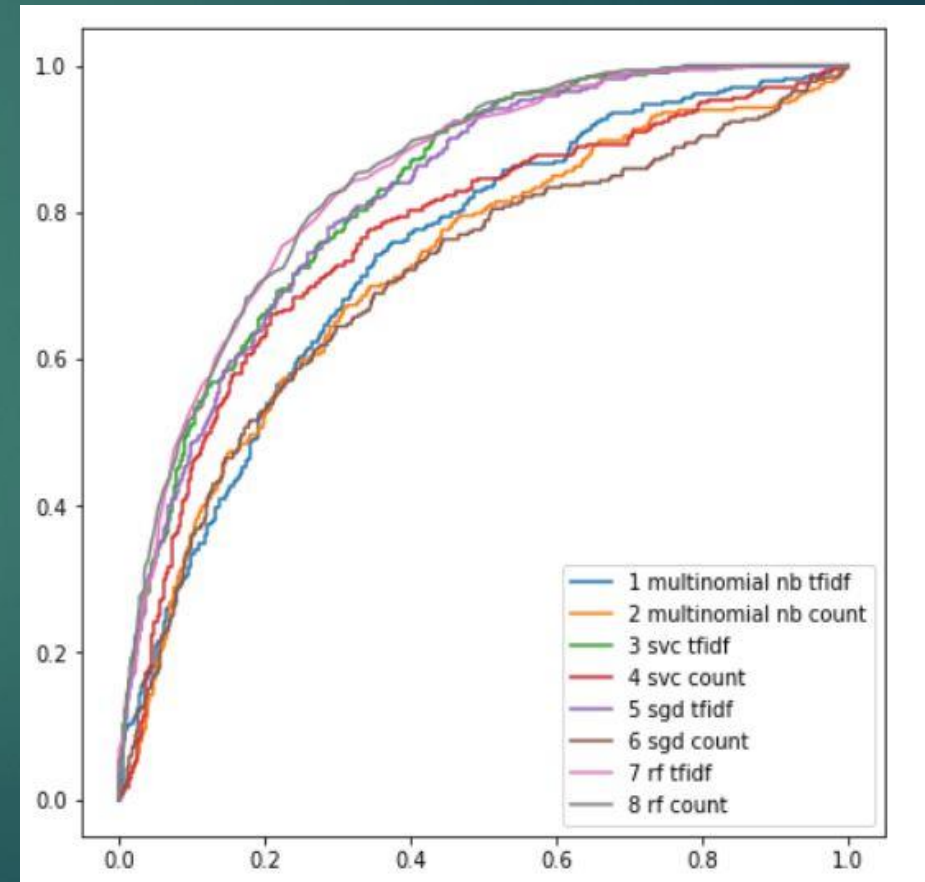
Model	Feature Generator	True Accuracy	False Accuracy	F1 Score AVG
Naïve Bayes	N-grams	0.65	0.71	0.64
Support Vector Machine	Tf-idf	0.66	0.82	0.78
Stochastic Gradient Descent	Tf-idf	0.64	0.83	0.78
Random Forrest	Tf-idf	0.78	0.78	0.74

# Receiver operating characteristic

PolitFact



Snopes



# Introspecting models

## SVC with TF-IDF

```
most_informative_feature_for_binary_classification
```

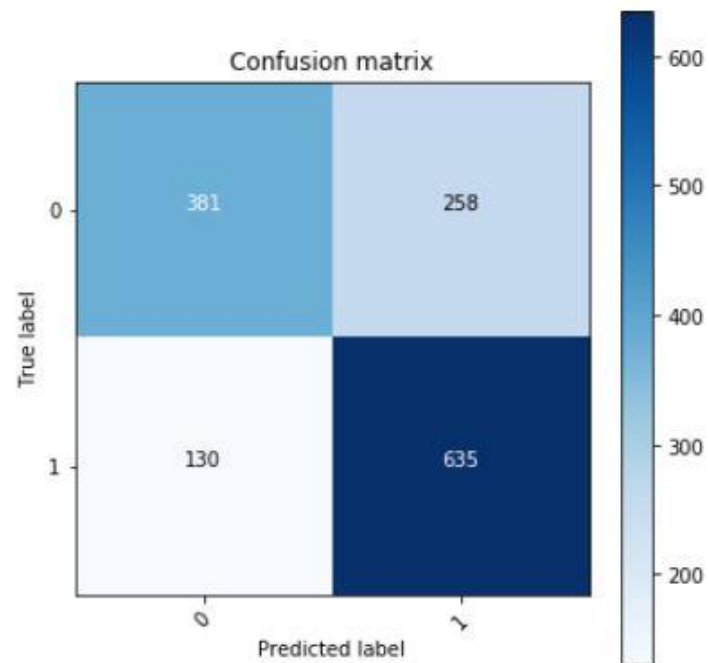
```
{0: [(-3.3080855552934567, 'rumor'),  
      (-2.039782915067682, 'hoax'),  
      (-1.9076314385830677, 'supposedly'),  
      (-1.8983009610970296, 'claim'),  
      (-1.6794732965581247, 'actually'),  
      (-1.5651242904050962, 'legend')],  
1: [(1.2259442705515842, 'sept'),  
     (1.2400878673910565, 'gulf'),  
     (1.2435471257056927, 'tony'),  
     (1.249899535261048, 'unusual'),  
     (1.2606305017735522, 'newspaper'),  
     (1.2710704478740549, 'broadcast'),  
     (1.3344750750756482, 'accurate'),  
     (1.780527326306749, 'additional')]]}
```

