

Mini Project 3

Can we use Natural Language Processing and Machine Learning Models to Detect 'Fake' vs 'Real' News?



Gayle Ferguson

The Data

- Fake.csv 23482 rows
- True.csv 21417 rows
- 4 columns: ‘Title’, ‘Text’, ‘Date’, ‘Subject’
- Target = ‘True’ (1) or ‘Fake’ (0)
- Dataset from Kaggle: <https://www.kaggle.com/datasets/clmentbisaillon/fake-and-real-news-dataset>
- Clément Bisaillon (Owner)

Publications:

Ahmed H, Traore I, Saad S. (2018) “Detecting opinion spams and fake news using text classification”, *Journal of Security and Privacy*, Volume 1(1).

Ahmed H, Traore I, Saad S. (2017) “Detection of Online Fake News Using N-Gram Analysis and Machine Learning Techniques. In: Traore I., Woungang I., Awad A. (eds) *Intelligent, Secure, and Dependable Systems in Distributed and Cloud Environments*. ISDDC 2017. Lecture Notes in Computer Science, Vol 10618. Springer, Cham (pp. 127-138).

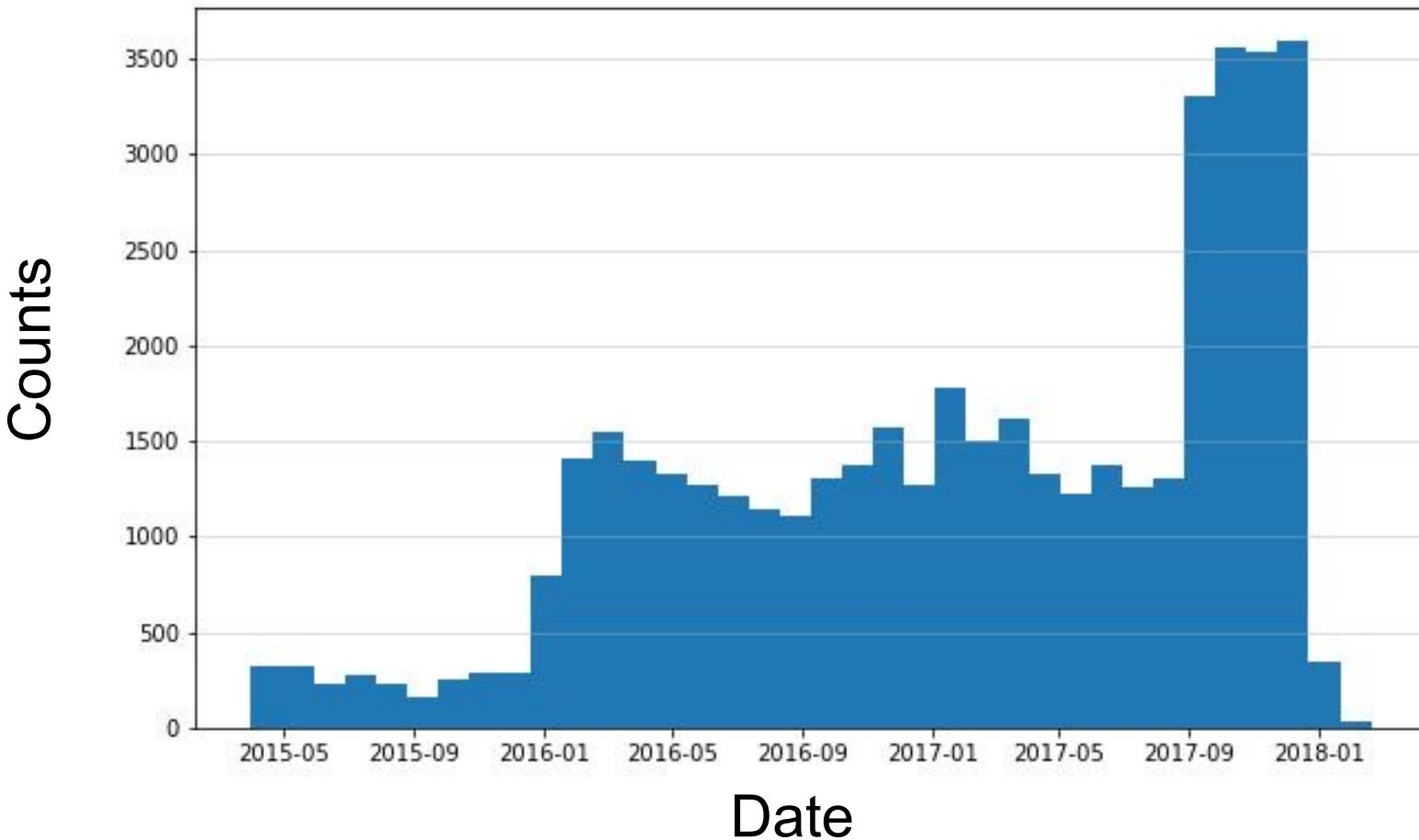
The Data

	title	text	subject	date	target
0	Donald Trump Sends Out Embarrassing New Year's Card	Donald Trump just couldn't wish all Americans a happy new year.	News	December 31, 2017	0
1	Drunk Bragging Trump Staffer Started Russian Hoax	House Intelligence Committee Chairman Devin Nunes has been accused of starting a fake news story about Russian election interference.	News	December 31, 2017	0
2	Sheriff David Clarke Becomes An Internet Joke	On Friday, it was revealed that former Milwaukee Sheriff David Clarke had been secretly working for Donald Trump.	News	December 30, 2017	0
3	Trump Is So Obsessed He Even Has Obama's Name	On Christmas day, Donald Trump announced that he would be changing his name to "The Donald".	News	December 29, 2017	0
4	Pope Francis Just Called Out Donald Trump During His Christmas Message	Pope Francis used his annual Christmas Day message to call out Donald Trump.	News	December 25, 2017	0

