# Enhancing Resilience to Adversarial Attacks in Fake News Detection

**Ankur Aggarwal**
New York University
Email: `aa10336@nyu.edu`
Net ID: aa10336

**Chirayu Bhatt**
New York University
Email: `cb5143@nyu.edu`
Net ID: cb5143

**Muskan Gandhi**
New York University
Email: `mng9349@nyu.edu`
Net ID: mng9349

## Abstract

Advancements in machine learning (ML) security are essential to counter the sophistication of adversarial attacks, especially in fake news detection. This study builds upon the UPFD model, integrating adversarial training to enhance model resilience. By incorporating the Fast Gradient Sign Method (FGSM) into the training of a Graph Neural Network (GNN), we significantly improve the robustness against adversarial attacks. Results indicate our model's enhanced accuracy in adversarial conditions, contributing to the field of ML security and misinformation mitigation. The modifications and training process are thoroughly documented for replication and further research.

**Code Repository:** `https://github.com/ankur928/ml-for-cybersec`

## 1 Introduction

### 1.1 Background

In the modern digital landscape, the phenomenon of fake news has become a pervasive challenge, significantly impacting public discourse and trust in media. The rapid proliferation of misinformation, facilitated by social media and other digital platforms, has led to a growing concern about its influence on political, social, and public health matters. Machine Learning (ML) has emerged as a key player in the battle against fake news, with models being developed to automate the detection of false information. The UPFD model, as presented in [1], represents a significant advancement in this arena. However, the evolving sophistication of fake news, coupled with the emergence of advanced adversarial techniques used to create and spread such content, poses new challenges that necessitate continual innovation in ML security.

### 1.2 Motivation for Extension

The motivation behind this project stems from the escalating complexity of adversarial tactics employed to bypass traditional fake news detection systems. Modern adversarial techniques are not only becoming more nuanced but also more adept at mimicking genuine content, thus posing a significant threat to the efficacy of existing detection models. In response, there is a critical need to develop models that go beyond mere detection and are capable of effectively countering these advanced adversarial methods. Our project seeks to address this gap by enhancing the robustness of existing fake news detection models, thereby contributing to more reliable and secure information dissemination in the digital age.

## 1.3 Contribution

This paper presents a substantial extension of the UPFD model through the integration of adversarial training methodologies, specifically the Fast Gradient Sign Method (FGSM). FGSM is renowned for its effectiveness in enhancing model resilience against adversarial attacks. By incorporating FGSM into the UPFD's Graph Neural Network (GNN) architecture, we have developed a more fortified model capable of maintaining high accuracy in the face of sophisticated adversarial manipulations. This paper details the modifications made, the rationale behind these changes, and the observed impact on model performance. The extended model not only demonstrates improved resilience in detecting fake news but also serves as a framework for future research in ML security against adversarial attacks.

## 1.4 Code Availability and Reproducibility

Ensuring the reproducibility of our findings and facilitating further research in this domain, we provide comprehensive documentation and source code for all implementations and modifications. The complete codebase, along with detailed documentation, is publicly available in our GitHub repository: `https://github.com/ankur928/ml-for-cybersec`. This repository includes scripts for model training, adversarial attack simulations, and performance evaluation, enabling researchers and practitioners to replicate our study and extend our work.

# 2 Key Ideas

## 2.1 Advancing Fake News Detection

### 2.1.1 The Challenge of Adversarial Attacks

Discuss the growing challenge of adversarial attacks in fake news detection. Explain how these attacks are becoming more sophisticated, making it harder for traditional detection methods to keep up.

### 2.1.2 Importance of Robust Models

Emphasize the necessity for models that are not just accurate but also robust against such adversarial manipulations. Detail why this robustness is critical in the current digital landscape.

## 2.2 Integration of Adversarial Training

### 2.2.1 Adopting FGSM

Explain the choice of the Fast Gradient Sign Method (FGSM) for adversarial training. Discuss its effectiveness in enhancing model resilience and why it was chosen over other adversarial training methods.

### 2.2.2 Impact on Model Performance

Highlight how FGSM integration improves the model's performance in detecting manipulated content. Discuss any specific findings that illustrate this improvement.

## 2.3 Contributions to ML Security

### 2.3.1 Novelty in Approach

Outline the novel aspects of your approach, such as specific adaptations made to the UPFD model or unique ways in which adversarial training was implemented.

