

## Story Line

---

Friday, January 24, 2025 1:49 PM

### Title: The Rogue Reviewer – A Data Poisoning Adventure

**Synopsis:** You are a machine learning detective tasked with investigating a growing threat to sentiment analysis models. Reports suggest that an unknown entity, "The Rogue Reviewer," has been injecting malicious data into public datasets.

The goal:

- To undermine the integrity of AI systems and spread misinformation.

Your mission is to :

- analyze the attack,
- measure its impact and;
- implement defenses to thwart future sabotage attempts.

### Tasks:

1. **Inject Malicious Data**
2. **Train the Model on the Poisoned Dataset**
3. **Mitigate the Poisoning Attack**

Determined to neutralize the Rogue's sabotage, you can employ countermeasures:

- **Outlier Detection:** Spotting and isolating suspicious reviews through clustering techniques.
- **Data Sanitization:** Cleaning the dataset by removing identified anomalies.
- **Robust Defenses:** Introducing adversarial training and validating data before training to bolster resilience against future attacks.

After retraining the model on the sanitized dataset, its performance improves, marking a significant step toward mitigating the impact of data poisoning.

### Epilogue: Lessons from the Rogue Reviewer

Through this hands-on investigation, you've unraveled the mechanics of data poisoning attacks and gained valuable experience in combating them. While the Rogue Reviewer remains at large, your expertise ensures that AI systems are better equipped to resist future attacks, safeguarding the integrity of machine learning models everywhere.

