

Detection of Fake News With Machine Learning

Gianluca Turcatel, PhD

Capstone Project, September 2020

The Problem

- Definition: *Fabricated stories with no verifiable facts, sources or quotes*
- *Not a new problem (Condorcet 1795)*
- Goal: manipulate and mislead the reader
- Fake news spread faster than real news
- Widespread problem:
 - 2018, Statistica Survey: 52% of respondents said that “*online website reports fake news regularly*” (2018, Statistica)
 - <https://www.tweetgen.com/>

Who Might Care?

Social Media Platforms



Media Outlets



The Reader



The Government

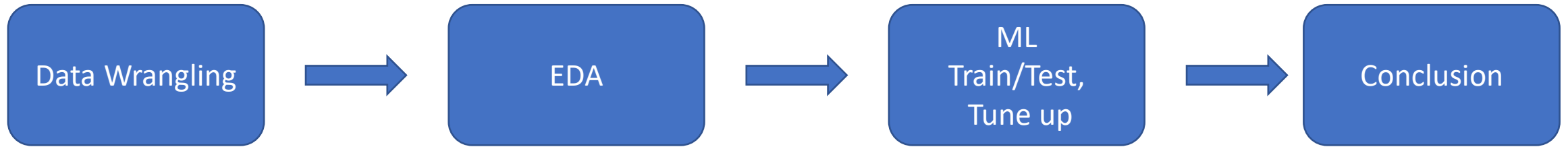


The Data

- Data Source :
 - Real News: Reuters (24417)
 - Fake News: Kaggle dataset(23502)
- Data Structure:
 - Date: March 31st , 2015, to February 19th, 2018
 - Title
 - Text
 - Category (Government News, Middle-east, News, US_News, left-news, politics, politicsNews or worldnews)

<https://www.kaggle.com/clmentbisailon/fake-and-real-news-dataset>

How the problem was tackled



Data Wrangling - Overview

STEP 1:

- Merge datasets
- Drop duplicates, incomplete, missing values



STEP 2:

- Removal of cross site scripts
- Tokenize most frequent hyperlinks
- Tokenize vulgarities
- Expand out verbal contraction
- Linguistic indexes calculation
 - Lemmatize text

Data Wrangling - Removal of cross site scripts

.....The children s story was so inspiring that it was picked up by the popular Facebook page Have a Gay Day. // < ![CDATA[// < ![CDATA[// < ![CDATA[(function(d, s, id) { var js, fjs = d.getElementsByTagName(s)[0]; if (d.getElementById(id)) return; js = d.createElement(s); js.id = id; js.src = "//connect.facebook.net/en_US/sdk.js#xfbml=1&version=v2.3"; fjs.parentNode.insertBefore(js, fjs);}(document, 'script', 'facebook-jssdk')); //]]>This.Posted by Have A Gay Day on Tuesday,

.... Congress just lacks a spine to do it! (function(d, s, id) { var js, fjs = d.getElementsByTagName(s)[0]; if (d.getElementById(id)) return; js = d.createElement(s); js.id = id; js.src = "//connect.facebook.net/en_US/sdk.js#xfbml=1&version=v2.3"; fjs.parentNode.insertBefore(js, fjs);}(document, 'script', 'facebook-jssdk'));Senator Elizabeth Warren goes....

Data Wrangling: Hyperlinks Tokenization

Common hyperlinks:

- Facebook
- Bitly
- Youtube
- Twitter
- Tmsnrt



Replaced with **11-**
character long tokens

....we were an enemy of the people. Here s the video via
YouTube:<https://www.youtube.com/watch?v=KyipRvdPyLE&feature=youtu.be>Donald Trump is the true

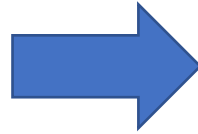


....we were an enemy of the people. Here s the video via
YouTube: **youtubelink**. Donald Trump is the true

Data Wrangling – Linguistic Indexes Calculation

STEP 1: Calculate number of:

Sentences
Words
Characters,
Syllables,
Polysyllabic words
Long word*
Complex Words**
Monosyllabic words



STEP 2: Compute the linguistic scores:

Flesh Kincaid
Flesh Read Easy
Simple Measure of Gobbledygook (SMOG)
Gunning Fox
Lasbarhetsindex (LIX)
Coleman-Liau (CLI)
Automated Readability Index (ARI)
Gulpease

