

# EE5183 FINTECH FINAL PROJECT REPORT

*Tadeo Hepperle, Student ID: A11922105*

NTU National Taiwan University

## ABSTRACT

Amazon reviews are important for purchase decisions and their integrity is essential to the online platforms reputation. Because of their importance, immoral businesses have an incentive to publish fake reviews. In this paper we explore different machine learning models to detect fake reviews based on a combination of textual and numerical features. Our models use pretrained BERT embeddings as a base and achieve AUROC and AUPRC scores of up to 0.933 and 0.970 respectively. As a main innovation we allow our models to take a look at contextual reviews for the same Amazon product to better understand if a review is fake or not.

**Index Terms**— BERT, NLP, fraud detection, neural networks, sentiment analysis

## 1. INTRODUCTION

In this project we take a look at a dataset consisting of Amazon reviews that were either real or fake. Real meaning that the review was genuine and reflects the customers "real" opinion. Fake meaning that the review is probably forged with other motivations in mind, for example to boost or damage a products' reputation. Because the majority of online shop customers rely on reviews when making their purchase decisions it is important for a platform like Amazon to ensure the reviews reflect the real qualities of a product and are not misused as a marketing instrument by unethical businesses. To detect fake reviews we propose two different neural network models that both make heavy use of a BERT architecture to generate text embeddings that are useful in the binary classification. To see how our best model performs in action you can check out this little online game, where you can compete against our model yourself in real-time.

### 1.1. Previous Work

Previous work on fake review detection was done by other researchers however most focussed on either the use of textual or non-textual features. For examples Gu and Budhkar (2021) used transformer based models in a way that they transformed all features into text, reducing the problem domain to a mere binary text classification. Mohawesh et al.

(2021) used transformers combined in ensemble methods to classify texts. They also take the context of a review into consideration. Weng et al. (2022) also detected Chinese fake reviews using pre-trained language models. Combining these approaches is one of the aims of our project.

## 2. MATERIAL AND METHODS

### 2.1. Data and Preprocessing

We used a dataset uploaded by Github user aayush210789 for our report, published here on Github. The data is from a corpus of reviews published by amazon and was internally labeled as fake or real. That means the labeling was done by Amazon and we sadly do not have advanced insight into the proceedings how exactly the labels were created. We do know however that the user behavior (for example if one user writes 100 reviews in one day) and the time of writing the review (e.g. 100 days after the product was bought) played a big role in the labeling. The dataset consists of a csv file with the columns shown in table 1.

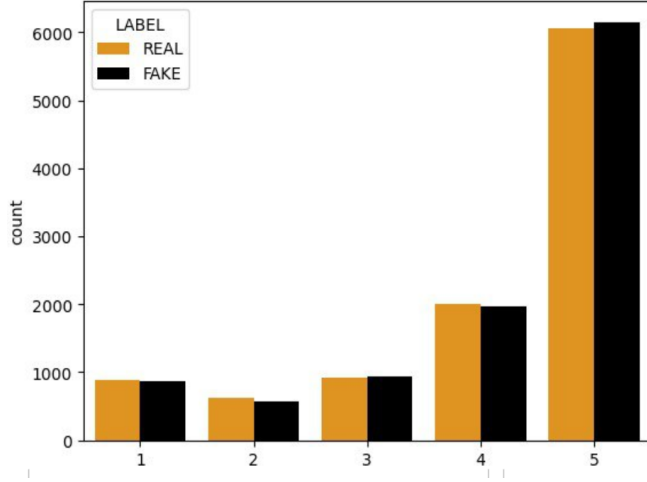
**Table 1:** Original features

Feature	Description
review_title	title of the review
review_text	content/text of the review, mostly moderately long text
rating	1-5 star rating for the product
verified	1 if there was a verified purchase before the review was submitted, else 0
label	0 for a real review, 1 for a fake review
product_id	every Amazon product has a unique alphanumeric id
product_title	title of the product
category	the category of the product, in total 30 categories were present, such as PC, Baby, Books, ...