## SVC with Count

```
most_informative_feature_for_binary_classification
```

```
{0: [(-0.38556785823005535, 'inaccurate'),  
      (-0.3703933635210672, 'element'),  
      (-0.3006726878957095, 'beginning'),  
      (-0.27297926876036493, 'fallen'),  
      (-0.26510110448557705, 'possibility'),  
      (-0.26469262985093467, 'panel'),  
      (-0.2611315168622069, 'chain'),  
      (-0.2608380587688582, 'truth')],  
1: [(0.2484141322132792, 'region'),  
     (0.24984931903452887, 'reasonable'),  
     (0.25037964510446453, 'caveat'),  
     (0.25499626677951687, 'finance'),  
     (0.2619618279721396, 'unemployed'),  
     (0.2879749171822839, 'context'),  
     (0.31765734337113016, 'careful'),  
     (0.3867982835069873, 'balance')]]}
```

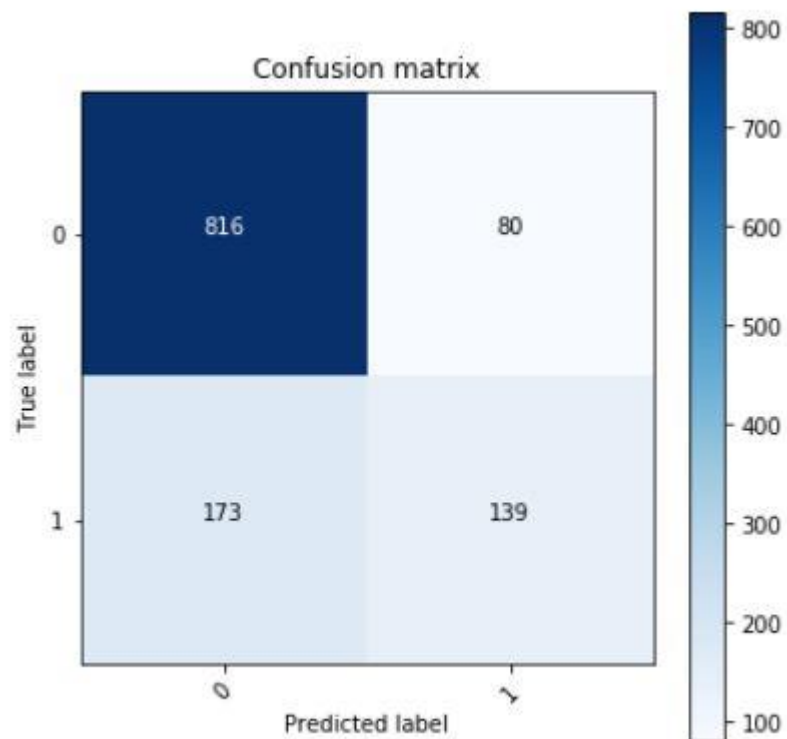
# Random Forrest with TF-IDF



	precision	recall	f1-score	support
false	0.75	0.60	0.66	639
true	0.71	0.83	0.77	765
micro avg	0.72	0.72	0.72	1404
macro avg	0.73	0.71	0.71	1404
weighted avg	0.73	0.72	0.72	1404

Confusion  
Matrix and  
Classification  
Report

# SGD with TD-IDF Snopes



	precision	recall	f1-score	support
false	0.83	0.91	0.87	896
true	0.63	0.45	0.52	312
avg / total	0.78	0.79	0.78	1208

Confusion  
Matrix and  
Classification  
Report

# Main Challenge Faced

- ▶ Content based text classification is just a part of a bigger picture.

All attributes from PolitFact dataset:

**Credibility**

**Speaker**

**Description**

**StatementMetadata**

**EditedBy**

**Subjects**

**Published**

**ReferredLinks**

**ResearchBy**





# Thank You!

► Thanks to Vinay Jayarama Setty for guidance!