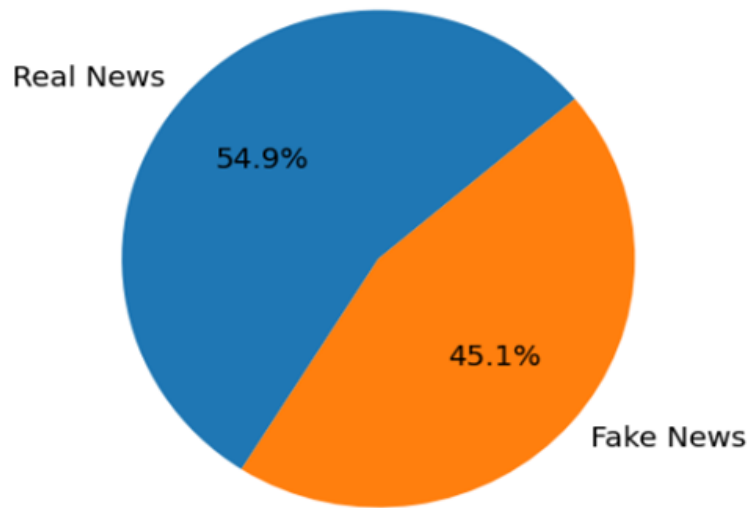
*Words with more than 6 characters;

** Words with more 3 syllables

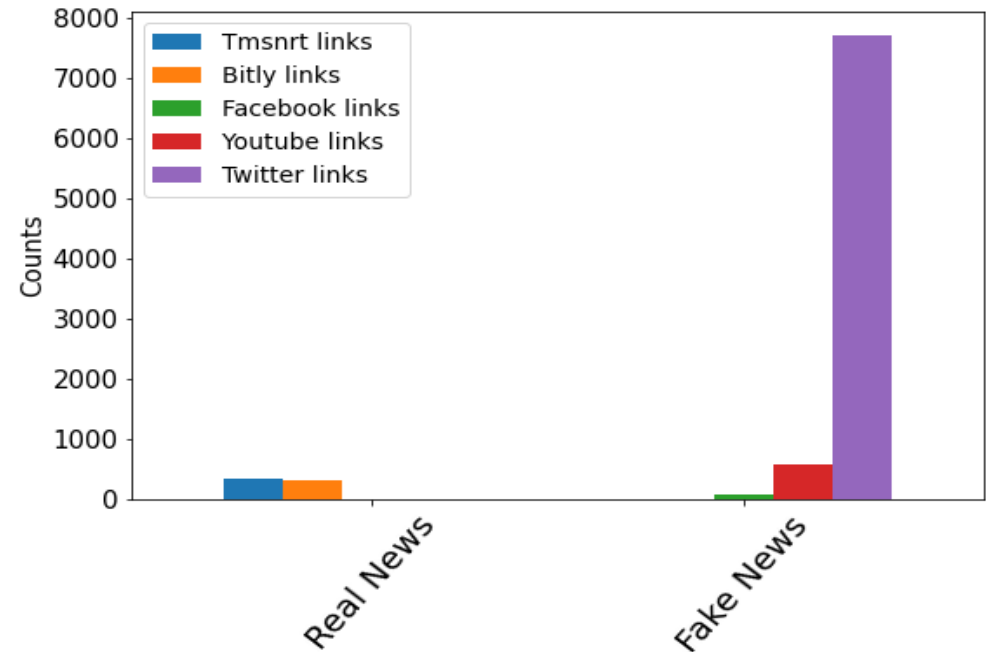
Example of linguistic index:

$$CLI = 0.0588 \frac{\text{characters}}{\text{words}} - 0.296 \frac{\text{sentences}}{\text{words}} - 15.8$$

EDA – News and Hyperlinks Distribution



The two classes are well balanced



Dichotomic distribution of the hyperlinks

EDA – Linguistic Features Trimming

32 text, linguistic Features:

- 16 for the title
- 16 for the text



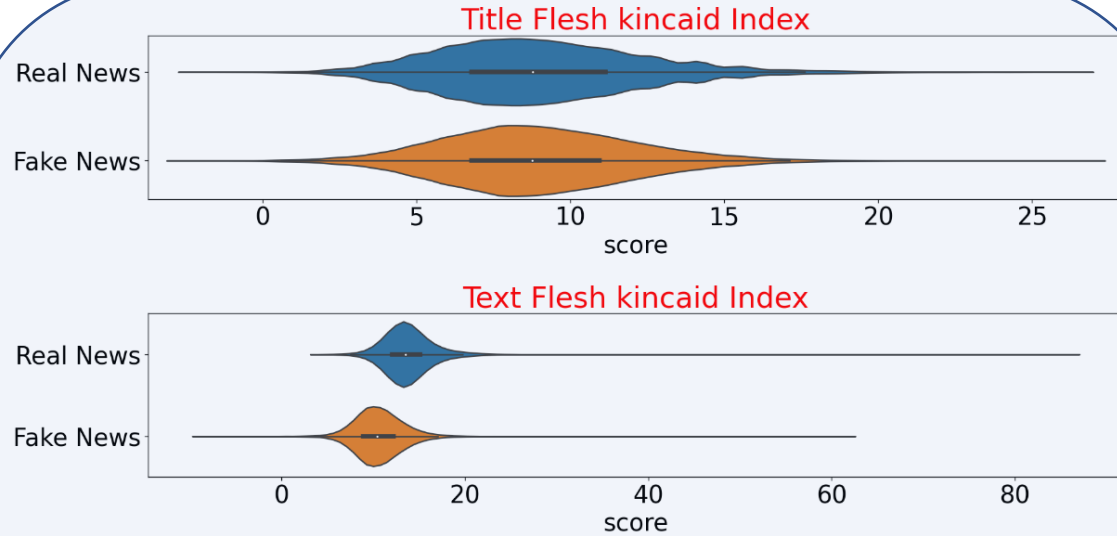
13 text, linguistic Features:

- 8 for the title
- 5 for the text

sent_title	1.00	0.25	-0.01	0.09	-0.21	-0.32	-0.27	0.56	-0.01	-0.05	-0.15	-0.14	0.18
Count_words_title	0.25	1.00	-0.00	0.29	-0.01	0.34	0.00	-0.27	0.11	-0.05	-0.28	-0.35	0.27
Complex_words_title	-0.01	-0.00	1.00	0.31	0.59	0.17	0.45	-0.31	0.02	0.12	0.07	0.10	-0.05
Long_words_title	0.09	0.29	0.31	1.00	0.63	0.59	0.69	-0.54	0.04	0.03	0.00	0.07	-0.01
Flesh_Kin_Grade_title	-0.21	-0.01	0.59	0.63	1.00	0.71	0.84	-0.63	-0.02	0.05	0.14	0.17	-0.12
smog_title	-0.32	0.34	0.17	0.59	0.71	1.00	0.66	-0.74	0.02	-0.02	0.01	-0.00	-0.01
ARI_title	-0.27	0.00	0.45	0.69	0.84	0.66	1.00	-0.76	0.01	0.06	0.10	0.18	-0.12
Gulpease_title	0.56	-0.27	-0.31	-0.54	-0.63	-0.74	-0.76	1.00	-0.05	-0.04	-0.04	-0.05	0.08
sent_text	-0.01	0.11	0.02	0.04	-0.02	0.02	0.01	-0.05	1.00	0.77	-0.27	-0.17	0.15
Complex_words_text	-0.05	-0.05	0.12	0.03	0.05	-0.02	0.06	-0.04	0.77	1.00	0.13	0.19	-0.16
Flesh_Kin_Grade_text	-0.15	-0.28	0.07	0.00	0.14	0.01	0.10	-0.04	-0.27	0.13	1.00	0.74	-0.76
CLI_text	-0.14	-0.35	0.10	0.07	0.17	-0.00	0.18	-0.05	-0.17	0.19	0.74	1.00	-0.79
Gulpease_text	0.18	0.27	-0.05	-0.01	-0.12	-0.01	-0.12	0.08	0.15	-0.16	-0.76	-0.79	1.00
	sent_title	Count_words_title	Complex_words_title	Long_words_title	Flesh_Kin_Grade_title	smog_title	ARI_title	Gulpease_title	sent_text	Complex_words_text	Flesh_Kin_Grade_text	CLI_text	Gulpease_text