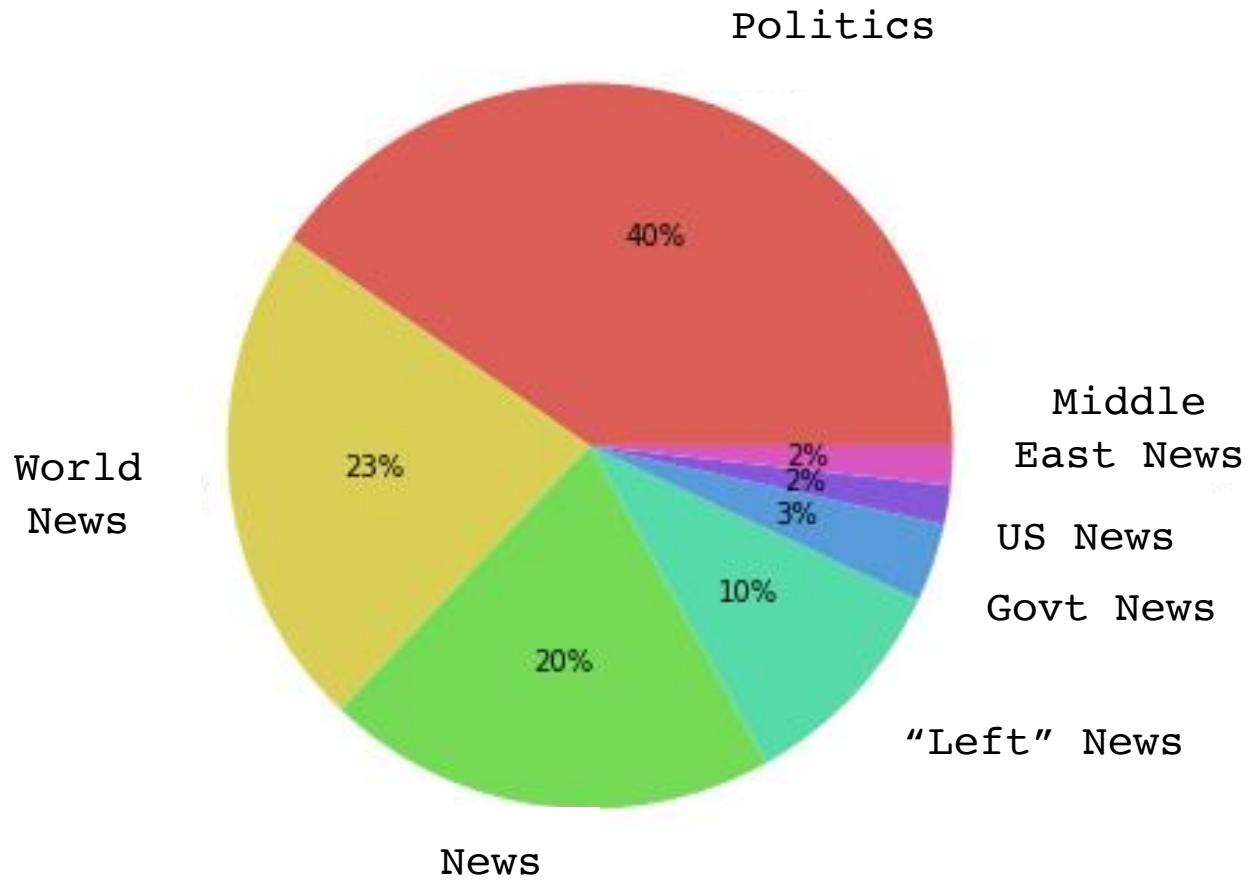
	title	text	subject	date	target
0	As U.S. budget fight looms, Republicans flip to Trump	WASHINGTON (Reuters) - The head of a conservative group that helped Republicans flip the Senate majority in the 2016 election has called on the party to support Donald Trump's budget proposal.	politicsNews	December 31, 2017	1
1	U.S. military to accept transgender recruits officially	WASHINGTON (Reuters) - Transgender people will be allowed to serve openly in the U.S. military for the first time, Defense Secretary Jim Mattis said on Tuesday.	politicsNews	December 29, 2017	1
2	Senior U.S. Republican senator: 'Let Mr. Mueller do his job'	WASHINGTON (Reuters) - The special counsel investigating Russian election interference has been given a wide mandate by the Senate, a senior Republican senator said on Tuesday.	politicsNews	December 31, 2017	1
3	FBI Russia probe helped by Australian diplomat	WASHINGTON (Reuters) - Trump campaign adviser George Papadopoulos' Australian diplomat friend helped the FBI's Russia probe, according to a source familiar with the matter.	politicsNews	December 30, 2017	1
4	Trump wants Postal Service to charge 'much more'	SEATTLE/WASHINGTON (Reuters) - President Donald Trump has told the Postmaster General to consider raising postage rates, a source close to the White House said on Tuesday.	politicsNews	December 29, 2017	1

