

LIVE ●

*Project by:*  
*Diogo Silva*  
*Gabriele*  
*Natasha*  
*Rafael*

# Natural Language Processing Challenge

IronHack Mini Project:  
Fake News Detection Using NLP  
Classifying Real vs Fake News Headlines

## FAKE NEWS





Develop an ML model to classify news headlines as real or fake



Explore TF-IDF and multiple ML algorithms (Logistic Regression, Naive Bayes, SVM, Random Forest, XGBoost)



We compare their performance, select the most accurate model, and use it to label the unseen test headlines.



Tune best performers for maximum accuracy.



Generate final predictions for the testing dataset.

***IronHack***  
***Mini***  
***Project***

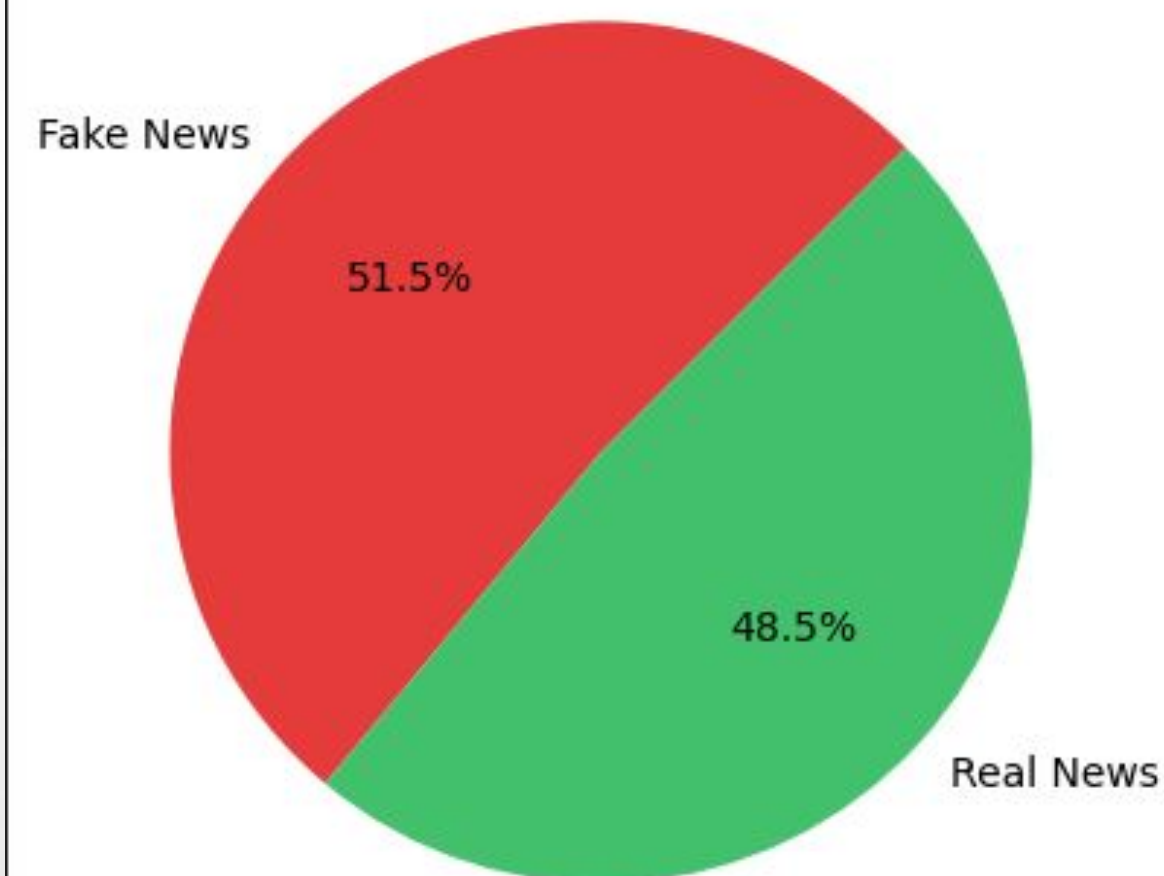
# PROJECT GOAL

## Dataset Sample

	fakeornot	headline
866	0	anderson cooper throws major shade at white ho...
4943	0	florida republican has an insane meltdown afte...
6532	0	trump promises to destroy the usa during inter...
16041	0	in his own words: stunning unofficial trump ad...
16903	0	the single chart to share that tells the truth...
32847	1	turkish nationalist opposition seeks to secure...
16578	0	so god made a patriot: ,i need a man who knows...
30193	1	mexican opposition leader anaya to seek presid...
3761	0	donald trump openly admits to america he has n...
27178	1	clinton loses to sanders in coal state of west...

