

Clickbait Detection

2025 NLP Course Project

Jakub Sawicki, Jędrzej Sokołowski, Wiktor Woźniak

What is Clickbait?

- Attention-grabbing headlines, often omitting key facts
- Sparks curiosity and encourages clicks using sensational language
- Difficult to detect automatically: can be subtle and context-dependent
- Exists on a spectrum, not just binary classification

Motivation & Challenges

- Manipulative clickbait undermines trust in news and social media
- Hard to draw the line between catchy and misleading
- Even humans disagree about what counts as clickbait
- Growing sophistication with AI-generated clickbait

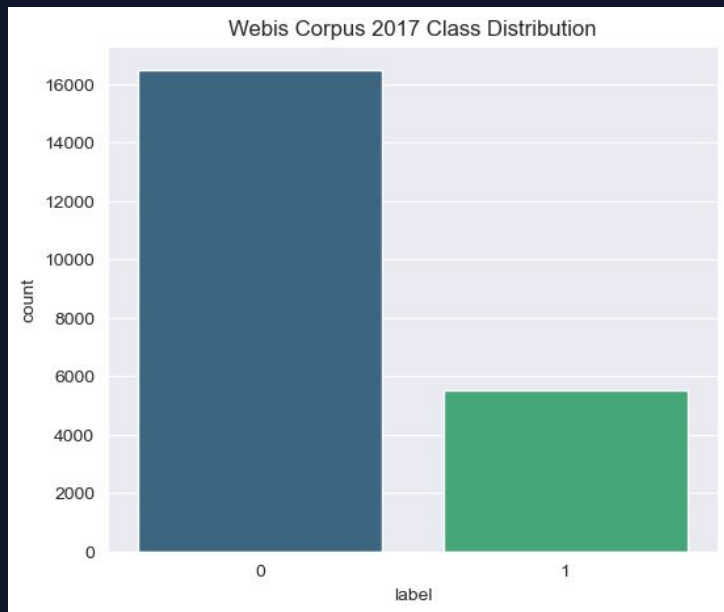
Related works, SOTA Approaches

- Early: Hand-crafted features (punctuation, length, headline-article overlap)
- Classical ML: SVM, Random Forests, high accuracy with engineered features
- Deep Learning: CNNs, RNNs using word embeddings (Word2Vec, GloVe)
- Transformers: BERT, RoBERTa, current state-of-the-art for both detection and “spoiling”

Key Open Datasets

- Webis Clickbait Corpus 2017 (Twitter, 38,517 posts, 9,276 clickbaits)
- Webis Clickbait Spoiling Corpus 2022 (Twitter, Reddit, Facebook, 5,000 posts, all clickbaits)
- Wikinews Clickbait Corpus (crowdsourced English news headlines)
- Thai & Chinese Headline Datasets (expands multilingual, cross-domain evaluation)

