

Real or Fake?

Classifying News Headlines with Natural
Language Processing

Laura Salakari - 17.01.2026

Table of Contents

1. Introduction to Data
2. Preprocessing
3. Models
 - a. Random Forest Classifier
 - b. Logistic Regression
 - c. LinearSVC
4. Final Testing



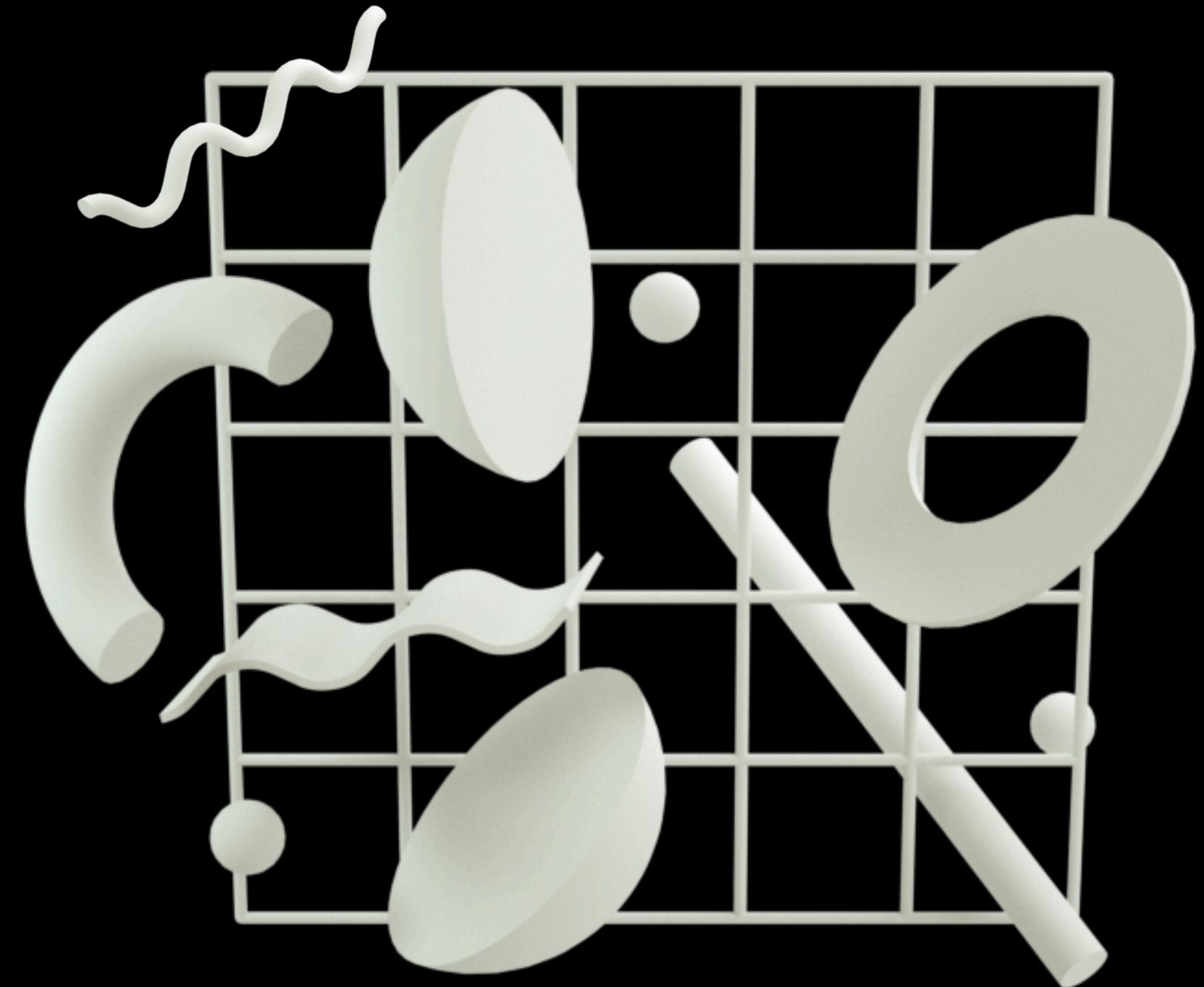
The Data

Training Data

- 34152 entries
- 17572 fake headlines (51.5%)
- 16580 real headlines (48.5%)

Testing Data

- 9984 unlabelled entries





Preprocessing the Data

Preparing the Data for NLP

1. Removing punctuation
2. Lowercasing the headlines
3. Removing stopwords
4. Splitting the training data (80/20)
5. Lemmatization

```
stopwords.py

1 stop_words = set(stopwords.words('english'))
2
3 def remove_stop_words(text):
4     return " ".join(
5         w for w in text.lower().split()
6         if w not in stop_words
7     )
8
9 df_train["headline"] = df_train["headline"].apply(remove_stop_words)
```