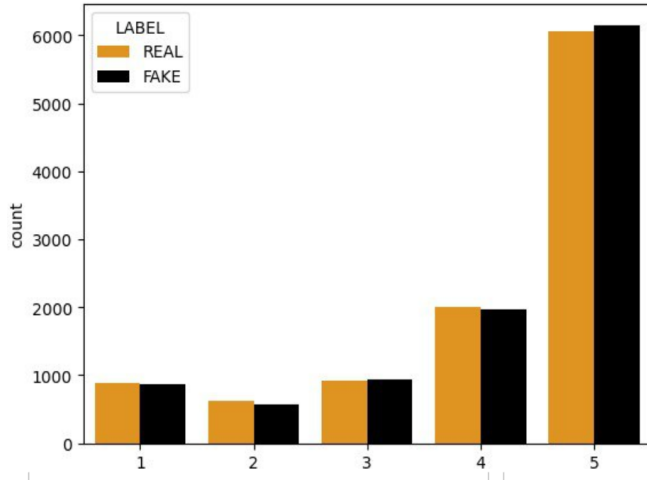
The CSV file contained exactly 10,500 fake and 10,500 real reviews such that imbalance was not an issue. The reviews encompassed 18857 Amazon products with unique IDs.

As Fig 1 shows, fake reviews tend to be less frequently verified purchases than real reviews. In contrast, there does

not be a relationship between real vs. fake reviews and the rating. In all brackets the distribution is about equal, as Fig 2 shows.



**Fig. 1: 3-BERT ROC**



**Fig. 2: 3-BERT ROC**

### 2.1.1. Scraping additional Data

We discovered that most of the products were still listed on Amazon under the same product\_id. So we wrote a script that made 18857 requests to the corresponding Amazon websites, that are available under "www.amazon.com/dp/product\_id". For more than 86% of the products (16,309 products) the requests went through, and we were able to save the entire HTML content. With these 26.4 GB of data we were able to extract much more information. For each of the products we could extract from the HTML the features from in table 2.

In total we were able to scrape 99,995 additional reviews from the product front pages. These reviews can be helpful

**Table 2: Scraped features**

Feature	Description
rating_count	how many ratings a product had in total
rating_avg	the average rating on a continuous scale from 1.0 to 5.0
rating1	percentage of ratings that were 1 star
rating...	percentage of ratings that were ... stars
rating5	percentage of ratings that were 5 stars
reviews	array of reviews on the front page of the product, with features displayed in table 3

**Table 3: Features of each Scraped Review**

Feature	Description
review_title	title of the review
review_text	content/text of the review, mostly moderately long text
rating	1-5 star rating for the product
verified	1 if there was a verified purchase before the review was submitted, else 0
helpfulness	how many people found this review helpful

in fake review detection because they help us to see a certain suspicious review in context: If a review is overly positive and has 5 stars it might be fake, but if we see that the most helpful reviews for this product are also all very positive it might be more likely that the review in question is in fact genuine.

Using only products, for which we were able to scrape additional data, 18,107 labeled reviews remained. There was a slight bias towards real reviews (52% real vs 48% fake), but that is not a big deal.

### 2.1.2. Splitting the Data

The data was randomly split into 3 parts: 70% of the reviews (12674) for the train set, 20% (3641 reviews) for the validation set, 10% (1792 reviews) for the test set.

## 2.2. Methods

We propose two different models for the detection of fake reviews, the first one is called 3-BERT the second one is called Context Embedding Model. The 3-BERT model ignored the scraped data and just used the columns from the initial dataset. The Context Embedding Model tried to use the context information of the review (e.g. what do other reviews say about the product) to produce better predictions. In theory the Context Embedding Model should be more powerful and achieve better performance, because it has access to a superset of the data that the 3-BERT model has access to. We use the Binary Cross Entropy Loss to train our models with gradient descent utilizing the Adam optimizer.