Merging datasets

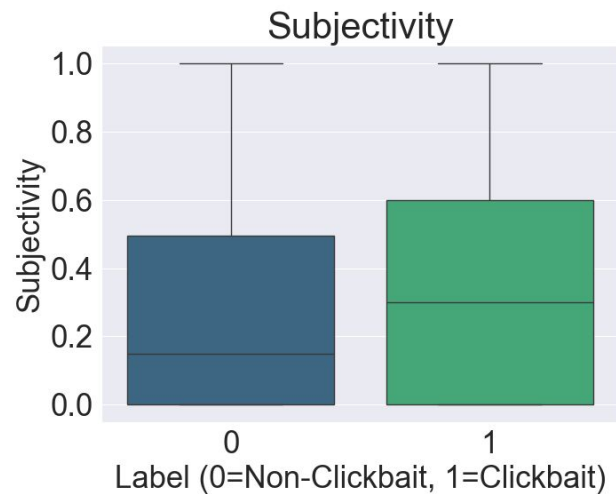
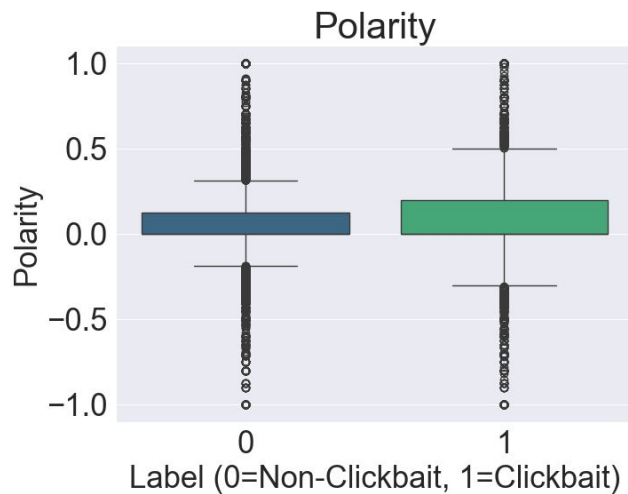
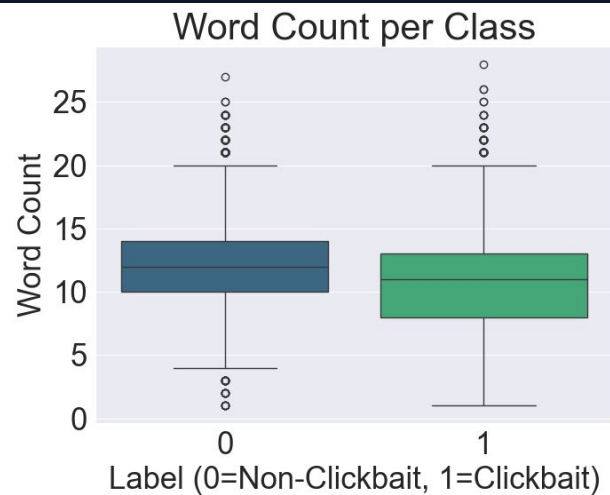
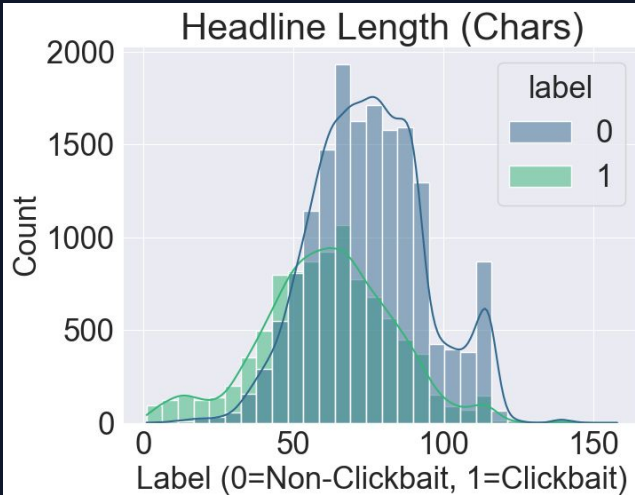


