

## **Toxic Comment Detection in Amazon Product Reviews**

### **1. Introduction and Project Statement**

Online product reviews are a cornerstone of digital commerce, influencing consumer decisions and brand reputation. However, the open nature of these reviews allows for the proliferation of toxic language, comments that are abusive, offensive, or aggressive in tone. This project addresses the detection of such toxic reviews using modern natural language processing (NLP) techniques. The goal is to develop a binary classification model that can distinguish between toxic and non-toxic comments in the “All Beauty” category of the Amazon Reviews 2023 dataset. This effort aims to support content moderation tools, improve user experience, and safeguard platforms from harmful discourse.

Toxic language in product reviews not only reflects poorly on the platform but also diminishes user trust. By implementing a system capable of identifying toxic reviews, we can support moderation efforts and flag potentially harmful content for review. Unlike platforms such as social media, where toxic content often includes overt hate speech, toxicity in reviews can be more subtle, requiring careful contextual understanding. This makes the task of automatic detection particularly well-suited for transformer-based models like BERT.

### **2. Data Sources and Technologies Used**

#### ***A. Data Source***

- Amazon Reviews 2023 Dataset (McAuley Lab via Hugging Face):
  - Subset: `raw_review_All_Beauty`
  - Includes review text, rating, product ID, timestamps, and other metadata

- Format: JSON-based dataset loaded via Hugging Face datasets library

This dataset was selected due to its size, credibility, and diversity of user-generated reviews. I chose the “All Beauty” subset to narrow the scope while maintaining textual variability.

### ***B. Tools and Technologies***

- Python 3.11: Primary language used for development and analysis
- Hugging Face Transformers: For loading pretrained models and fine-tuning BERT
- Scikit-learn: For evaluation metrics and data visualization
- NLTK: For preprocessing and stopwords removal
- PyTorch: Backend framework for training the transformer model
- Jupyter Notebook: Interactive development and result visualization
- Matplotlib/Seaborn: For plotting confusion matrices and performance metrics

## **3. Methods Employed**

### ***A. Data Preprocessing and Weak Labeling***

The original dataset contained millions of reviews. I filtered this down to a manageable sample of 2,000 reviews. A weak supervision approach was used to assign toxicity labels. The labeling heuristic involved searching for toxic terms such as “trash”, “hate”, or “idiot” using regular expressions. This approach allowed us to avoid manual annotation but introduced noise.

### ***B. Exploratory Text Analysis***

We tokenized the review text and removed common English stopwords using NLTK. A word frequency analysis helped highlight the most common terms in the corpus. This step

allowed us to observe differences in language usage between potentially toxic and non-toxic reviews.

### ***C. Data Splitting and Tokenization***

To evaluate model generalization, the data was split into an 80/20 train-validation set using stratified sampling to preserve class distribution. Each review was tokenized using bert-base-uncased with padding and truncation (max length 128 tokens) to prepare inputs for the transformer model.

### ***D. Model Training***

I fine-tuned the BERT base model for binary sequence classification. Training was done over three epochs using the Hugging Face Trainer API. The training used the following hyperparameters:

- Batch size: 8
- Learning rate: 2e-5
- Optimizer: AdamW with weight decay
- Evaluation strategy: Per epoch

### ***E. Prediction and Inference***

A custom prediction function was implemented to allow testing on individual reviews. This function tokenizes the input, performs inference with the fine-tuned model, and returns the predicted class and associated confidence score.

## **4. Results**

The objective of this project was to develop a binary classifier capable of detecting toxic comments in product reviews, using weak labels derived from keyword-based heuristics. I

trained a BERT-based sequence classification model using a subset of the “All Beauty” category from the Amazon Reviews 2023 dataset. The results below detail the model’s classification performance, class confidence behavior, and overall interpretability of the learned toxic features.

*A. Classification Performance*

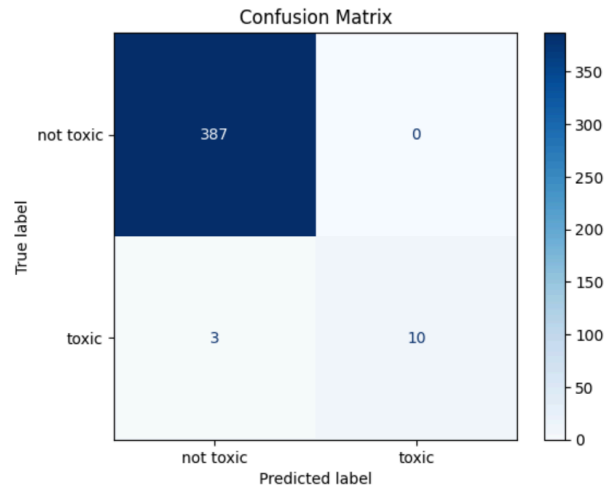
Using a stratified split (80/20) on a sample of 2,000 reviews (with 13 labeled toxic), the model achieved the following metrics on the validation set:

<i>Class</i>	<i>Precision</i>	<i>Recall</i>	<i>F1-Score</i>	<i>Support</i>
<i>Not Toxic</i>	0.99	1.00	1.00	387
<i>Toxic</i>	1.00	0.77	0.87	13
<i>Accuracy</i>			0.99	400

The toxic class, while underrepresented, was classified with 100% precision, indicating zero false positives. This makes the model conservative - it only flags reviews as toxic when very confident. However, recall was 77%, suggesting that 3 toxic reviews were misclassified as not toxic (false negatives). This tradeoff is acceptable in some domains, but it emphasizes the importance of high-quality labeling if recall is to be improved. In contrast, the non-toxic class was predicted with near-perfect performance, likely due to its high representation and relatively consistent linguistic patterns.