Stopword code



```
lemmatization.py

1 import spacy
2
3 nlp = spacy.load("en_core_web_sm", disable=["ner", "parser"])
4
5 def lemmatize(text):
6     if not isinstance(text, str):
7         return ""
8     doc = nlp(text)
9     return " ".join(token.lemma_ for token in doc if token.is_alpha)
10
11 X_train = X_train.apply(lemmatize)
```

Lemmatization code

NLP Models

Tokenizing in Pipelines



Common steps

All models were built with Scikit's Pipeline system and include GridSearches to improve model performance.
All Pipelines begin with Tf-Idf tokenization and vectorization

Random Forest Classifier

1. Tf-Idf with grid search to find best tokenization and vectorization
2. Fitting RandomForest with training data

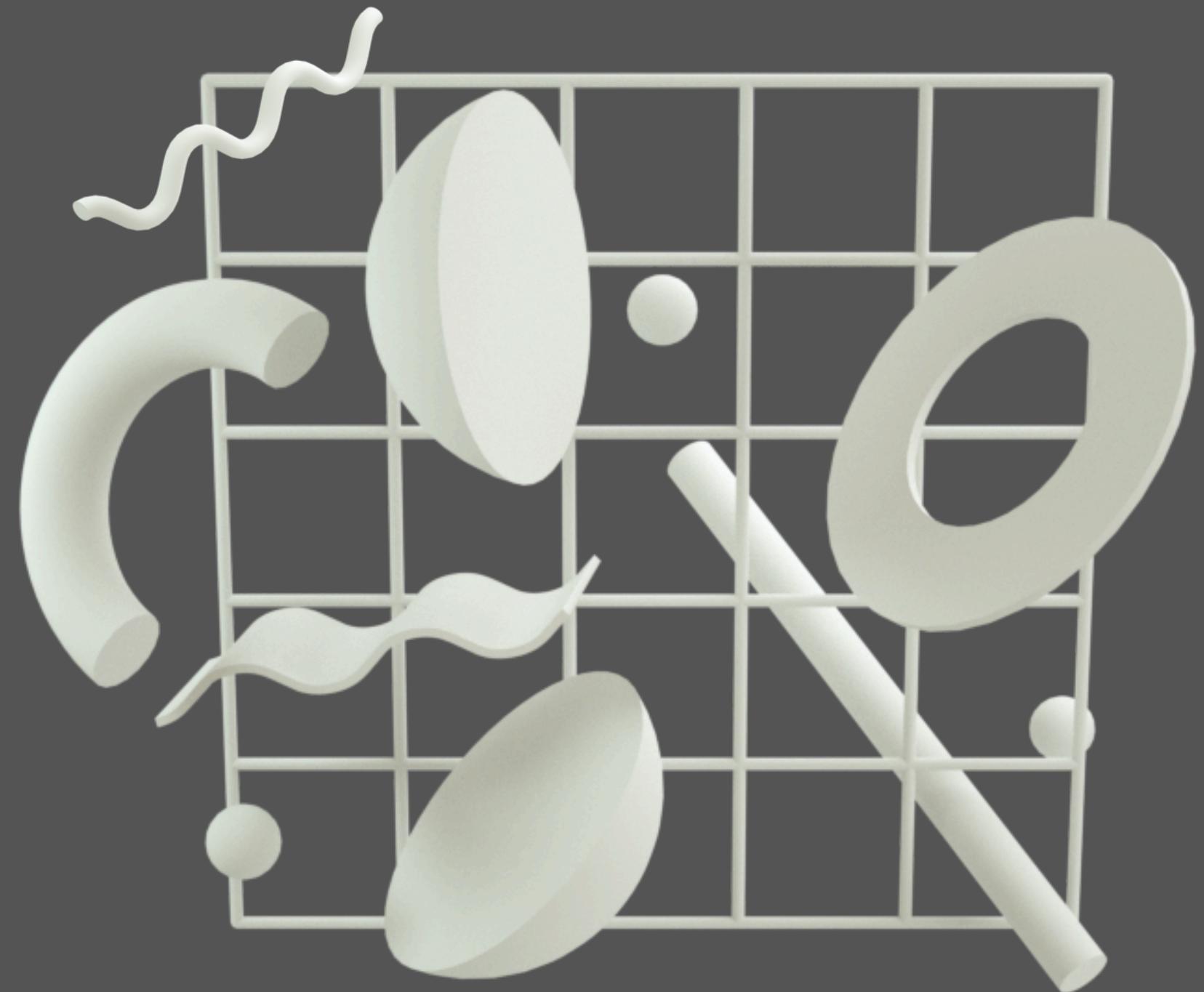
Predictions on the model with the best estimators from the pipeline yield an accuracy of **0.9212**

```
...  
rfc.py  
  
1 rf_pipeline = Pipeline([  
2     ("tfidf", TfidfVectorizer()),  
3     ("rfc", RandomForestClassifier(n_estimators=200, criterion="entropy",  
4         random_state=42))  
5 ])  
6  
6 rf_param_grid = {  
7     "tfidf__ngram_range": [(1,1), (1,2), (1,3)],  
8     "tfidf__min_df": [1,2,5],  
9     "tfidf__max_df": [0.85, 0.9, 0.95],  
10 }  
11  
12 rf_grid = GridSearchCV(rf_pipeline, rf_param_grid, scoring="accuracy", cv=5,  
13     verbose=2)  
13 rf_grid.fit(X_train, y_train)
```