- **Dataset Shape - 34.152 rows x 2 columns**
- **training\_data.csv (headline + label)**
- **testing\_data.csv (headline only)**
- **Train/validation split: 80/20**
- **Labels ("Fakeornot"): 0 = real, 1 = fake**

Fake vs Real News Distribution



# Dataset



## Text Preprocessing

- Lowercasing
- Removing punctuation
- Stripping whitespace
- TF-IDF handles tokenization + term importance

## TF-IDF Vectorization

### Parameters used:

- n-grams: 1-2
- max\_df: 0.9
- min\_df: 5
- stop words: English

Converts text into numerical features representing word relevance.

## Baseline Model

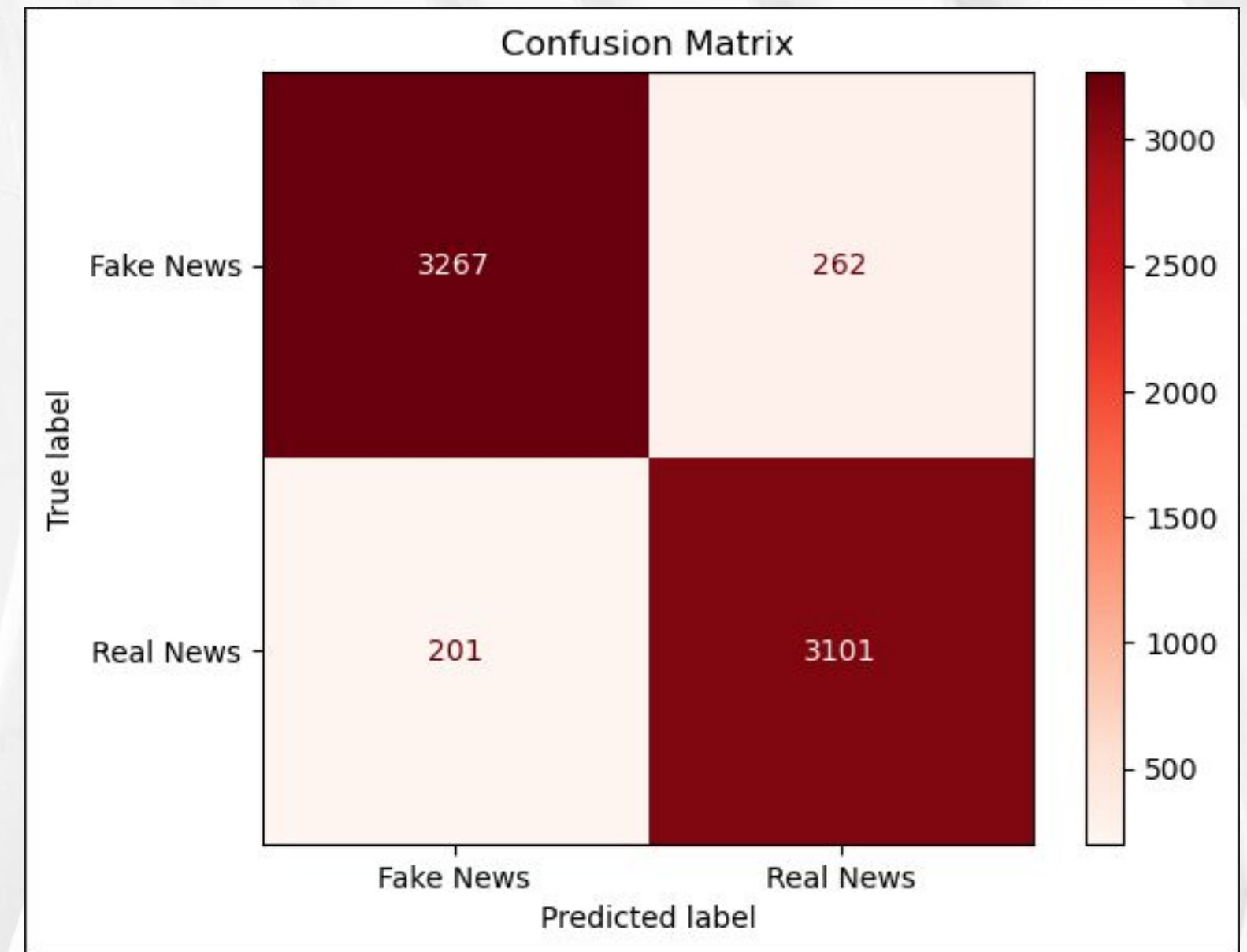
### Logistic Regression

Metrics Used: Accuracy, Precision, Recall, F1-Score