Article Counts over Time

The articles span May 2015 – Jan 2018,
with the majority from late 2017

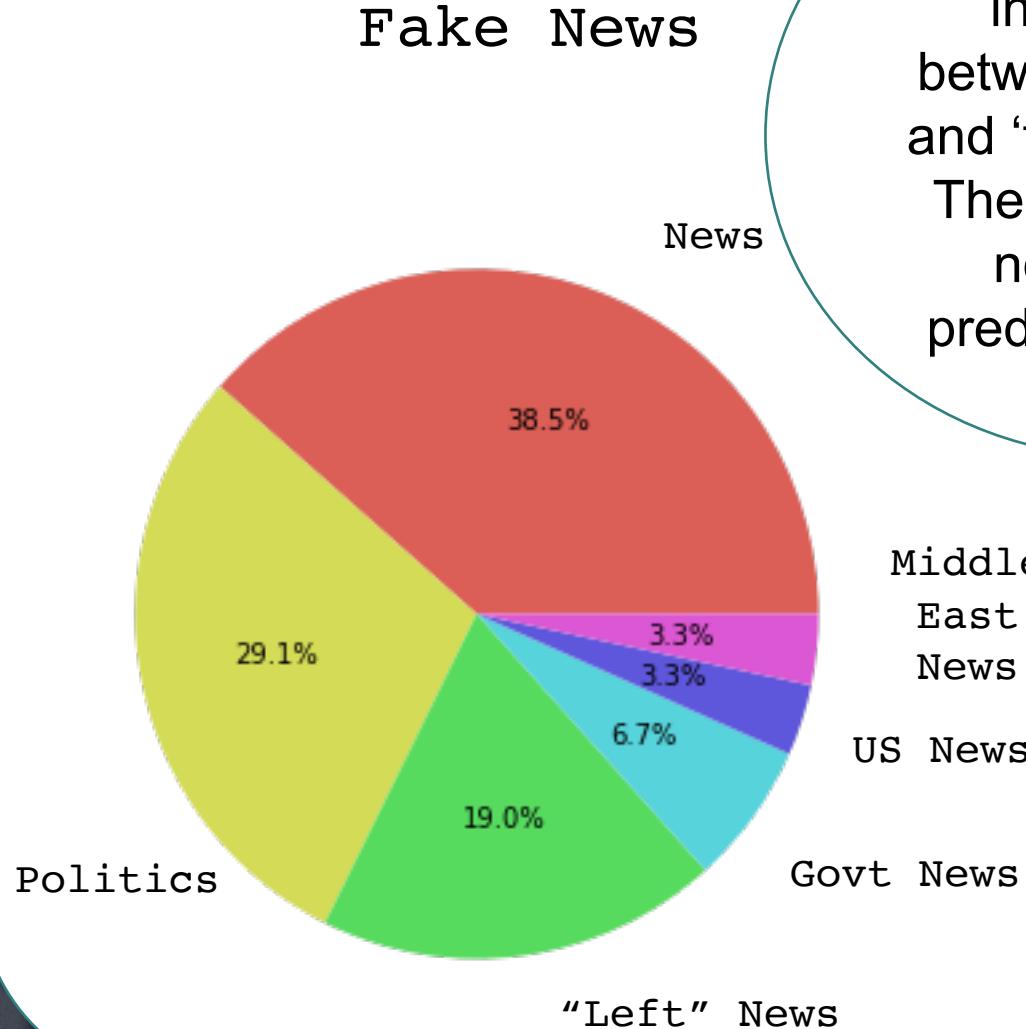


Corpus Classification by Subject



Could the 'subject' labels be a useful predictive feature?

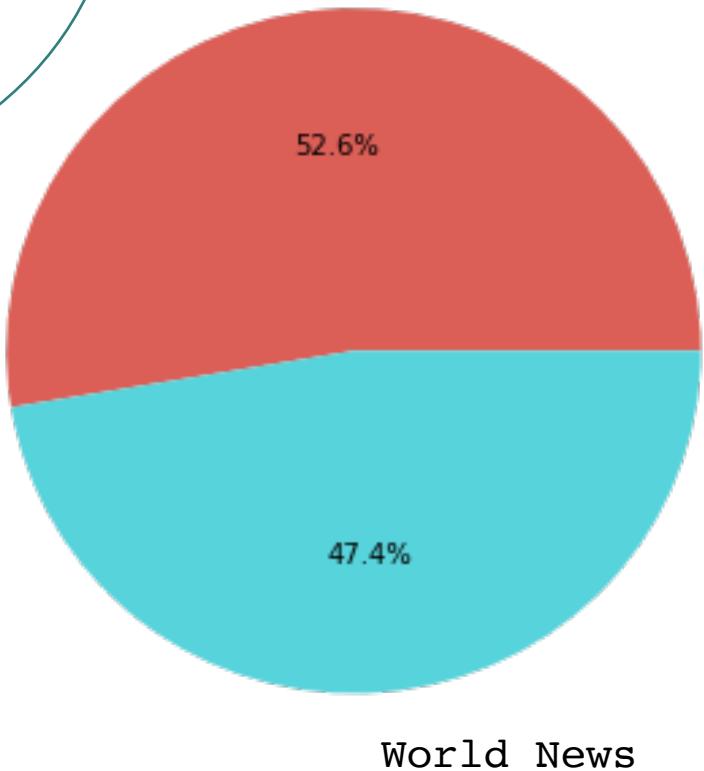
Corpus Classification by Subject



Subject labelling is inconsistent between the 'fake' and 'true' datasets. Therefore, this is not a useful predictive feature

True News

Politics



Do the two classes differ in word content?

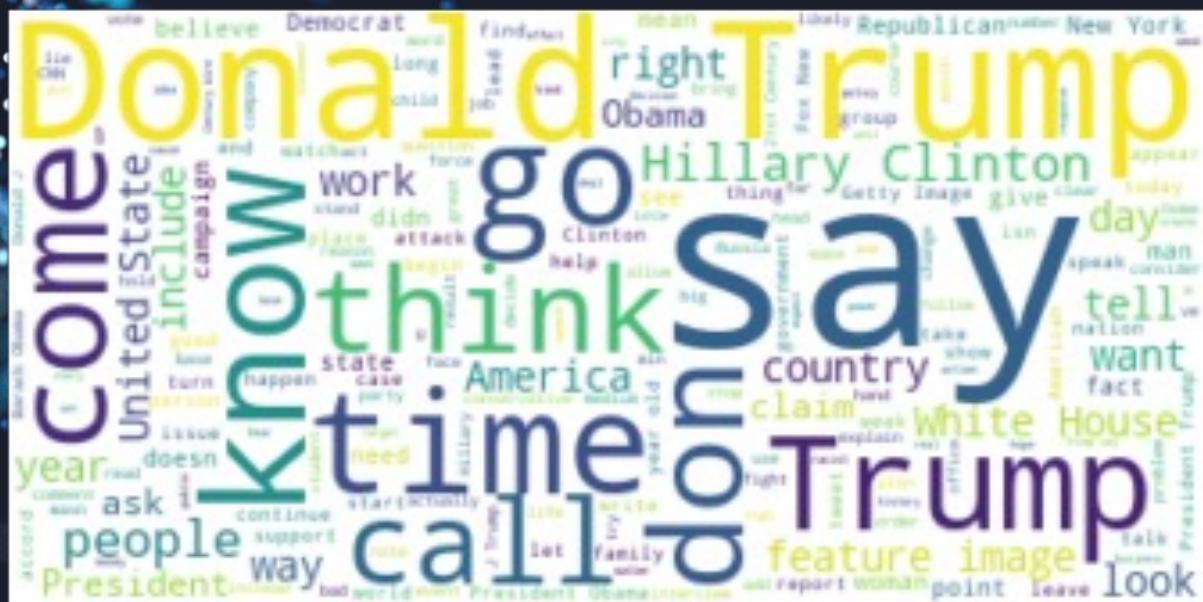
Fake News WordCloud!