### ***B. Confusion Matrix Analysis***

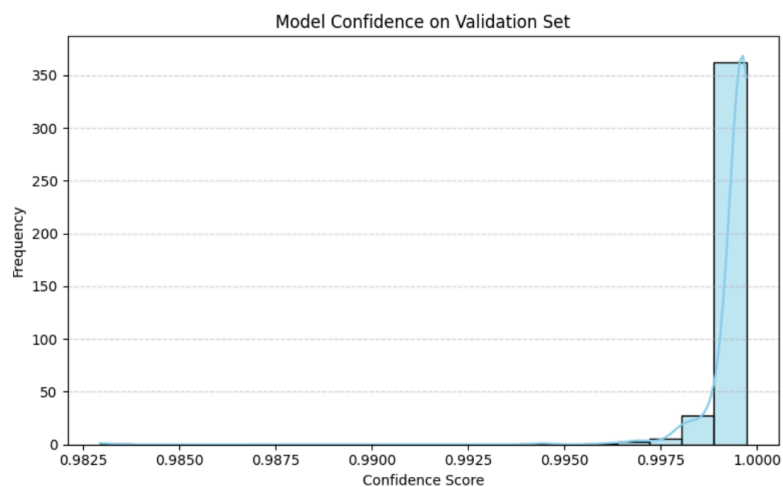
The confusion matrix below confirms the model's conservative classification approach:



This balance reflects a high true positive rate for both classes, and importantly, the absence of false alarms - a critical requirement for toxicity detection in production environments.

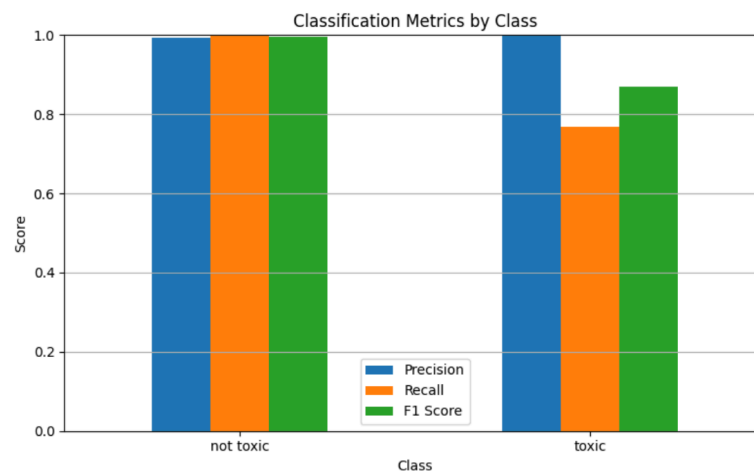
### ***C. Confidence Distribution***

The model's confidence in its predictions was visualized using a histogram of softmax outputs:



Most confidence scores lie extremely close to 1.0, with minimal variance. While this reflects the model's confidence and correlates with its accuracy, such overconfidence could signal limited model calibration or overfitting to patterns present in the weak-labeled training data. If deployed, additional techniques like temperature scaling could be considered to smooth the probability distribution and better reflect model uncertainty.

#### ***D. Per-Class Metric Visualization***



To further illustrate performance per class, a bar plot was generated showing precision, recall, and F1-score. This visualization reaffirms:

- Perfect precision for both classes
- Slightly reduced recall for the toxic class due to missed examples.
- Balanced F1 score - showing good generalization despite label imbalance.

#### ***E. Class Imbalance and Labeling Insight***

The dataset contained only ~0.65% toxic reviews after keyword-based weak labeling. This extreme imbalance, while realistic for many customer-facing platforms, made it challenging for the model to learn general toxic language patterns. The fact that BERT

still performed well is a testament to the effectiveness of transfer learning from pretrained language models.

However, this also suggests that further improvements could be made by:

- Using semi-supervised learning or active learning to refine labels.
- Augmenting toxic samples via data synthesis or bootstrapping
- Incorporating lexicon-based or model-driven labeling pipelines (e.g., Perspective API, BART zero-shot classification) to replace or supplement keyword-based heuristics.

## **5. Conclusion**

This project explored the application of natural language processing and transformer-based models to the task of toxic comment detection in product reviews. By leveraging a weakly-labeled subset of the “All Beauty” category from the Amazon Reviews 2023 dataset, we trained and evaluated a binary classifier using the pre-trained BERT architecture fine-tuned for sequence classification.

Despite the challenges posed by extreme class imbalance and the noisy nature of weak labels, the model achieved high performance, boasting 99% accuracy, perfect precision for both classes, and a notable F1 score of 0.87 for the minority toxic class. The model’s high confidence, supported by visual diagnostics and performance metrics, confirms its ability to generalize well to the validation set. However, slightly lower recall on the toxic class and the tightly skewed confidence distribution suggest that further refinement in labeling or calibration could improve sensitivity to subtler toxicity.

This work demonstrates that transformer-based models like BERT can perform robustly even with weak supervision, making them viable for scalable and automated moderation tasks. Future work can expand on this by incorporating more nuanced labeling techniques, exploring domain adaptation across product categories, or evaluating interpretability through attention visualization.

In sum, this project validates a practical pipeline for building toxicity detection models using publicly available product review data, while also highlighting the importance of thoughtful labeling and evaluation when dealing with sensitive or imbalanced tasks in NLP.

## 6. References

- McAuley, J., et al. (2023). Amazon Product Reviews Dataset. Retrieved from:  
<https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023>
- Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.
- Hugging Face Transformers Documentation: <https://huggingface.co/docs/transformers>
- scikit-learn Documentation:  
<https://scikit-learn.org/stable/modules/classes.html#module-sklearn.metrics>
- NLTK Documentation: <https://www.nltk.org>
- Wolf, T., et al. (2020). Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on EMNLP: System Demonstrations*.