### 2.3.2 Broader Implications

Discuss the broader implications of your work for machine learning security, especially in relation to the trustworthiness and credibility of information on digital platforms.

### 2.4 Practical Applications and Future Directions

#### 2.4.1 Real-World Impact

Describe potential real-world applications of your enhanced model. For instance, how it could be employed by social media platforms or news aggregators to filter out fake news.

#### 2.4.2 Path Forward

Suggest directions for future research that builds on your work. This could include testing the model on different datasets, exploring other forms of adversarial training, or integrating your approach into existing fake news detection systems.

### 2.5 Conclusion

Conclude by reiterating the significance of your contributions and their potential impact on the field of fake news detection and ML security.

## 3 Methodology

### 3.1 Model Foundation

Our methodology is founded on the Graph Neural Network (GNN) architecture established by the UPFD framework [1]. The original GNN effectively captures the relational information embedded within the data, a critical feature for differentiating between authentic and fake news.

### 3.2 Adversarial Training Integration

#### 3.2.1 Incorporating FGSM

To enhance the model's resilience, we integrate the Fast Gradient Sign Method (FGSM), an adversarial training technique known for its efficacy in defense mechanisms. FGSM introduces controlled perturbations to the input data, simulating potential adversarial attacks during training.

#### 3.2.2 Modifications to the Training Process

The training process is modified to include these adversarial examples in each epoch. This ensures that the model not only learns from the true distribution of the data but also gains resilience against input data that has been intentionally distorted.

### 3.3 Extension of the Original Work

Building upon the codebase from the existing GitHub repository [2], we extend the work by incorporating additional modules for generating and processing adversarial examples. This extension is aimed at providing a more comprehensive learning experience for the model.

### 3.4 Dataset Utilization

#### 3.4.1 Employing GossipCop and PolitiFact

Consistent with the original study, we utilize the GossipCop and PolitiFact datasets due to their rich and diverse content. These datasets offer a robust benchmark for training and evaluating the enhanced model.

#### 3.4.2 Generation of Adversarial Samples

We generate adversarial samples from these datasets using the FGSM approach. The samples are verified to ensure they present a realistic yet challenging scenario for the model, reflecting the sophistication of potential real-world attacks.

### 3.5 Training Protocol

#### 3.5.1 Epochs and Batching

The model is trained over multiple epochs with batch processing, allowing for efficient computation and effective gradient updates. The adversarial examples are batched with the genuine data to balance the learning focus between standard detection and adversarial robustness.

#### 3.5.2 Hyperparameter Optimization

Hyperparameters are fine-tuned to optimize the model's performance. This includes adjusting learning rates, the number of GAT layers, and the intensity of adversarial perturbations represented by the epsilon value.

### 3.6 Evaluation Metrics

#### 3.6.1 Accuracy and Loss

The primary metrics for evaluating the model are accuracy and loss. Accuracy measures the model's ability to correctly classify news articles, while loss provides insight into the model's predictive confidence.

#### 3.6.2 Robustness Assessment

Additionally, we assess the model's robustness by evaluating its performance on adversarially perturbed data. This assessment is critical for understanding the model's defense capability against attacks.

### 3.7 Reproducibility and Code Availability

#### 3.7.1 Documentation and Sharing

Ensuring reproducibility, all code modifications, and the extended training protocols are thoroughly documented. The complete codebase is made available in our GitHub repository [4], providing the community with the means to replicate our study and further the research.

## 4 Code Implementation and Results

### 4.1 Adaptation and Extension of Existing Framework

#### 4.1.1 Utilizing the UPFD Framework

Our methodology significantly builds upon the User Preference-aware Fake News Detection (UPFD) framework. We chose the UPFD model, implemented using Graph Neural Networks (GNNs) in Pytorch-Geometric, for its proven effectiveness in the domain of fake news detection. This choice allowed us to focus on enhancing an already robust model with new capabilities.

#### 4.1.2 Integrating Adversarial Training

The key extension in our work involves the integration of adversarial training, particularly using the Fast Gradient Sign Method (FGSM). This approach is pivotal in enhancing the model's resilience against more sophisticated forms of misinformation that employ adversarial techniques for evasion.