38850618 words in the combination of all documents.

True News WordCloud!

34881474 words in the combination of all documents.

They differ surprisingly little!



Workflow

test/train split (20:80)

1. Data Cleaning

- Removal of STOPWORDS & punctuation
- Stemming
- Lemmatization

2. Feature Engineering

- Word count, character count, word density, punctuation count, Title Case count, UPPER CASE count on both '**title**' and '**text**'
- Count Vectorization ('**title**' and '**text**')
- TF-IDF Vectorizer ('**title**' and '**text**')
 - bi-grams, tri-grams on both words and characters
- Topic modelling (LDA) ('**title**' only)

Libraries:

Pandas and NumPy
Regex
SpaCy
SciKitLearn
Matplotlib / Seaborn

Dataset 1

Dataset 2

Dimension reduction
TruncatedSVD

Concatenation

Results of preliminary modelling (non-optimized)

Dataset 1

(35910, 26)

Dataset 2

(35910, 200)

	Dataset 1 (X_train)	Dataset 2 (X_train_combined_2)
Naïve Bayes	0.849744	0.853976
Logistic Regression	0.922143	0.973936
LinearSVC	0.715972	0.972934
Random Forest	0.963021	0.980062
Gradient Boosting	0.946536	0.917465

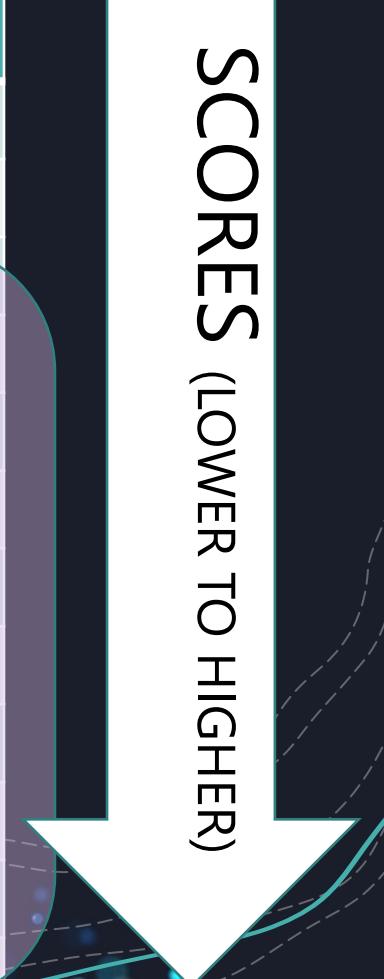
Accuracy Scores for Cross-Validated Models on Dataset 1

Model	Set	Accuracy Score	Meta Classifier
Linear SVC	Train	0.81 +/- 0.01	
	Test	0.70 +/- 0.09	
Naïve Bayes (Gaussian)	Train	0.81 +/- 0.01	
	Test	0.83 +/- 0.00	
Logistic Regression	Train	0.89 +/- 0.01	
	Test	0.90 +/- 0.01	
Gradient Boosting	Train	0.95 +/- 0.00	
	Test	0.94 +/- 0.01	
Random Forest	Train	0.96 +/- 0.00	
	Test	0.94 +/- 0.01	
Stacking Classifier	Train	0.96 +/- 0.00	Random Forest
	Test	0.94 +/- 0.01	

SCORES (LOWER TO HIGHER)

Accuracy scores for cross-validated models on Dataset 2

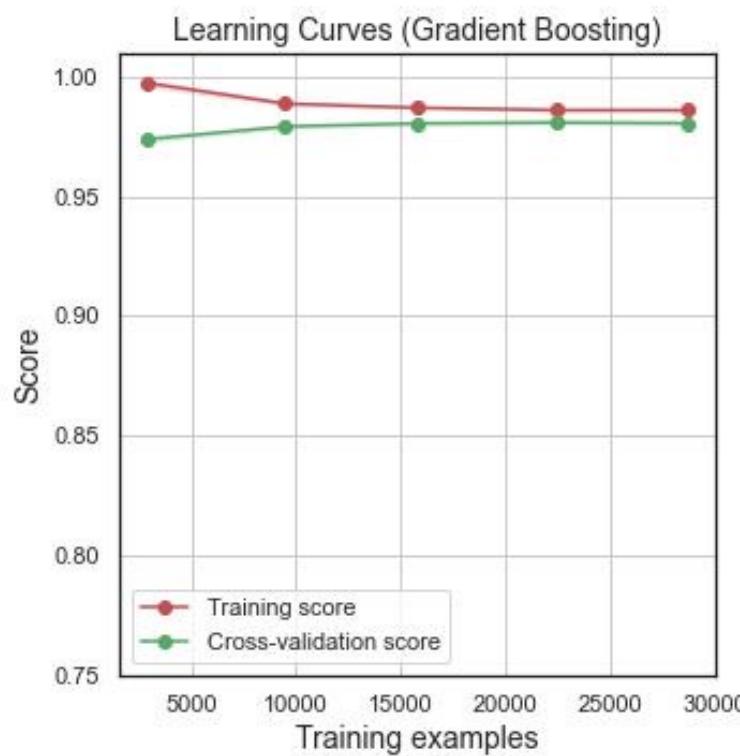
Model	Set	Accuracy Score	Meta Classifier
Naïve Bayes (Gaussian)	Train	0.85 +/- 0.00	
	Test	0.85 +/- 0.01	
Logistic Regression	Train	0.97 +/- 0.00	
	Test	0.97 +/- 0.00	
Random Forest	Train	0.98 +/- 0.00	
	Test	0.97 +/- 0.00	
Stacking Classifier	Train	0.98 +/- 0.00	Random Forest
	Test	0.97 +/- 0.00	
Linear SVC	Train	0.98 +/- 0.00	
	Test	0.98 +/- 0.00	
Gradient Boosting	Train	0.98 +/- 0.00	
	Test	0.98 +/- 0.00	

SCORES (LOWER TO HIGHER) 