EDA – Indexes' Insight

Flesh Kincaid Index

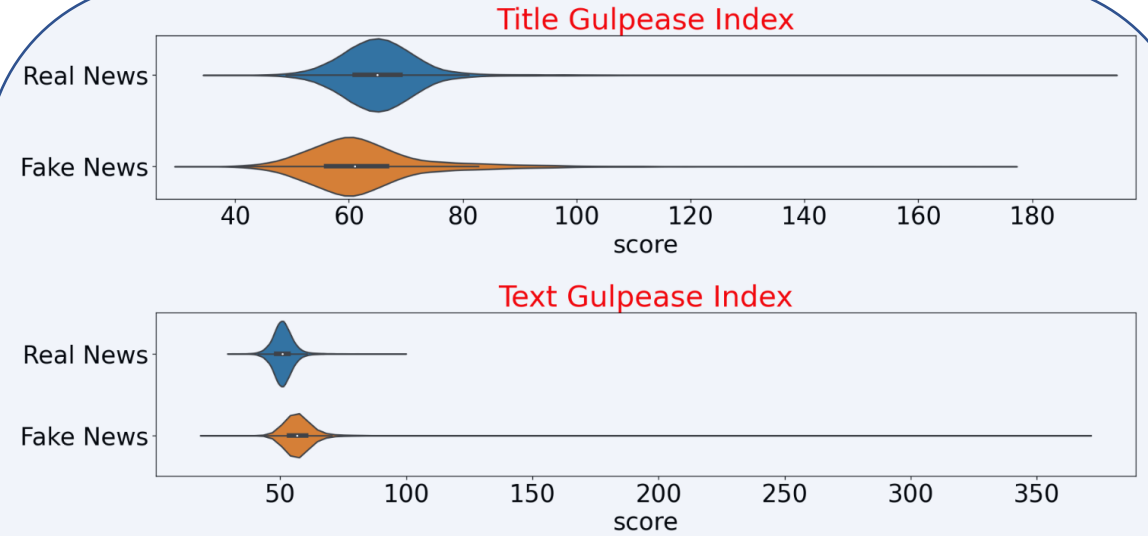


$$0.39 \frac{\text{Total words}}{\text{Total sentences}} + 11.8 \frac{\text{Total Syllables}}{\text{Total words}} - 15.59$$



Text of real news is more complex/sophisticated

Gulpease Index



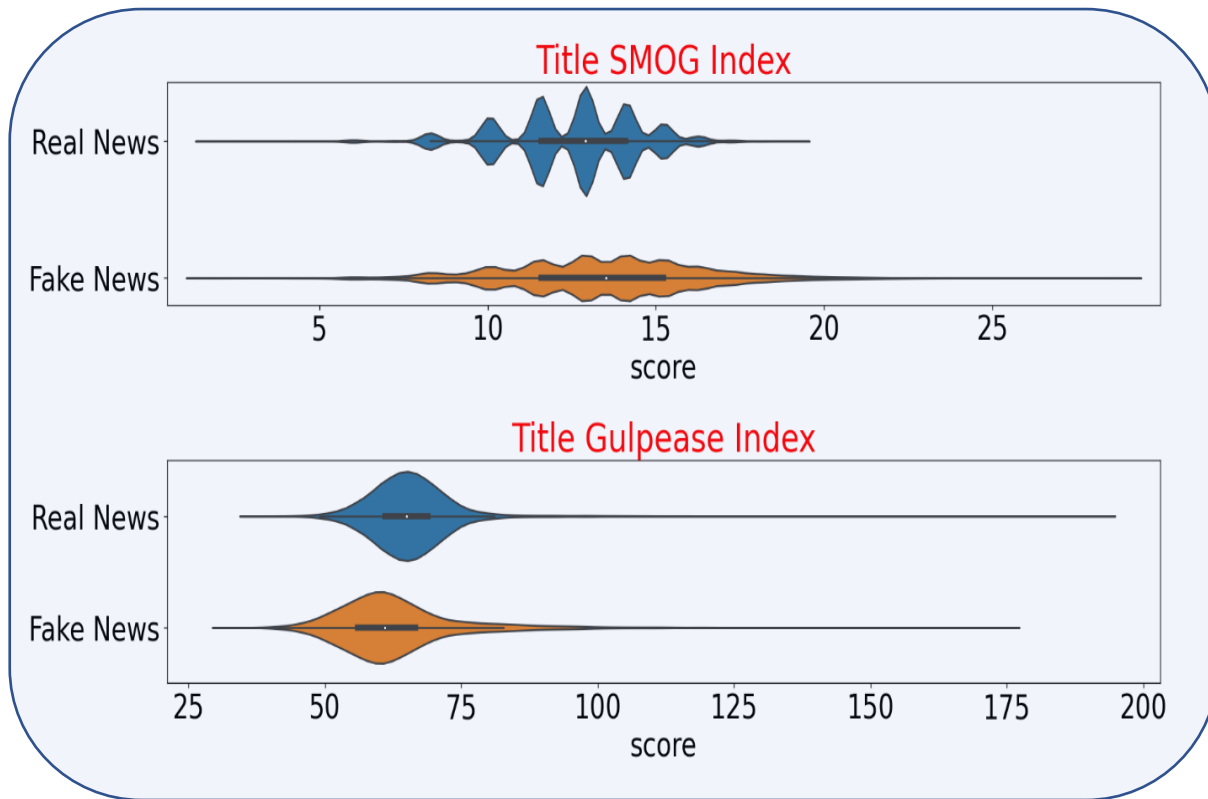
$$89 + \frac{300 * \text{sentences} - 10 * \text{characters}}{\text{words}}$$