### 2.2.1. 3-BERT Model

The 3-BERT model architecture can be seen in 7. It uses a pretrained bert model from huggingface, that is fine tuned to our task during the training process. This bert model converts the review text, review title and product title from the original features into a 768 dimensional embedding each. This is done by running the bert model over the sequence, which encodes each word as a 768-dimensional embedding. This is done by taking the embedding of the [CLS] token, that is inserted before the text tokens itself, as the text embedding. The 30 categories and 5 possible rating values were one-hot encoded and form together with the verified-flag a 36-dimensional vector. This is fed through a 6 layer deep fully connected neural network with dropout together with the BERT embedding of the product and review title. The result is joined with the review text embedding and again put through a fully connected neural network.

### 2.2.2. Context Embedding Model

The context embedding model uses all features available from the tables presented above. A central idea of this model is to create not only text embeddings but consistent embeddings for an entire review. This "review embedding network" could then be used with the same weights on reviews in question (fake vs. real) and scraped reviews because they share the same features. Fig 9 shows the architecture of this: BERT is used to get 768-dimensional embeddings for the review title and review text each. The sentiment analysis library "textblob" is used to form features "valence" and "subjectivity" for both title and text. Then all these vectors are concatenated with "rating" and "verified". This big vector is then fed through an MLP layer to form the reviews embedding. This review embedding network is only one part of the context embedding model though. It shares weights in all instances that is occurs in our model. To get a joined representation of the context we use a transformer encoder layer to feed the "review embedding" of each scraped review through it. This gives us a context embedding, that is then joined with all other product features to get a product embedding, see Fig 8. This in turn is then then joined with the original reviews embedding for this product and put through a binary classification MLP to predict the final label: fake vs. real.

## 3. RESULTS

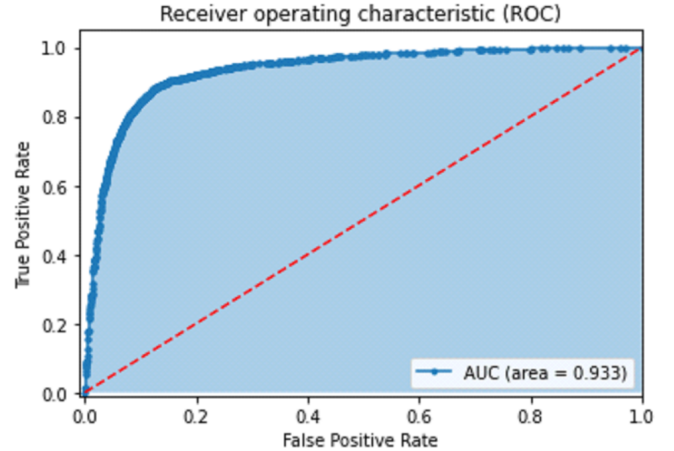
All of our models were trained with the Adam optimizer and used the huggingface model "distilbert-base-uncased" as a base. The 3-BERT model used a learning rate of 0.000023 and 32 epochs with a batch size of 4. The total training time was 16.8 hours on a GPU-accelerated machine. The context embedding model (CEM) had a learning rate of 0.001 and trained for 10 EPOCHS with a batch size of 256 which took 7 hours. Performance differences could be due to different

training times, but our computing and time resources were sadly not infinite. We compared the models to a few simpler models we trained: Unimodal: a model that transforms all features to text features and use only one BERT for the predictions. MLP + BERT: uses MLP on categorical data and BERT for text. BERT + NN: Fine-tuned BERT (reviews only) + simple 5 layer neural network on categorical features. Table 4 shows the performance of all models in terms of AUROC and AUPRC values on unseen test data.

**Table 4:** Features of each Scraped Review

Model	AUROC	AUPRC
Unimodal	0.620	0.610
MLP + BERT	0.812	0.820
BERT + NN	0.839	0.819
3-BERT	0.933	0.925
CEM	0.731	0.970