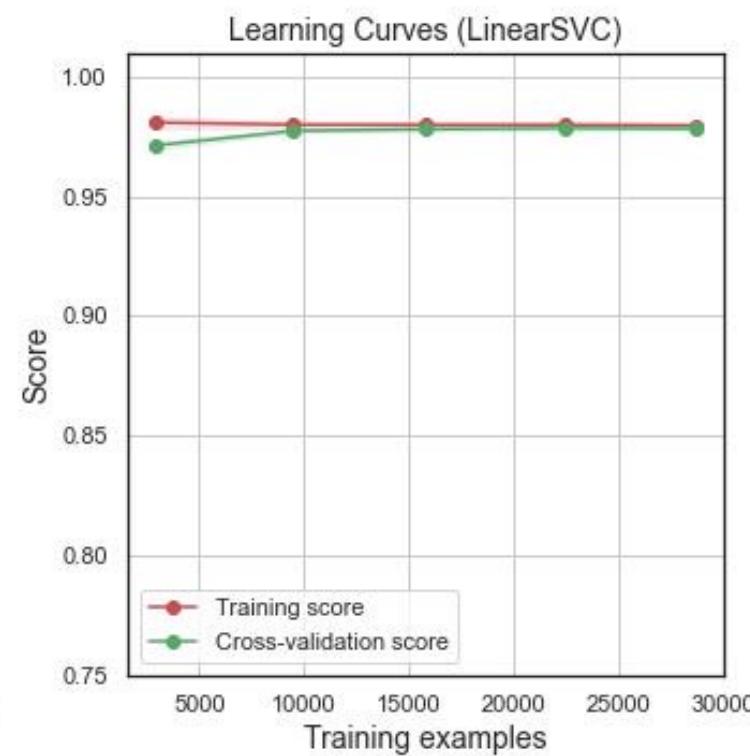
Model Evaluation

Dataset 2

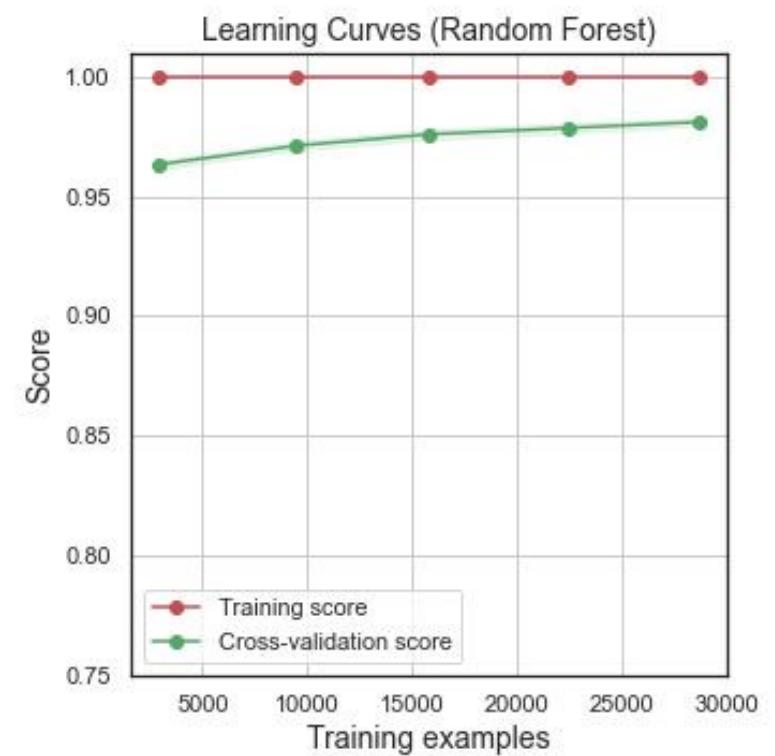
Gradient Boosting



Linear SVC



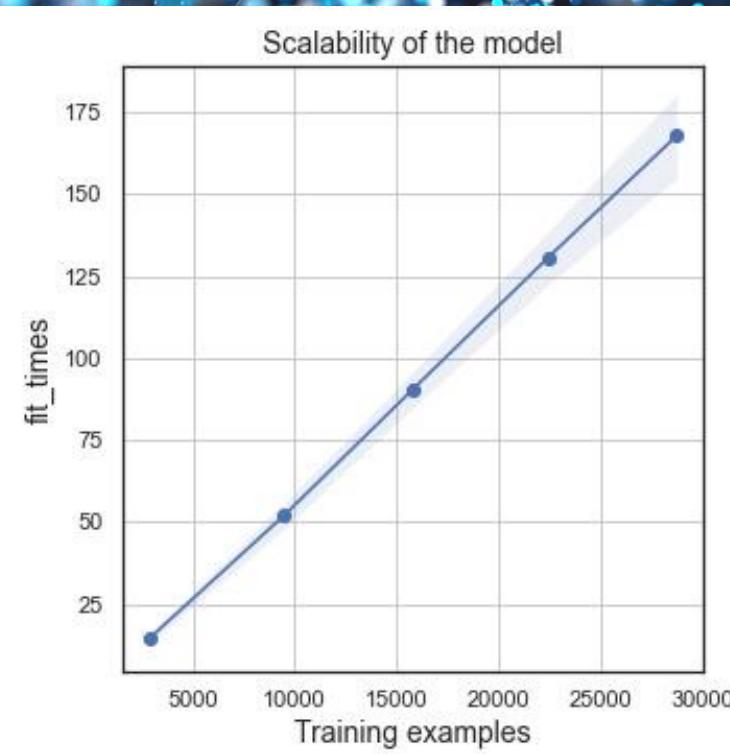
Random Forest



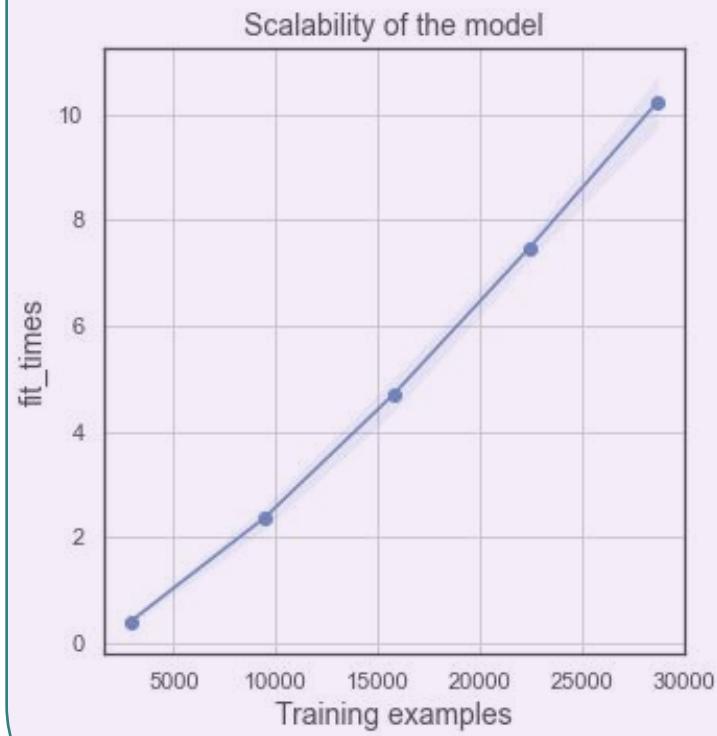