Title of real news is shorter/simpler
Text of real news is longer/more complex

EDA – Indexes comparison

SMOG index vs Gulpease Index



$$1.0430 \sqrt{\text{number of polysyllables} \frac{30}{\text{number of sentences}}} + 3.1291$$

$$89 + \frac{300 * \text{sentences} - 10 * \text{characters}}{\text{words}}$$



Each index amplifies one or more specific linguistic properties in the word and sentence and, thus, generates its **own unique distribution**

EDA - Predictive words

Real News



Real news – key takeaways:

- **Bitly/tmsnrt hyperlinks.**
- Name of **non-mainstream politicians** (Odinga - Kenyan politician, Rajoy - Spanish politician, Hariri - Lebanese politician, Gulen - Turkish Islamic scholar) and **name of cities and countries rarely mentioned by the media** (Kirkuk - city in Iraq, Nairobi -Kenya's capital, Marawi - Islamic city, Harare - capital of Zimbabwe, Rakhine - state in Myanmar).

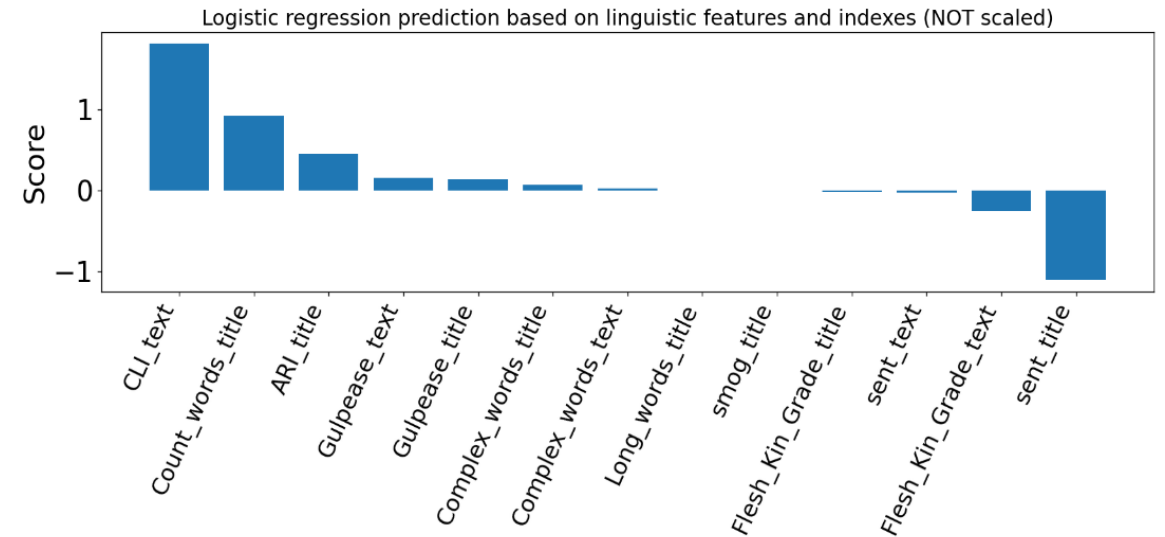
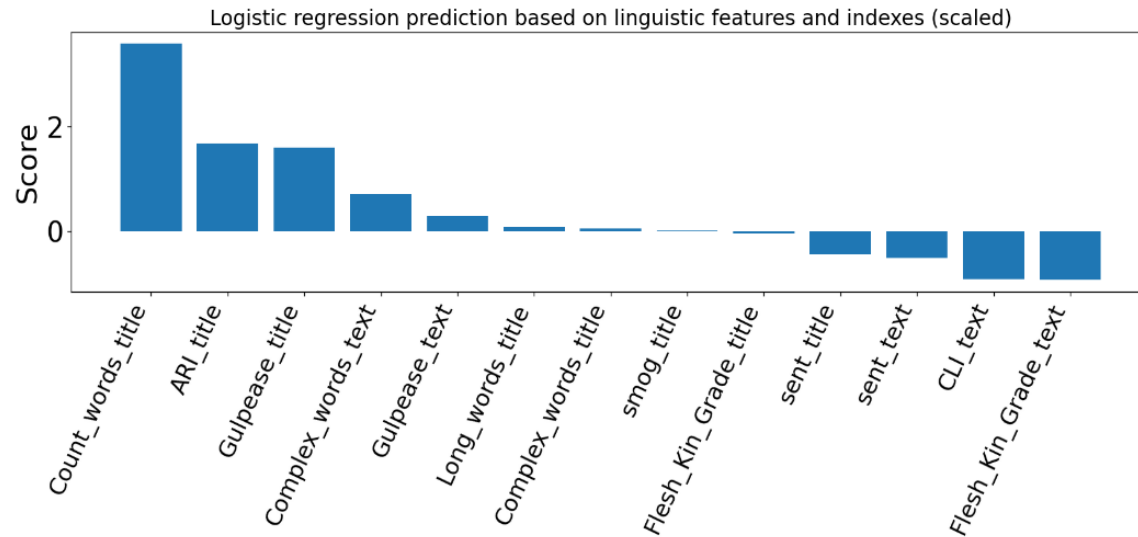
Fake News