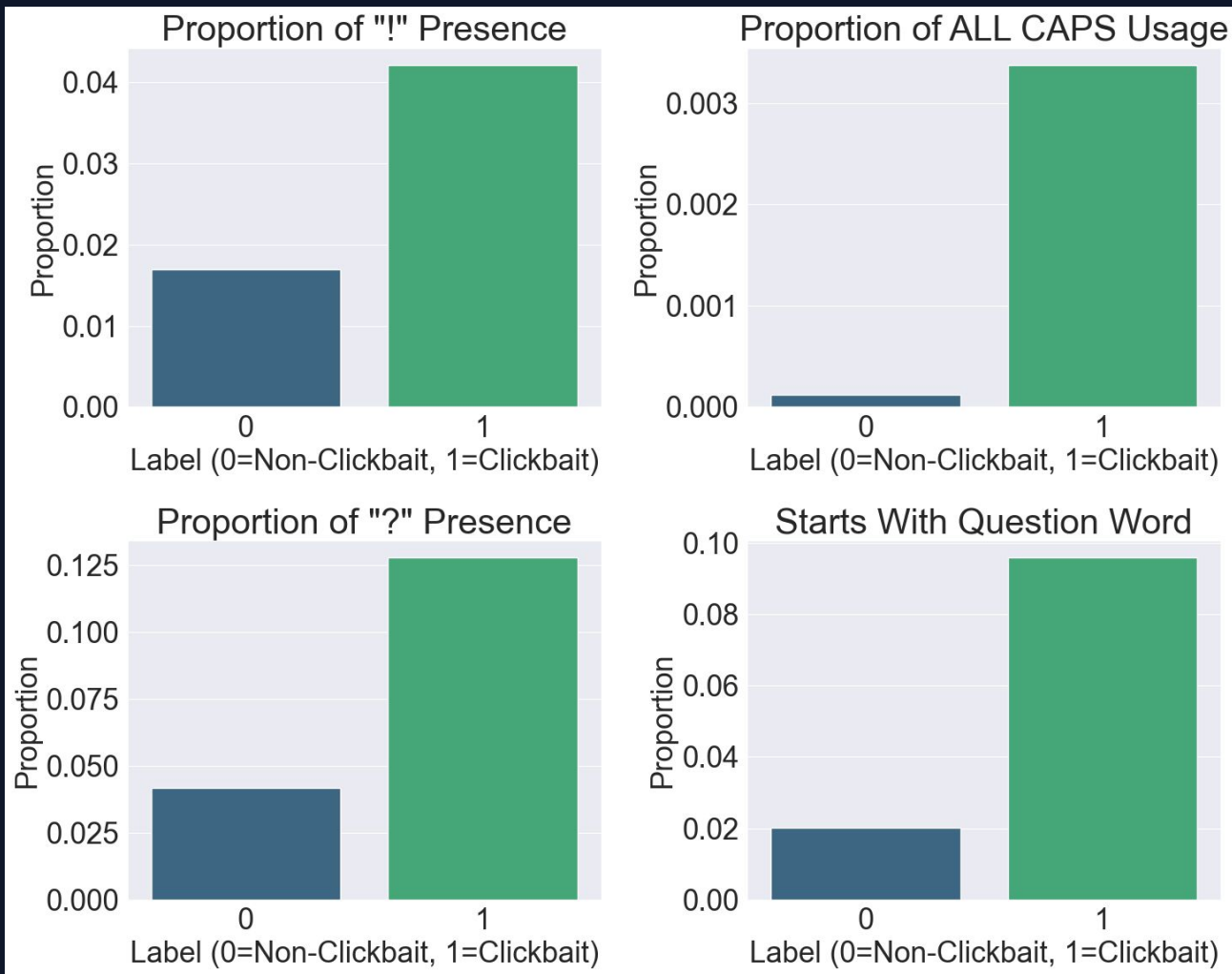
Project Goals & Methods

- Build interpretable, effective clickbait detection system using open benchmarks
- Compare classical ML, neural, and transformer models for detection
- Evaluate Precision, Recall, F1-score

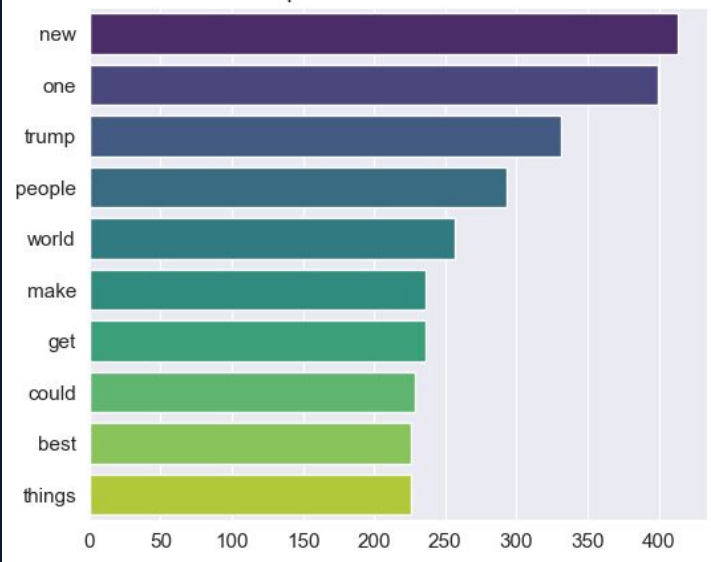
Data Processing

- Remove empty rows
- Train/Validation/Test split - 0.7/0.15/0.15
- Training dataset undersampling
- Add new features:
 - Character count
 - Word count
 - Exclamation marks (!)
 - Question marks (?)
 - ALL CAPS presence
 - Sentiment polarity
 - Sentiment subjectivity
 - Question word markers