Model Evaluation

Dataset 2

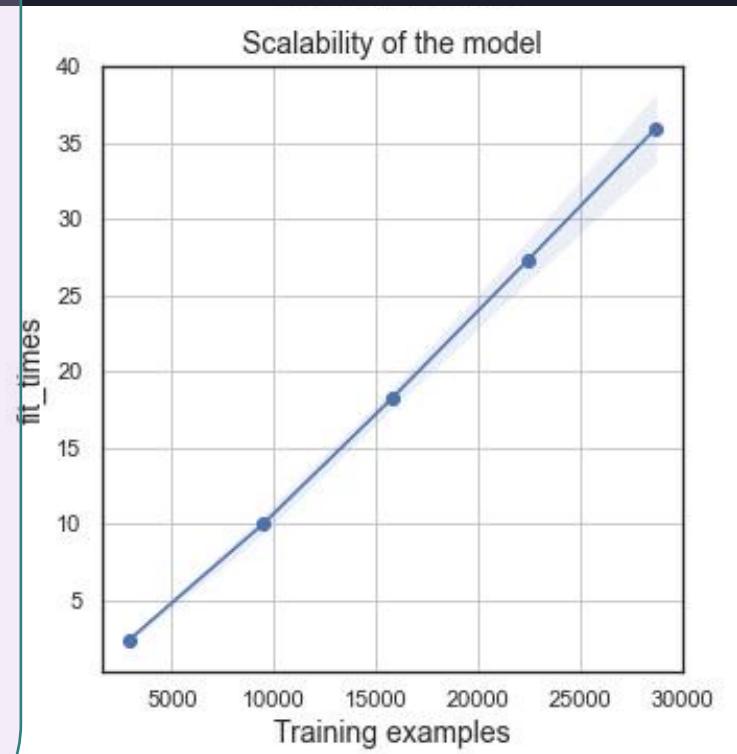
Gradient Boosting



Linear SVC



Random Forest

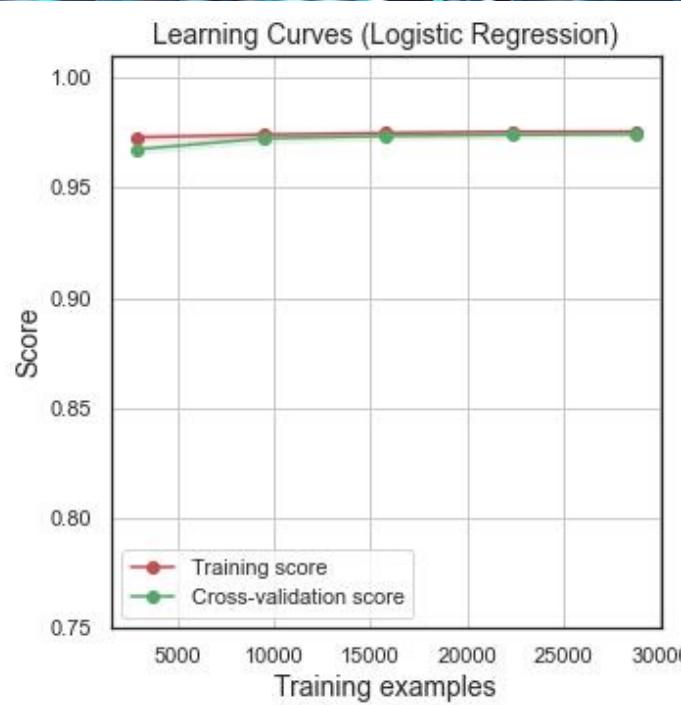


All things being equal, linear SVC wins on speed!