Fake news – key takeaways:

- **Twitter and YouTube hyperlinks.**
- **Photography related words** (*getty, somodevilla, screenshot, flickr, getty, raedle, angerer*).
- **Vulgarities** (*f**k, s**t, bigots*).

EDA - Most predictive linguist indexes, text features



Words in the title:

➡ The longer the title, the higher the prob the news is fake

Text Flesh Kincaid index :

➡ The higher the Flesh Kincaid, the higher the prob the news is real

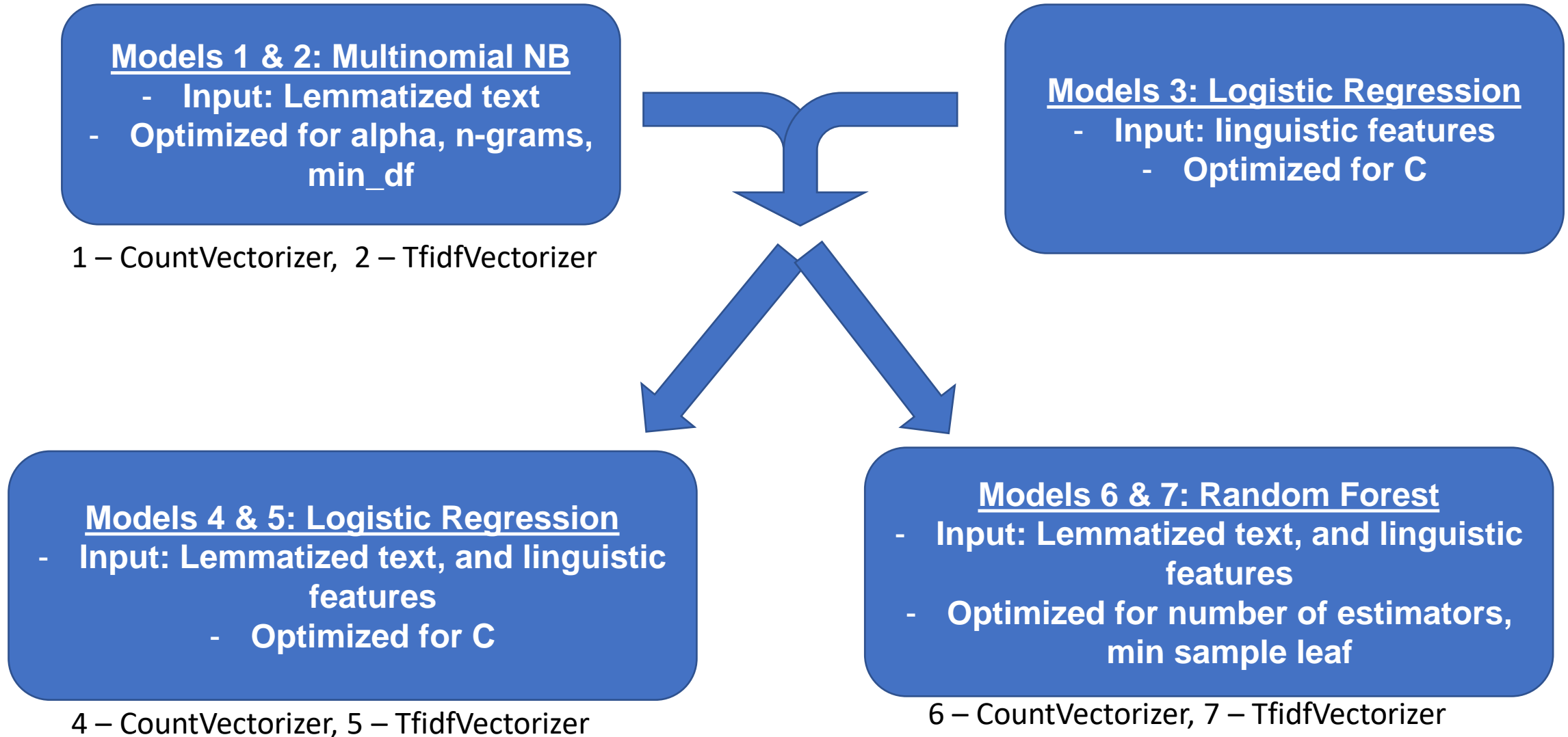
CLI index in the text:

➡ The higher the complexity of the text, the lower the prob the news is fake

Number of sentences in title:

➡ The longer the title, the higher the prob the news is fake

Machine Learning Modeling - Overview



Machine Learning Modeling - Results

	Classifier	Extraction Technique	Linguistic scores & text features	Balance Accuracy	AUC Score
Model 1	MultinomialNB (alpha =0.01)	CountVectorizer (min_df = 50, max_df =0.65, ngram_range = 2,3)	YES	0.9336	0.9922
Model 2	MultinomialNB (alpha =0.01)	TFIDFVectorizer (min_df = 50, max_df =0.65, ngram_range = 2,3)	NO	0.9262	0.9934
Model 3	Logistic Regressor (C=0.1)	None	NO	0.7649	0.9431
Model 4	Logistic Regressor (C=1)	CountVectorizer (min_df = 50, max_df =0.65, ngram_range = 2,3)	YES	0.9748	0.9987
Model 5	Logistic Regressor (C=200)	TFIDFVectorizer (min_df = 50, max_df =0.65, ngram_range = 2,3)	YES	0.9755	0.9992
Model 6	Random Forest (n_estimators = 100, min_samples_leaf = 1)	CountVectorizer (min_df = 50, max_df =0.65, ngram_range = 2,3)	YES	0.9632	0.9978
Model 7	Random Forest (n_estimators = 100, min_samples_leaf = 1)	TFIDFVectorizer (min_df = 50, max_df =0.65, ngram_range = 2,3)	YES	0.9637	0.9981



Model 5 – False Negatives

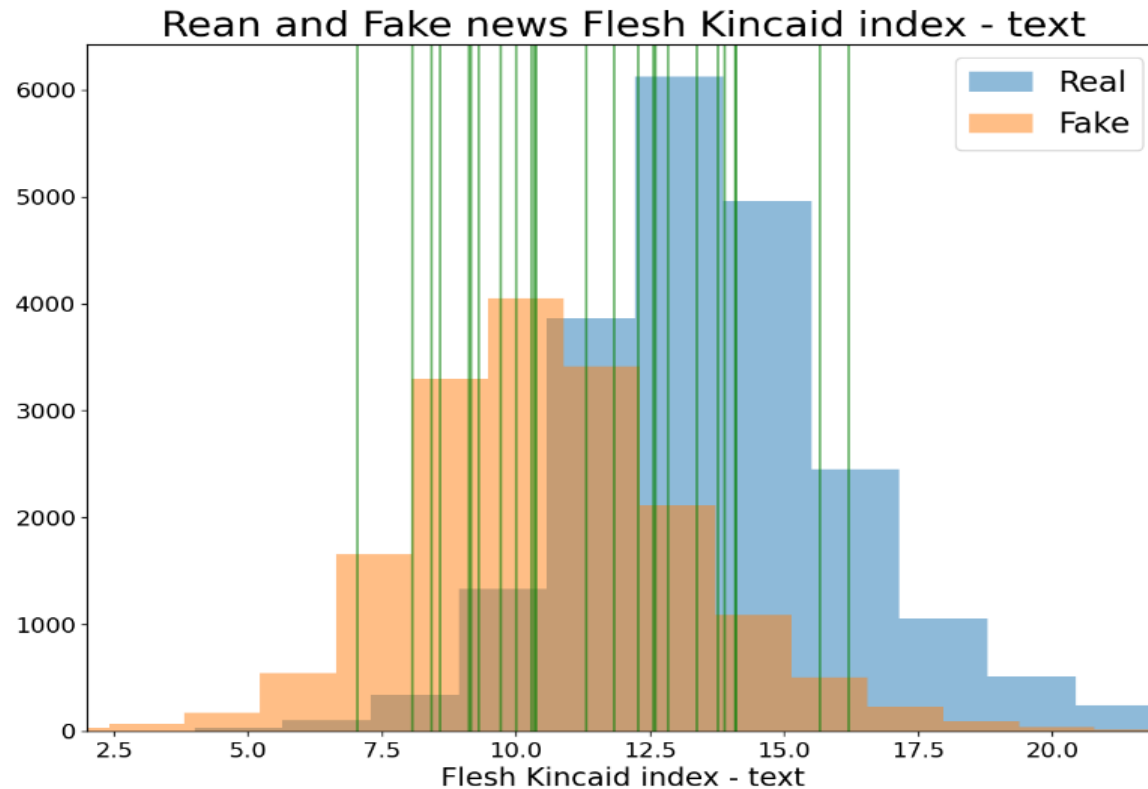
Confusion Matrix

[4186	25]
[65	3440]