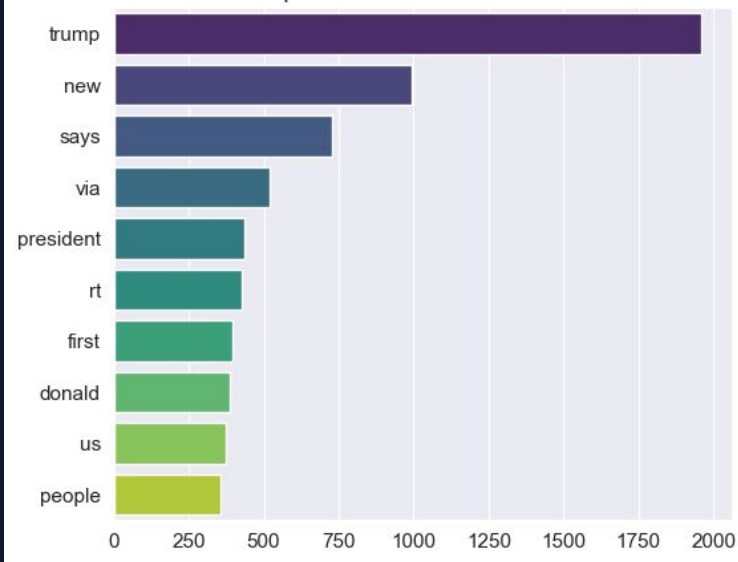




Top 10 Words in Clickbait



Top 10 Words in Non-Clickbait



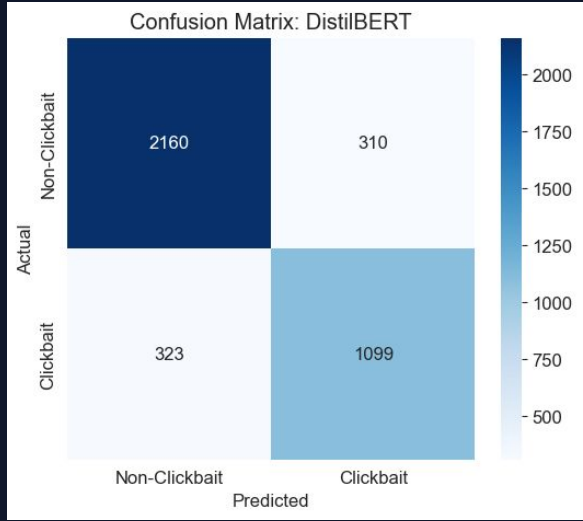
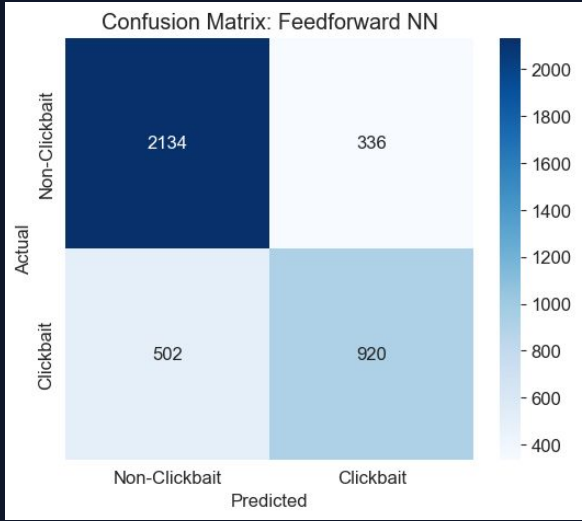
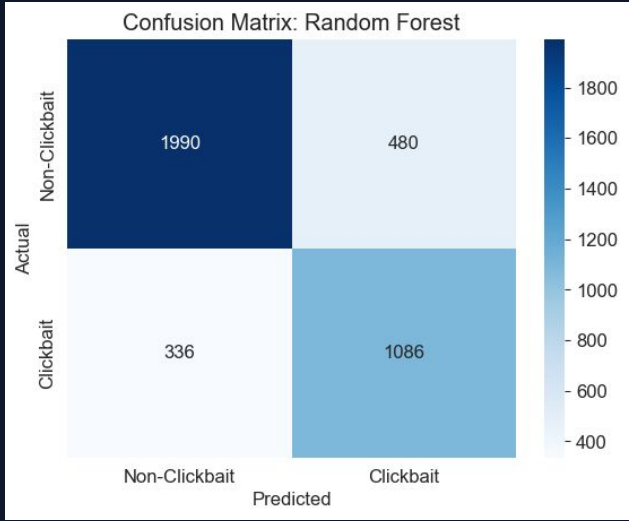
Feature & Model Summary