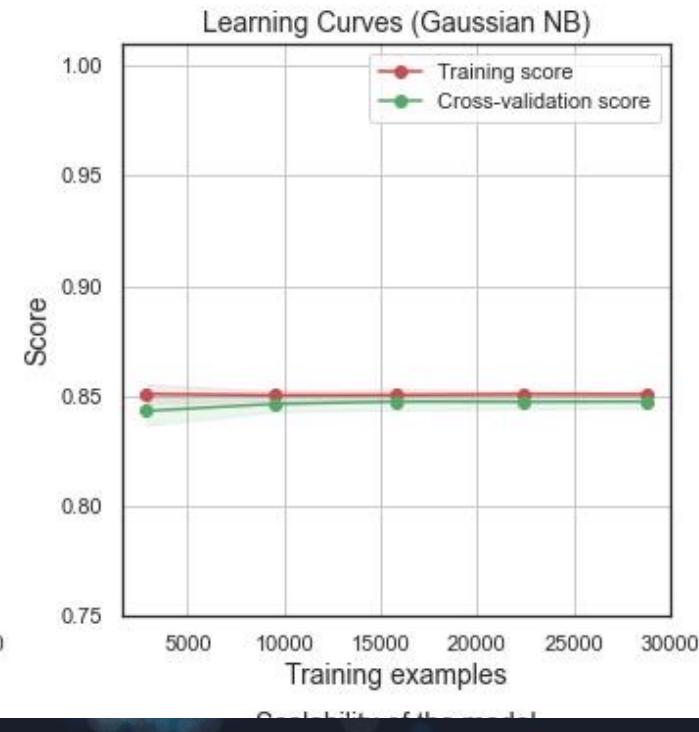
Model Evaluation

Dataset 2

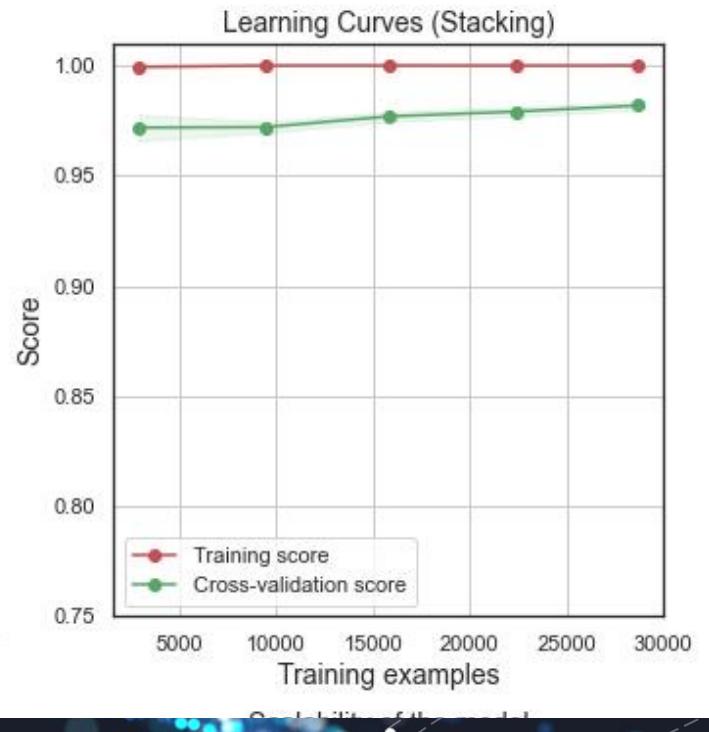
Logistic Regression



Gaussian NB



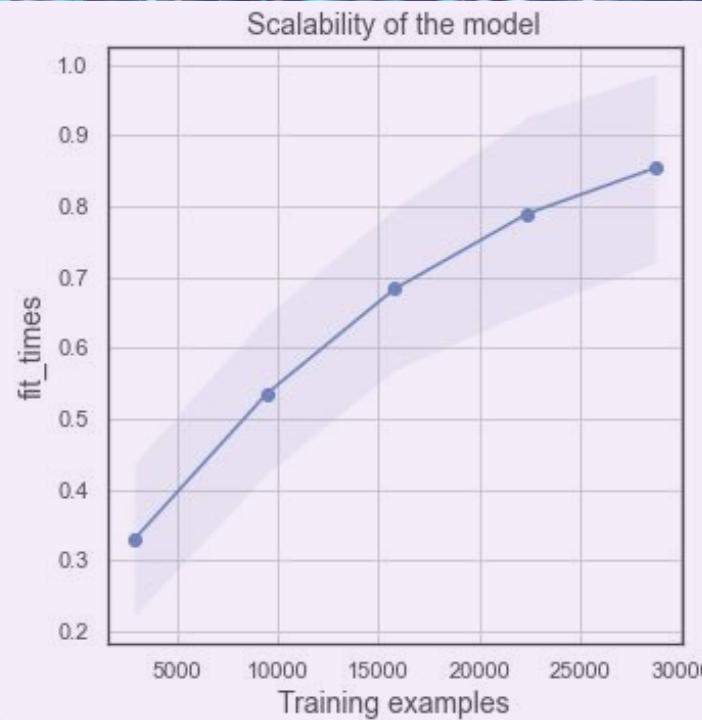
Stacking



Model Evaluation

Dataset 2

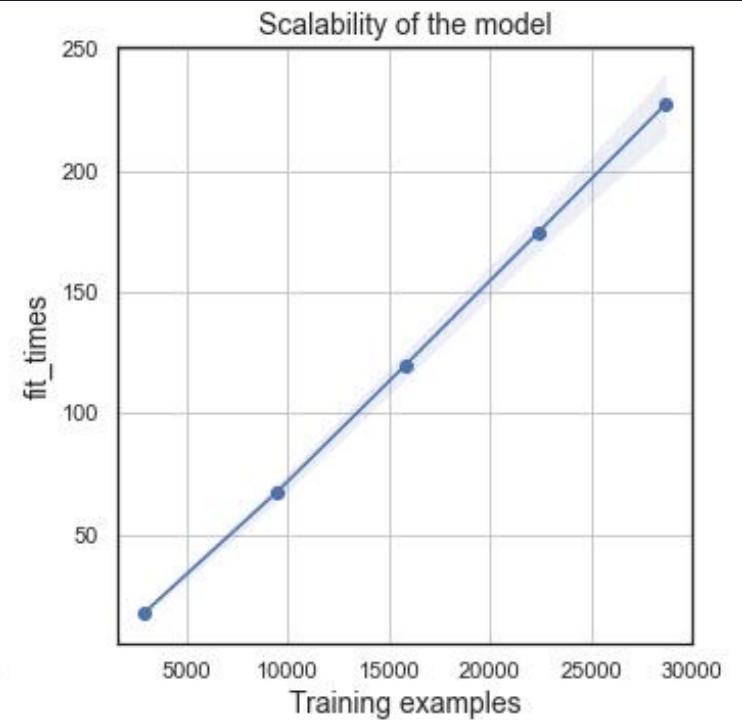
Logistic Regression



Gaussian NB



Stacking



Logistic Regression wins on speed!