| | | Predicted | |
|--------|------|-----------|------|
| | | Fake | Real |
| Actual | Fake | 3224 | 305 |
| | Real | 233 | 3069 |

Tf-Idf Params

- ngram_range = (1, 3)
- min_df = 2
- max_df = 0.85



Logistic Regression

```
lr.py

1 lr_pipeline = Pipeline([
2     ("tfidf", TfidfVectorizer(max_df=0.85, min_df=2, ngram_range=(1, 3))),
3     ("lr", LogisticRegression(random_state=42, max_iter=500))
4 ])
5
6 lr_params = {
7     "lr__C": [0.1, 1, 10],
8 }
9
10 lr_grid = GridSearchCV(lr_pipeline, lr_params, scoring="accuracy", cv=5, verbose=2)
11 lr_grid.fit(X_train, y_train)
```

1. Tf-Idf with best params from RFC to find best tokenization and vectorization
2. Grid Search for best C value in Logistic Regression (10)
3. Fitting with best params

Predictions on the model with the best estimators from the pipeline yield an accuracy of **0.9413**

| | | Predicted | |
|--------|------|-----------|------|
| | | Fake | Real |
| Actual | Fake | 3311 | 218 |
| | Real | 183 | 3119 |

LinearSVC

```
lsvc.py

1 lsvc_pipeline = Pipeline([
2     ("tfidf", TfidfVectorizer(max_df=0.85, min_df=2, ngram_range=(1, 3))),
3     ("lsvc", LinearSVC(random_state=42, class_weight="balanced"))
4 ])
5
6 lsvc_grid_params = {
7     "lsvc__C": [0.01, 0.1, 1, 10]
8 }
9
10 lsvc_grid = GridSearchCV(lsvc_pipeline, lsvc_grid_params, scoring="accuracy", cv=5,
11                           verbose=2)
12 lsvc_grid.fit(X_train, y_train)
```

1. Tf-Idf with best params from RFC to find best tokenization and vectorization
2. Grid Search for best C value in LinearSVC (1)
3. Fitting with best params

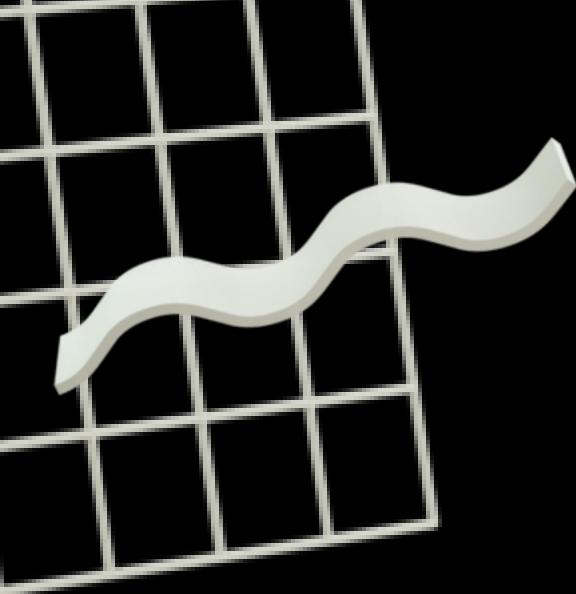
Predictions on the model with the best estimators from the pipeline yield an accuracy of **0.9417**

| | | Predicted | |
|--------|------|-----------|------|
| | | Fake | Real |
| Actual | Fake | 3308 | 221 |
| | Real | 177 | 3125 |

Comparison

| Model | Accuracy |
|---------------------|---------------|
| LinearSVC | <u>0.9417</u> |
| Logistic Regression | 0.9413 |
| Random Forest | 0.9212 |