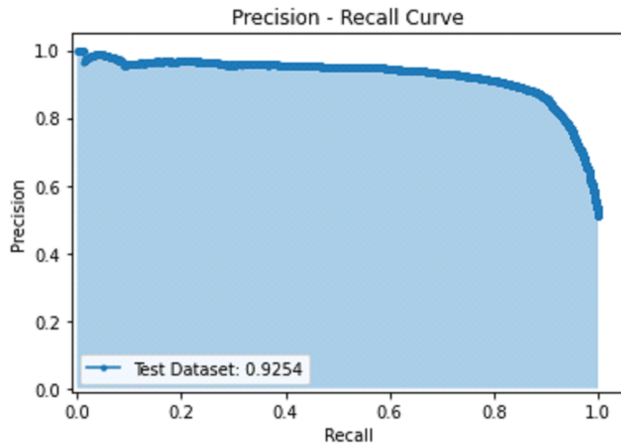
Figure 3 and 4 show the ROC and PRC curves for the final 3-BERT model. This seems to be our best model, it also had the longest training time. In comparison, the vastly more complicated CEM surpasses it in terms of AUPRC but fails to provide a high quality ROC curve. Figure 5 and 6 show the curves for that model.



**Fig. 3:** 3-BERT ROC

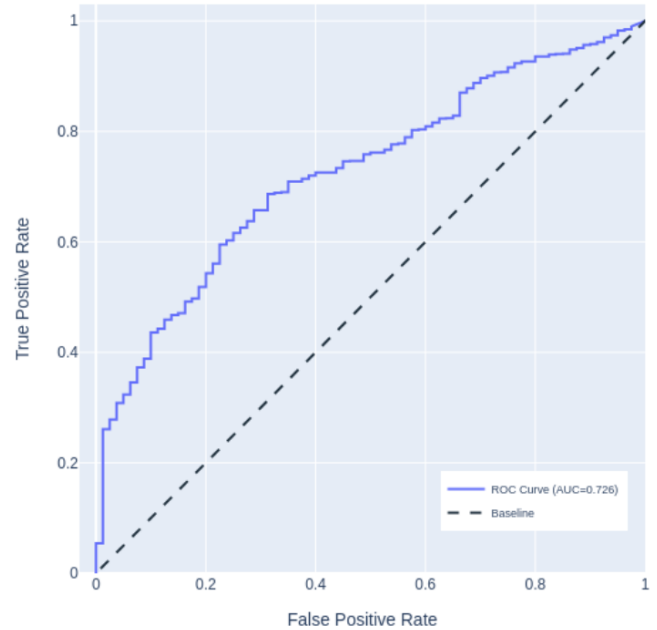
## 4. DISCUSSION

We were surprised that the 3-BERT model surpassed the Context Embedding model in terms of AUROC, because it uses less information. It could be the case that the tranformer encoder layer in the Context Embedding Model simply adds too much complexity and makes it hard to learn the dependencies in the data. A consideration is also that our language models could quickly run into overfitting because there are way more parameters in the model than we have training examples. This



**Fig. 4:** 3-BERT PRC

is usually prevented though by stopping early, and we did not train the models further when we realized that the validation error went up again. In the future it would be interesting to see what happens if we have more data available or different labels. We also feel like the models cannot be 100% accurate in all cases, as they are a result of a different model used by Amazon. This leaves us wondering if our model could perform very differently given labels produced by a different model.



**Fig. 5:** Context Embedding Model ROC

## 5. REFERENCES

### Bibliography

Ken Gu and Akshay Budhkar. A package for learning on tabular and text data with transformers. In *Proceedings of the Third Workshop on Multimodal Artificial Intelligence*, pages 69–73, Mexico City, Mexico, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.maiworkshop-1.10. URL <https://aclanthology.org/2021.maiworkshop-1.10>.

Rami Mohawesh, Shuxiang Xu, Matthew Springer, Muna Al-Hawawreh, and Sumbal Maqsood. Fake or genuine? contextualised text representation for fake review detection. *arXiv preprint arXiv:2112.14343*, 2021.

Chia-Hsien Weng, Kuan-Cheng Lin, and Jia-Ching Ying. Detection of chinese deceptive reviews based on pre-trained language model. *Applied Sciences*, 12(7):3338, 2022.