## Baseline Model: Logistic Regression



Metric	Value
Accuracy	93,12%
Precision	91,92%
Recall	94,03%
F1-Score	92,96%



# Text Preprocessing & Baseline Model Development

LOGISTIC REGRESSION 93,1% ▲

RANDOM FOREST 91,1% ▼

LINEAR SVC 93,2% ▲

XGBOOST 90,5% ▼

NAIVE BAYES 92,8% ▲





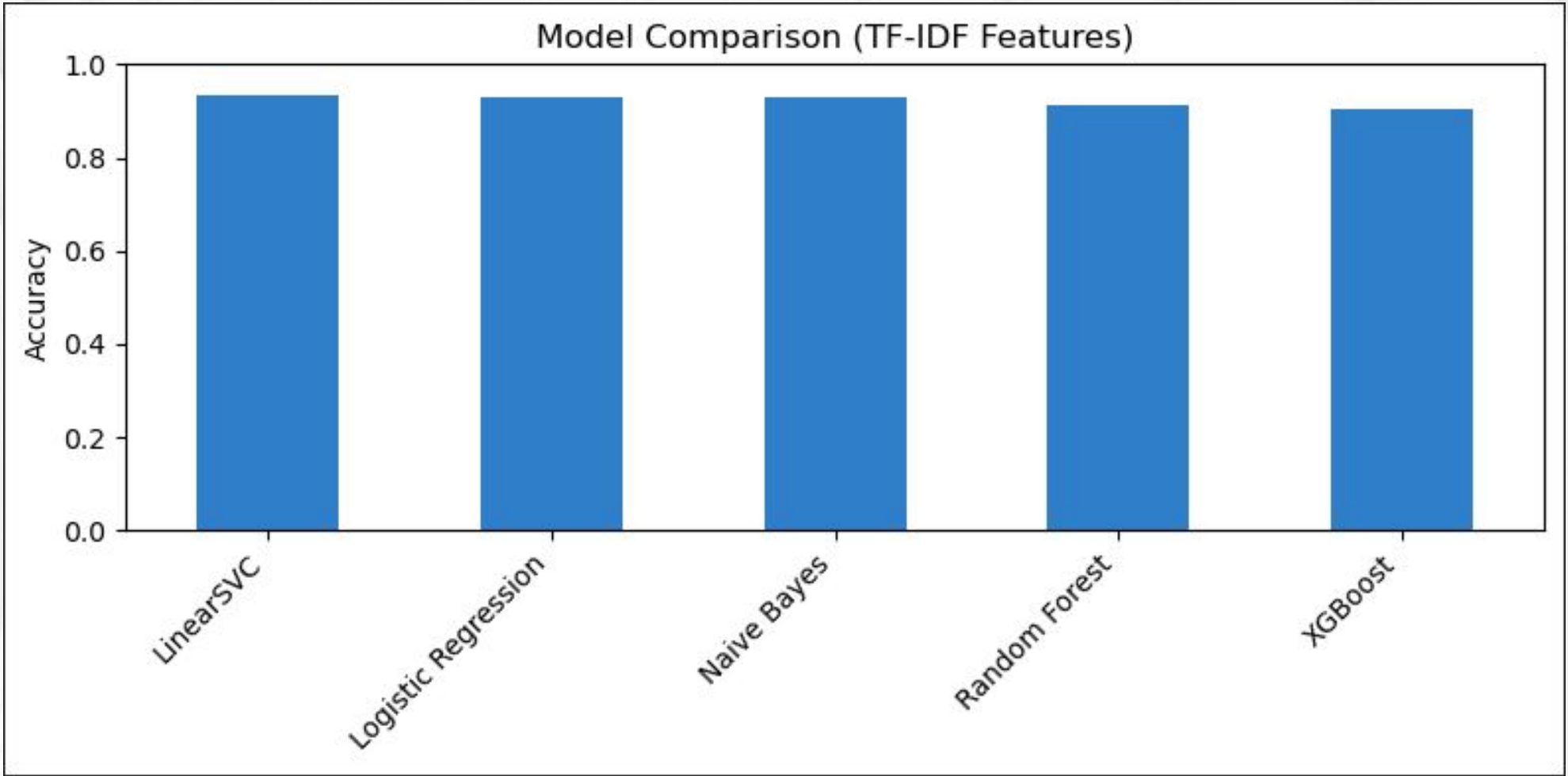
# Trying Different Classifiers

- Logistic Regression
- LinearSVC
- Naive Bayes (MultinomialNB)
- Random Forest
- XGBoost

## Initial Results

Logistic Regression and LinearSVC showed strongest baseline performance.

Model	Accuracy Evaluation
LinearSVC	93.4%
Logistic Regression	93.2%
Naive Bayes	92.3%
Random Forest	91.6%
XGBoost	90.3%