* LinearSVC fared slightly better with classifying real headlines correctly, while Logistic Regression had an edge with fake headline



Tackling the Testing Data



Preparing for classification



Before proceeding with classification, the testing data needs to be preprocessed the same way as the training data was:

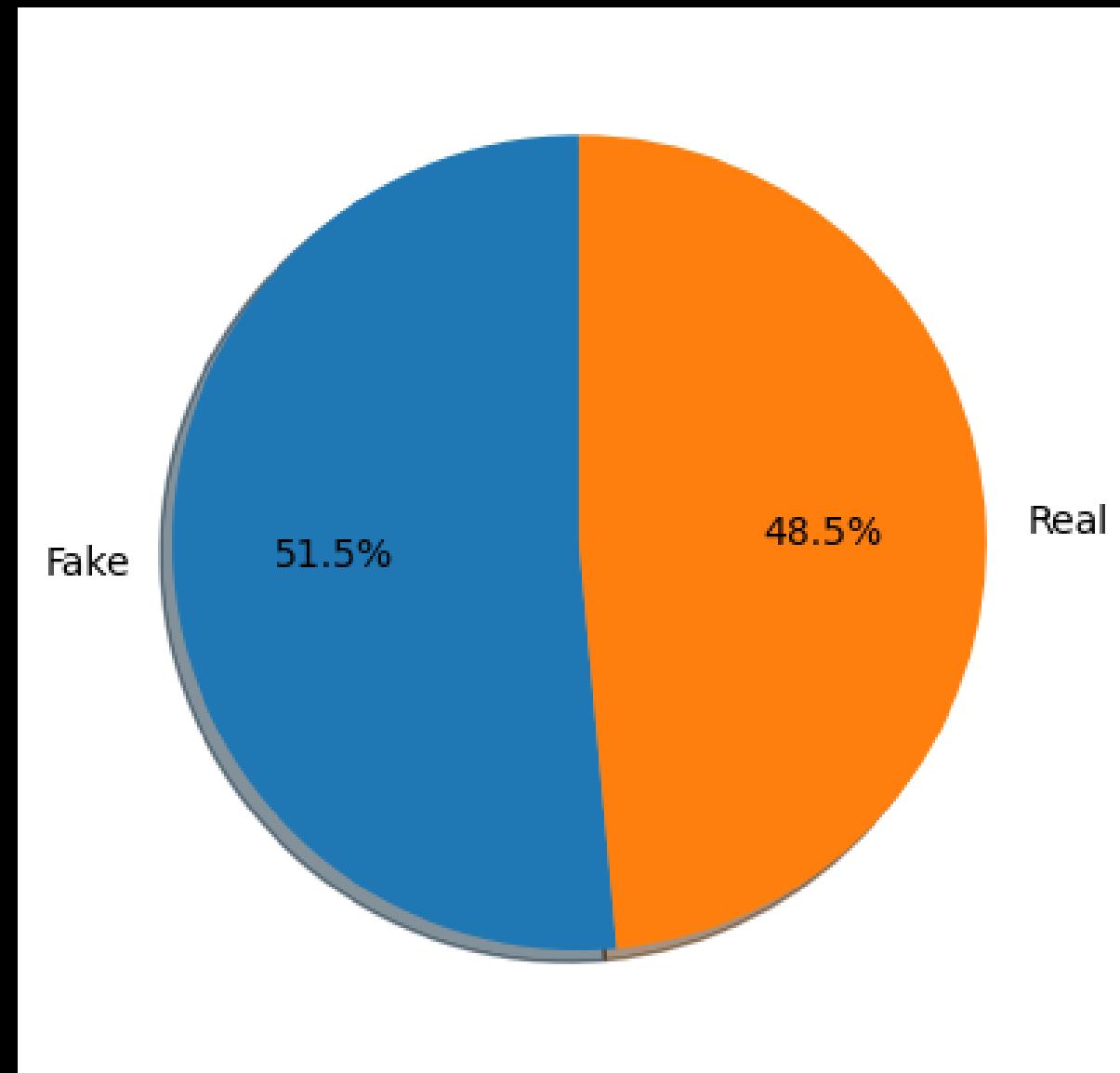
1. Punctuation
2. Lowercasing
3. Stopwords
4. Lemmatization

Classification is conducted with the best model from LinearSVC as it had the highest accuracy

testing_data.csv is overwritten with the new labels accompanying the original, unprocessed headlines

Testing Results

51.5% of the headlines (5138) are classified as fake, leaving 48.5% (4846) as real headlines.



Thank you for
your attention