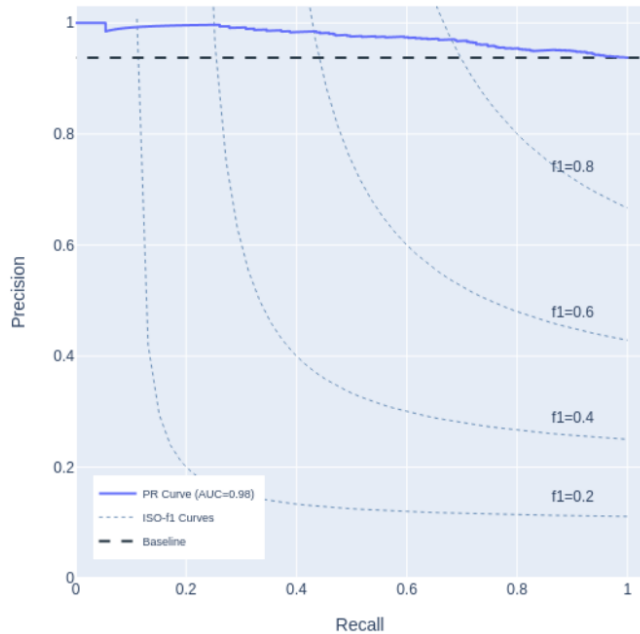


Fig. 6: Context Embedding Model PRC

## 6. APPENDIX

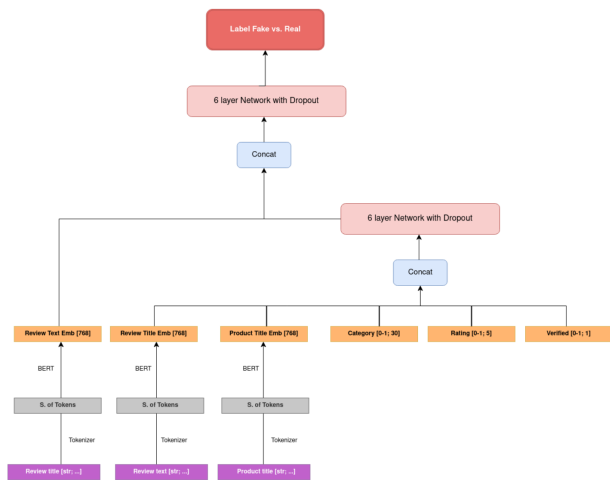


Fig. 7: Review Embedding Network

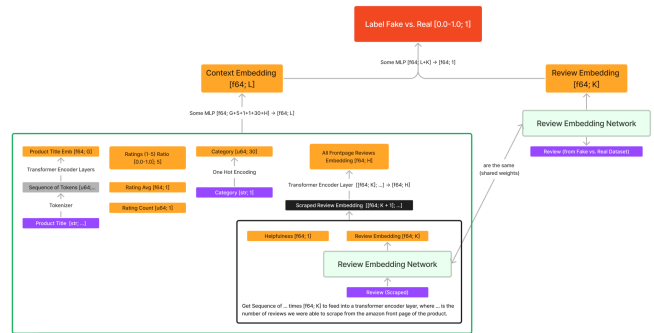


Fig. 8: Context Embedding Model

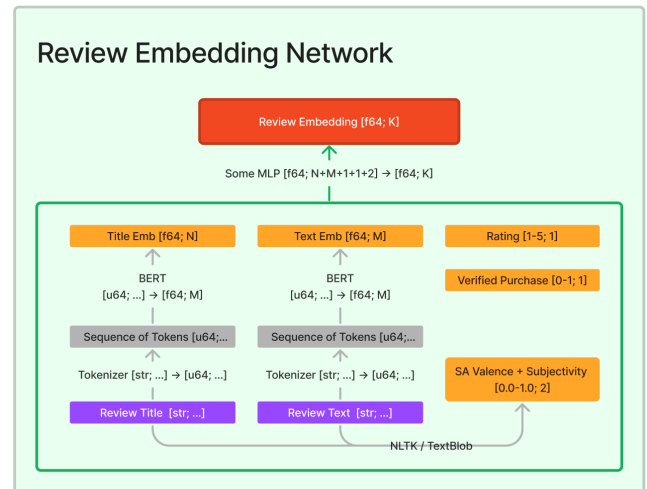


Fig. 9: Review Embedding Network