### 4.2 Dataset and Model Configuration

#### 4.2.1 Utilizing GossipCop and PolitiFact Datasets

We employed the GossipCop and PolitiFact datasets for their comprehensive and diverse range of news articles. These datasets offer a rich ground for testing the enhanced capabilities of our model in distinguishing real news from fake under varied conditions.

### 4.2.2 Feature Incorporation

Our model extends the feature set used in the original UPFD framework by integrating BERT and spaCy word2vec encodings, Twitter profile features, and content features. This comprehensive feature set allows for a more nuanced understanding of the data, enhancing the model's ability to detect subtle forms of fake news.

## 4.3 Implementation Details

### 4.3.1 Model Training and Fine-tuning

We adapted and fine-tuned the training protocols to effectively incorporate adversarial training. This involved not only training the model on standard data but also on data that had been manipulated using the FGSM technique. Hyperparameter tuning was conducted to find the optimal balance between model performance and training efficiency.

### 4.3.2 Code Documentation and Accessibility

Our code modifications and the extended training process are thoroughly documented for clarity and ease of replication. The entire codebase, including the scripts for model training and evaluation, is publicly available in our GitHub repository [4], ensuring accessibility for the research community.

## 4.4 Results and Evaluation

### 4.4.1 Performance Metrics

The extended model was evaluated on two primary metrics: accuracy and training loss. The results indicate an improvement in the model's performance, particularly in its ability to maintain accuracy in the presence of adversarial perturbations.

| Epsilon | Train Loss | Validation Accuracy |
|---------|-----------|---------------------|
| 0.00 | 0.9276 | 0.8515 |
| 0.01 | 0.9211 | 0.8515 |
| 0.02 | 0.7730 | 0.8515 |
| 0.05 | 1.0842 | 0.8515 |
| 0.10 | 0.6221 | 0.8515 |

Table 1: Model performance across different epsilon values for FGSM adversarial training

### 4.4.2 Comparative Analysis

A comparative analysis with the baseline UPFD models demonstrates that our extended model not only retains the original accuracy but also shows enhanced resilience to adversarial attacks. This is a clear indication of the effectiveness of FGSM adversarial training in improving the model's robustness.

### 4.4.3 Implications of Results

The results have significant implications for the field of ML security, particularly in fake news detection. They suggest that adversarial training can be a potent tool in augmenting the capabilities of existing models to handle the evolving nature of misinformation.

## 4.5 Conclusion

In conclusion, our extension of the UPFD model with FGSM adversarial training techniques represents a significant advancement in the resilience of fake news detection models. This study not only demonstrates the potential of such approaches in enhancing ML security but also provides a foundation for future research in this critical area of information integrity and digital trust.

# 5 Discussion and Analysis

## 5.1 Efficacy of Adversarial Training

### 5.1.1 Resilience to Adversarial Attacks

The integration of the Fast Gradient Sign Method (FGSM) into the UPFD framework marks a significant advancement in the robustness of fake news detection models. Our results demonstrate that the model maintains high accuracy across various levels of adversarial perturbation, indicating its enhanced ability to withstand adversarial attacks compared to the baseline model.

### 5.1.2 Implications for Model Generalization

The consistent performance across different epsilon values suggests that the model has not only learned to identify fake news in standard scenarios but has also generalized to effectively counter manipulated content. This generalization is critical for practical applications where the nature of adversarial attacks can be diverse and unpredictable.

## 5.2 Comparison with the Original UPFD Model

### 5.2.1 Performance Metrics

While the original UPFD model provided a strong foundation in detecting fake news using GNNs, our extended model shows improved resilience when subjected to adversarial conditions. This improvement is evident in the model's ability to maintain validation accuracy even in the presence of adversarial perturbations.

### 5.2.2 Adversarial Training as a Necessary Enhancement

The results affirm that adversarial training, particularly FGSM, is a necessary enhancement for models operating in security-sensitive domains like fake news detection. It equips the model to cope with evolving tactics used in misinformation campaigns.

## 5.3 Challenges and Limitations

### 5.3.1 Computational Overheads

One of the challenges observed in our approach is the increased computational overhead due to adversarial training. This aspect warrants consideration, especially for deployment in resource-constrained environments.

### 5.3.2 Scope of Adversarial Attacks

While FGSM provides a robust method for adversarial training, it represents a specific type of attack. Future work could explore the model's resilience against other forms of adversarial attacks to ensure a comprehensive defense mechanism.