# Models Evaluation & Initial Results

***Breaking  
News***

**LIVE** 



***WAIT, LET'S TUNE OUR  
MODELS***

LOGISTIC REGRESSION 93,1% 

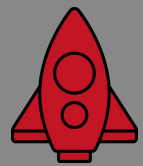
RANDOM FOREST 91,1% 

LINEAR SVC 93,2% 

XGBOOST 90,5% 

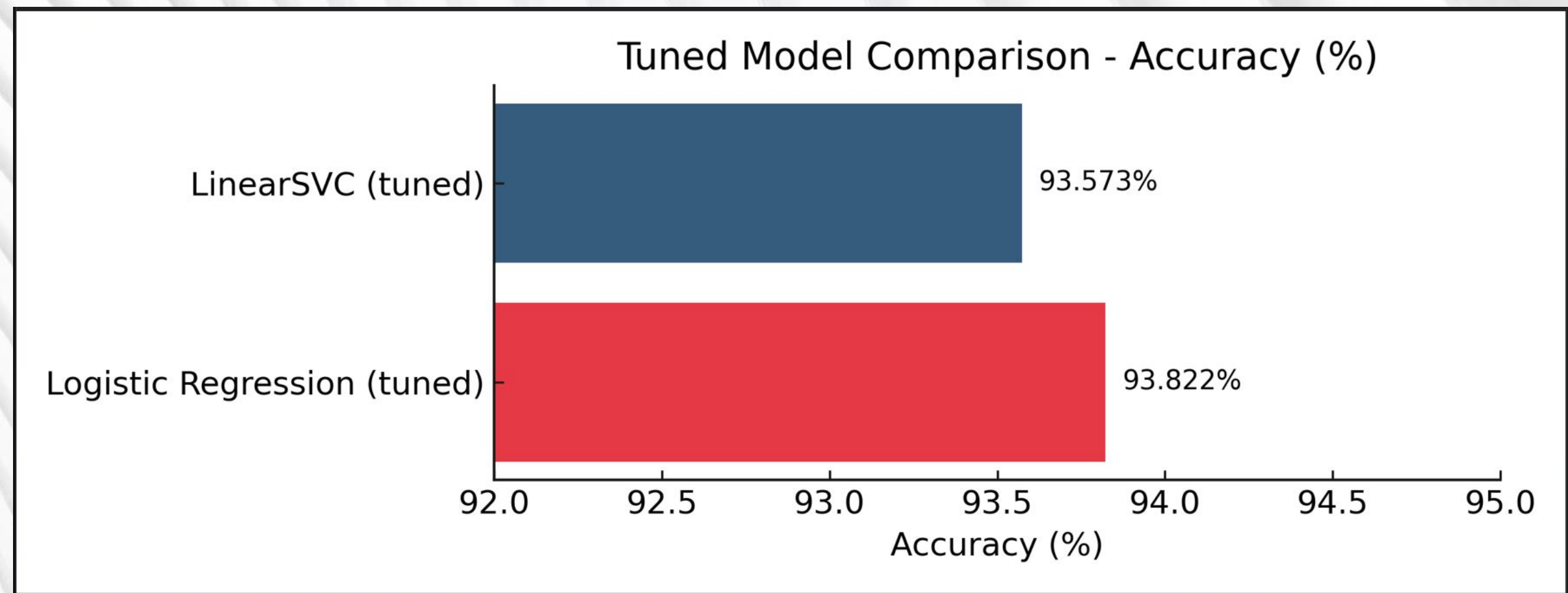
NAIVE BAYES 92,8% 





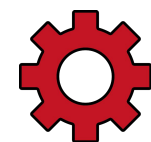
## Tuning the Best Performance Models

- Used GridSearchCV
- 3-fold cross-validation
- Tuned: Logistic Regression and LinearSVC
- Improved accuracy beyond baselines