- Classical Baseline: Random Forest (TF-IDF)
- Feedforward Neural Network (3 hidden layers with ReLU)
- Transformer: DistilBERT fine-tuned for headline classification
- Metrics: Precision, Recall, F1-score

Random Forest Model Performance (Test Set)				
Class	Precision	Recall	F1-score	Support
0 (Non-Clickbait)	0.858 ± 0.003	0.821 ± 0.009	0.839 ± 0.005	2470
1 (Clickbait)	0.711 ± 0.010	0.765 ± 0.005	0.737 ± 0.006	1422

Feed-Forward Neural Network Performance (Test Set)				
Class	Precision	Recall	F1-score	Support
0 (Non-Clickbait)	0.827 ± 0.015	0.838 ± 0.023	0.832 ± 0.005	2470
1 (Clickbait)	0.713 ± 0.018	0.693 ± 0.039	0.702 ± 0.012	1422

DistilBERT Model Performance (Test Set)				
Class	Precision	Recall	F1-score	Support
0 (Non-Clickbait)	0.881 ± 0.014	0.861 ± 0.031	0.871 ± 0.011	2470
1 (Clickbait)	0.770 ± 0.032	0.797 ± 0.032	0.782 ± 0.008	1422



Clickbaitness measure

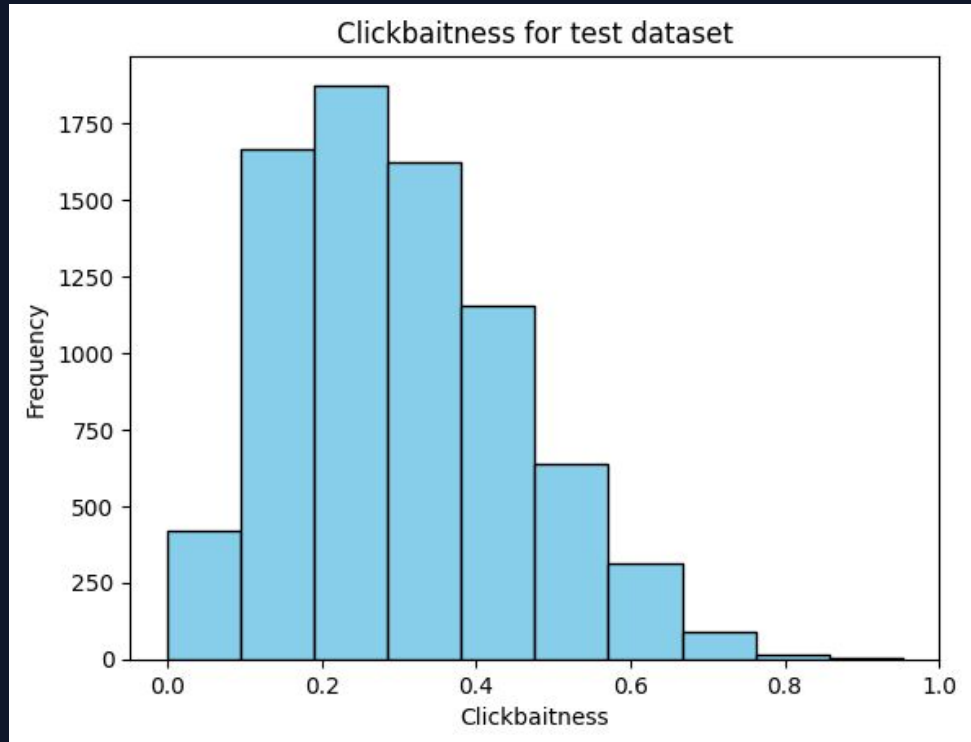
$$\text{Baitness} = (\text{EC} + \text{C} + \text{S} + \text{EOT}) / 4$$