## 5.4 Future Research Directions

### 5.4.1 Exploring Other Adversarial Techniques

Further research could investigate the effectiveness of other adversarial training techniques, potentially offering a broader defense spectrum against different types of adversarial attacks.

### 5.4.2 Application to Diverse Datasets

Expanding the model's application to other datasets, beyond GossipCop and PolitiFact, could provide insights into its adaptability and scalability in various contexts.

## 5.5 Conclusion

In conclusion, our analysis indicates that the addition of adversarial training to the UPFD framework significantly enhances the model's capability to detect and resist fake news. The approach not only improves accuracy in adversarial conditions but also brings to light important considerations for future research in the domain of ML security and fake news detection.

# 6 Conclusion

## 6.1 Summary of Contributions

This study represents a substantial advancement in machine learning security, with a specific focus on the increasingly vital field of fake news detection. By extending the UPFD model through the integration of the Fast Gradient Sign Method (FGSM) for adversarial training, we have significantly enhanced the model's resilience. Our comprehensive experiments demonstrate that this enhancement not only maintains but, in some cases, improves the accuracy of the model in detecting fake news, even when faced with sophisticated adversarial manipulations. This contribution is particularly notable given the growing prevalence and sophistication of misinformation campaigns.

## 6.2 Advancements in Fake News Detection

The application of FGSM adversarial training within the UPFD framework marks a critical development in the struggle against digital misinformation. This advancement is not merely a technical achievement but also a strategic one, addressing the need for detection models that are robust against evolving and increasingly subtle adversarial attacks. Our work underscores the importance of anticipatory and adaptive approaches in the realm of information integrity.

## 6.3 Challenges and Future Directions

### 6.3.1 Technical Challenges

While the model exhibits promising results, it also introduces new challenges, particularly in computational efficiency. The increased computational load due to the adversarial training process necessitates further optimization for practical, large-scale applications.

### 6.3.2 Exploring Broader Applications

Future research could explore the application of our approach to a wider array of datasets, encompassing different types of news and social media content. Additionally, investigating the effectiveness of other adversarial methods could provide a more robust framework that is adaptable to various forms of misinformation.

## 6.4 Implications for ML Security

Our research contributes significantly to the field of ML security. By demonstrating the effectiveness of integrating adversarial training into fake news detection models, we pave the way for more secure and reliable information filtering systems. This work has potential applications in a range of domains beyond news, including social media content monitoring and public opinion analysis, where the integrity of information is crucial.

## 6.5 Acknowledgments

## 6.6 Concluding Remarks

In closing, the integration of adversarial training into fake news detection models, as illustrated in our study, represents a significant step forward in combating digital misinformation. This approach not only enhances the accuracy and resilience of detection models but also contributes to the larger goal of fostering a more trustworthy digital information environment. We are optimistic that our contributions will serve as a catalyst for further innovation and research in machine learning security, particularly in the ongoing battle against fake news and misinformation.

## References

[1] Y. Dou, K. Shu, C. Xia, P. S. Yu, and L. Sun, "User Preference-aware Fake News Detection," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 1199-1208. DOI: 10.1145/3404835.3463030. [Online]. Available: https://arxiv.org/pdf/2104.12259.pdf

[2] Y. Dou, K. Shu, et al., "GNN-based Fake News Detection - UPFD," GitHub repository, last updated July 2021. [Online]. Available: https://github.com/safe-graph/GNN-FakeNews

[3] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and Harnessing Adversarial Examples," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. [Online]. Available: https://arxiv.org/abs/1412.6572

[4] A. Aggarwal, "Enhancing Resilience to Adversarial Attacks in Fake News Detection" GitHub repository, [Online]. Available: https://github.com/ankur928/machine-learning

[5] C. Shu, S. Wang, and H. Liu, "Beyond News Contents: The Role of Social Context for Fake News Detection," in *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, 2019, pp. 312–320. DOI: 10.1145/3289600.3290994.

[6] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake News Detection on Social Media: A Data Mining Perspective," *ACM SIGKDD Explorations Newsletter*, vol. 19, no. 1, pp. 22-36, 2017. DOI: 10.1145/3137597.3137600.

[7] DeepFindr, "Fake News Detection using Graphs with Pytorch Geometric," YouTube, uploaded on Jan 17, 2022. [Video]. Available: https://www.youtube.com/watch?v=QAIVFr24FrA