*2 Best Performing Models*

# Hyperparameter Tuning





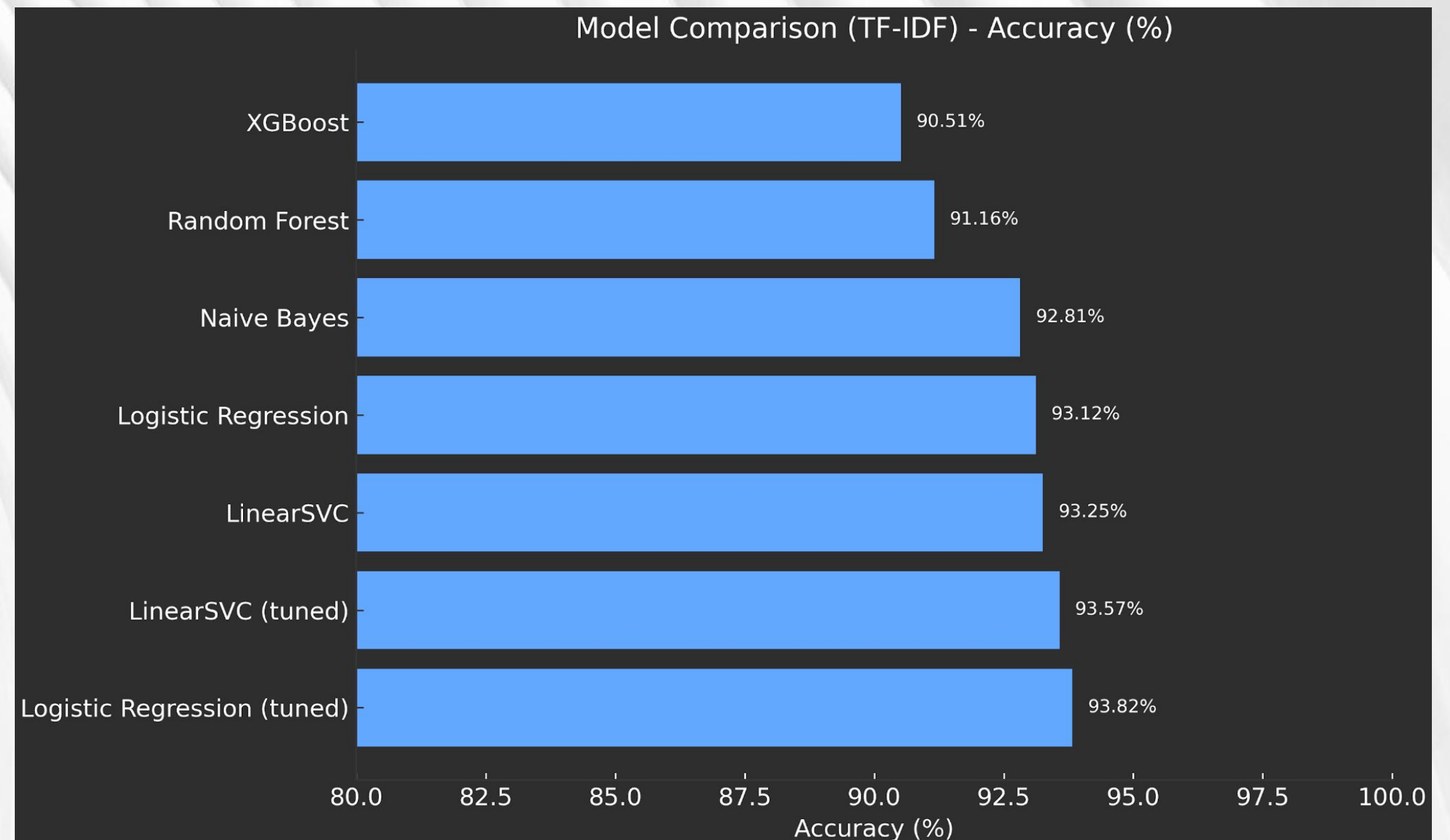
## Overall Evaluation

Comparison of accuracy across all evaluated TF-IDF models.

Linear models (Logistic Regression and LinearSVC) outperform tree-based models.

The **tuned Logistic Regression achieves the highest accuracy (93.82%)**, emerging as the top performer.

Hyperparameter tuning provided consistent improvements, especially for linear models.



## Final Model Comparison

### NLP MINI PROJECT

LOGISTIC REGRESSION 93,8% ▲

RANDOM FOREST 91,6% ▲

LINEAR SVC 93,4% ▲

XGBOOST 90,3% ▼

NAIVE BAYES 92,9% ▲



*Logistic Regression  
(Tuned)*



**93,8%**

Accuracy

*LinearSVC  
(Tuned)*

**93,6%**

Accuracy

*LinearSVC*

**93,4%**

Accuracy

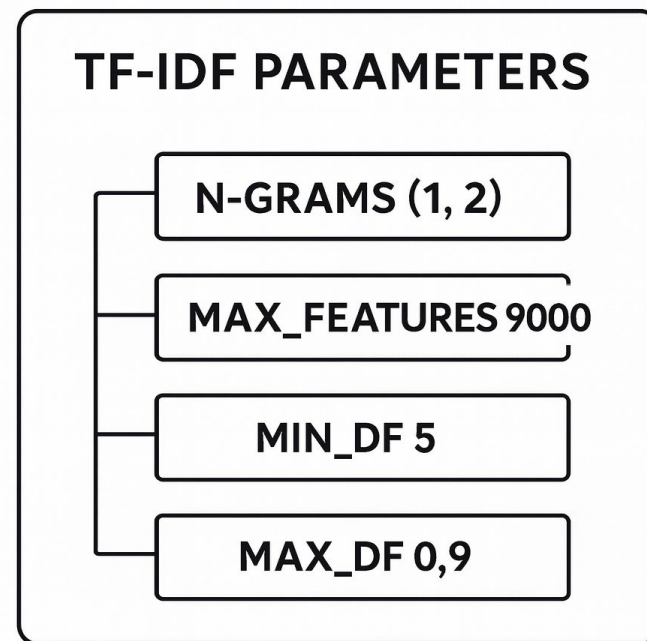
***Best Overall***









## Training Final Model

-  TF-IDF retrained on full dataset



-  Final Logistic Regression (tuned) trained on all labeled data
-  Saved vectorizer + model for prediction

## Predictions for Test File

-  Applied full TF-IDF: transform it with the tfidf\_full vectorizer, predict,
-  Predicted labels 0/1: overwrite the first column (which originally contains 2s) with 0/1.
-  Updated first column accordingly
-  Saved as testing\_predicted.csv



# FINAL MODEL TRAINED



“

# Conclusion

LIVE ●

**S** Classical models can outperform deep learning on short text

- TF-IDF was highly effective
- Logistic Regression achieved 93.8 percent accuracy
- Reliable approach for fake news classification

— Next Steps for future improvement:

- Fine-tune transformer models
- Add character-level features
- Combine TF-IDF with embeddings
- Deploy as API or dashboard

”

***Thank You!***