EC = Eye catchness

C = Curiosity

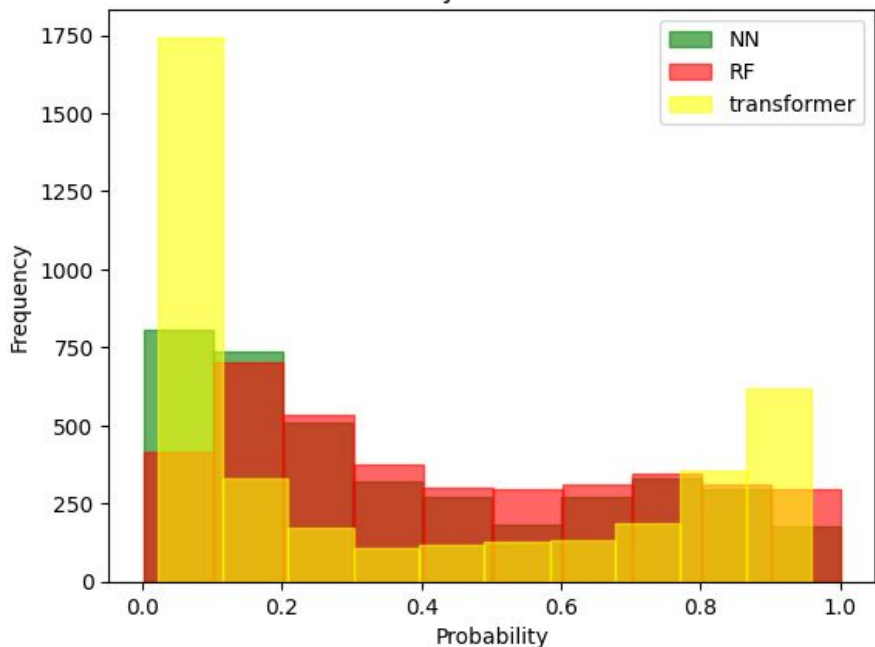
S = Sentiment

EOT = Ease of text

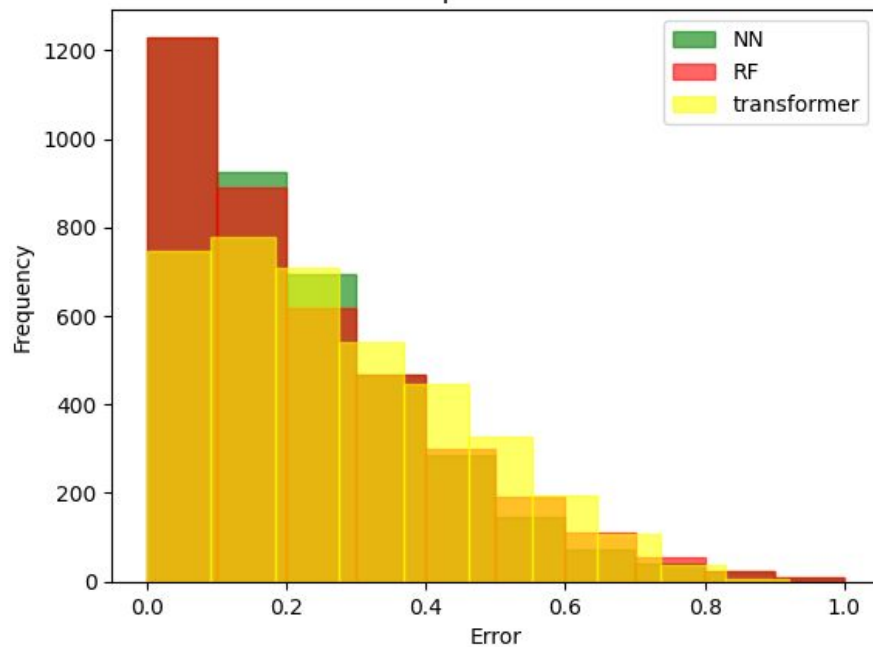


Clickbaitness measure

Probability for test dataset



Absolute error of probabilities and baitness



Clickbaitness measure

Model	MAE	MSE
Neural Network	0.217 ± 0.174	0.077 ± 0.118
Random Forest	0.226 ± 0.188	0.086 ± 0.130
Transformer	0.269 ± 0.184	0.106 ± 0.127

Summary

- DistilBERT was a superior model
- FF-NN and RF approaches achieved comparable performance
- Clickbait probabilities are not a good method for predicting baitness of the
post

Questions?