Excerpt of fake news classified as real

“...i will still have that same big-shouldered chicago lust for power that drove me from greenwood avenue in hyde park to pennsylvania avenue in three short years. but if hillary replaces me in the oval, she and bill will take control of the democratic party it will become the clinton party once again and they will block me from having any future influence. i will end up like jimmy carter hammering away in appalachia for habitat for humanity. on the other hand, if a republican wins in 16, the clintons will be finished their foundation and their speaking fees will dry up and they will be a thing of the past. but i will still be the titular head of the party. i will be able to continue my push to transform america into a european-style socialist state. so, personally, i would be better off if the next president is a republican...”

Model 5 – False Positives



Real news Flesh Kincaid index mean: **13.8**

Fake news Flesh Kincaid index mean : **10.6**

False positive Flesh Kincaid index mean: **11.39**



Real news with shorter, simpler texts are prone to misclassification

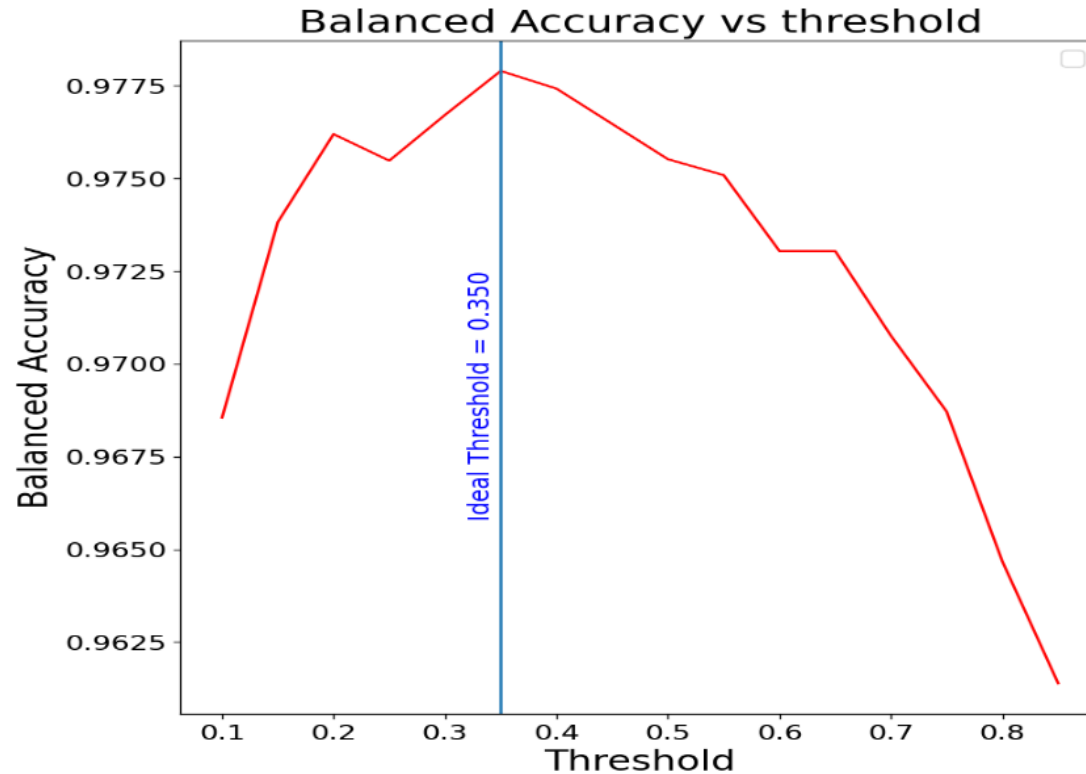
Model 5 – Hyperlinks weight in prediction

Tokens Excluded	Balanced accuracy	AUC
Yes	0.9741	0.999211
No	0.9738	0.999205



Drop of 0.0006%

Model 5 – Probability Threshold Tuning



Threshold=0.5: **Balanced Accuracy 0.9755**



Threshold=0.35: **Balanced Accuracy 0.9779**

Conclusion

- All models are very accurate, **actually too much accurate**: real and fake news were collected from two different sources.
- In future real and fake news articles (as labelled by a third party) could be collected from a similar distribution of sources, or we could pull articles from a few sources and use crowdsourcing to label them as real or fake.
- **Possible improvements of the classifier:**
 - Reduce min_df hyperparameter of the Vectorizer
 - Detection of punctuation patterns
 - Maintaining words semantic
 - Detection of spelling and grammar error