Model Evaluation

Dataset 2

Model	Accuracy (test dataset)	Precision (the extent to which we are not classifying fake news as true)	Recall / Sensitivity (the extent to which we are finding the true news)	AUC	SPEED (FAST TO SLOW)
Logistic Regression	0.9732	0.9696	0.9744	0.9955	
Linear SVC	0.9747	0.9759	0.9711	0.9964	
Random Forest	0.9810	0.9784	0.9818	0.9979	
Gradient Boosting	0.9738	0.9719	0.9735	0.9969	
Stacking Classifier	0.9792	0.9779	0.9786	0.9928	

Model Evaluation (Random Forest)

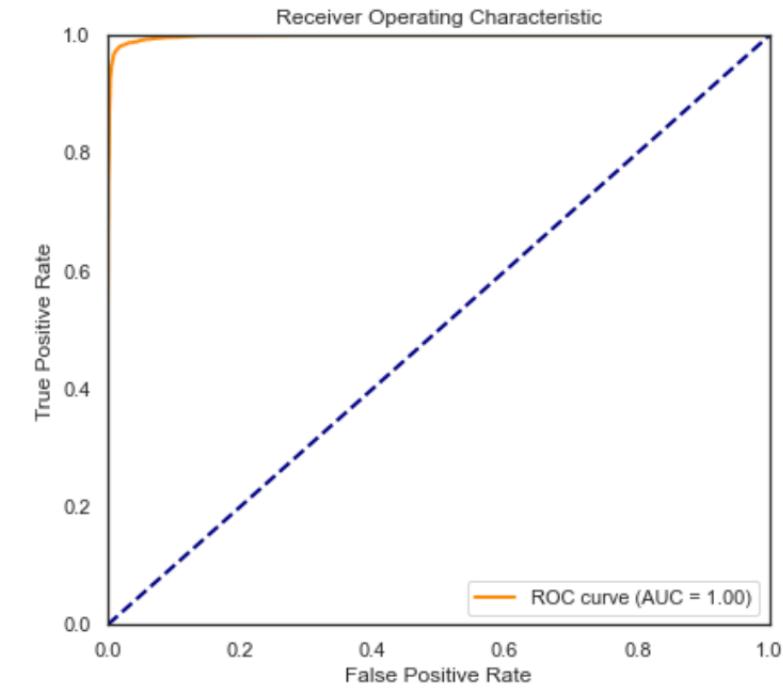
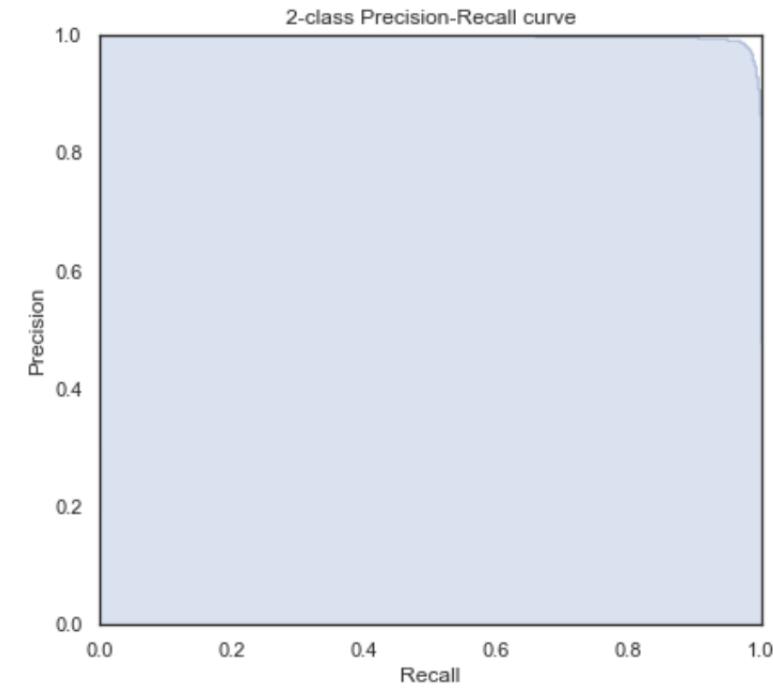
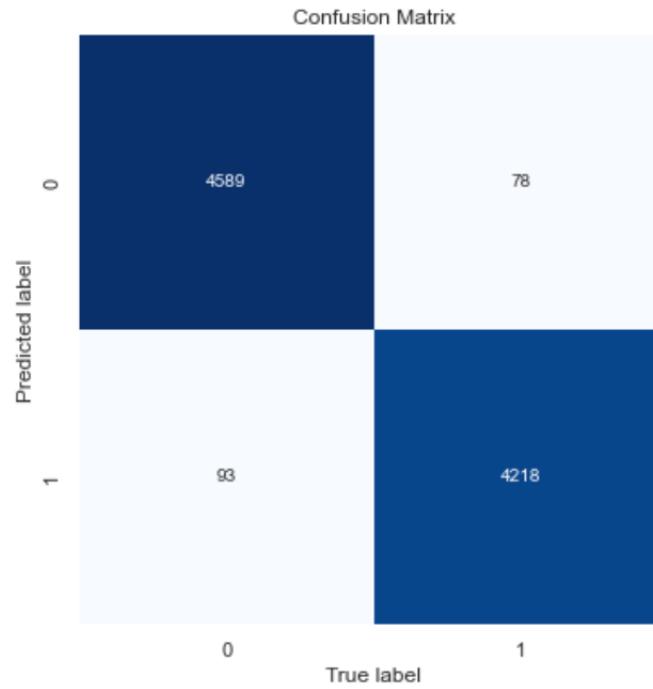
Dataset 2

Random Forest

Accuracy : 0.9810
Precision: 0.9784
Recall : 0.9818
ROC AUC : 0.9979

TP / N] Proportion of predicted labels that match the true labels. Best: 1, Worst: 0
TP / (TP + FP)] Not to label a negative sample as positive. Best: 1, Worst: 0
TP / (TP + FN)] Find all the positive samples. Best: 1, Worst: 0
Best: 1, Worst: < 0.5

TP: True Positives, FP: False Positives, TN: True Negatives, FN: False Negatives, N: Number of samples



Conclusions

- All models (with the exception of Naïve Bayes) were highly accurate, even without optimization; very little trade-off between precision and recall.
- If speed and computational demands are of primary concern, we'd recommend a Linear SVC model or Logistic Regression model
- For the highest accuracy, Random Forest is recommended: this model is a good compromise between speed/computational demands and accuracy but is recommended only for very large volumes of training data.

NEXT STEPS:

- Optimization of the hyperparameters for the Random Forest (and other models);
- Topic modelling on the full text may improve model accuracy, but this is unlikely to be worth the effort.
- Exploration of different meta-classifiers in a Stacking Model to see if accuracy can be further improved.

