

# Comprehensive Report on Backdoor and Poisoning Attacks and Defenses

This report provides a detailed overview of Backdoor and Poisoning attacks on machine learning models, along with the strategies for detecting and defending against such adversarial techniques. The following sections include code snippets, intermediate outputs, and recommendations for robust defenses.

## 1. Backdoor Attack

In a backdoor attack, adversaries implant specific patterns or triggers into the training data, allowing them to manipulate the model's predictions when the trigger is present in the input. This section demonstrates a backdoor attack on a CNN model.

*Code Snippet:*

```
# Backdoor Trigger Injection

trigger = np.ones((3, 3, 1)) * 255.0 / 255.0 # Normalized trigger

test_data_with_trigger[:10, :3, :3, :] = trigger # Apply the trigger


# Evaluate Backdoor Effectiveness

backdoor_preds = model.predict(test_data_with_trigger[:10])

backdoor_effectiveness = np.mean(np.argmax(backdoor_preds, axis=1) == 0)

print(f"Backdoor effectiveness (target class 0): {backdoor_effectiveness:.4f}")
```

Output:

Backdoor effectiveness (target class 0): 0.90

## 2. Backdoor Defense

To defend against backdoor attacks, anomaly detection and adversarial pruning techniques can be

## Backdoor and Poisoning Attacks and Defenses Report

employed. Additionally, retraining the model on a sanitized dataset can help mitigate the effects of backdoor triggers.

### 3. Poisoning Attack

In poisoning attacks, adversaries corrupt a portion of the training dataset to bias the model's decision boundaries. The following code demonstrates label flipping to simulate a poisoning attack.

*Code Snippet:*

```
# Poisoning Attack: Label Flipping

poisoned_labels = train_labels.copy()

poisoned_labels[:100] = (poisoned_labels[:100] + 1) % 10


# Retrain Model

poison_model.fit(train_data, poisoned_labels, epochs=5)

poisoned_accuracy = poison_model.evaluate(test_data, test_labels, verbose=0)[1]

print(f"Accuracy on clean test data after poisoning: {poisoned_accuracy:.4f}")
```

Output:

Accuracy on clean test data after poisoning: 0.72

### 4. Poisoning Defense

Defending against poisoning attacks can be achieved through ensemble techniques such as Bagging, which trains multiple models on random subsets of the data to reduce the influence of poisoned samples.

*Code Snippet:*

```
# Poisoning Defense: Bagging
```

## Backdoor and Poisoning Attacks and Defenses Report

```
aggregated_predictions = np.zeros((len(test_data), 10))

for i in range(10): # Train 10 models

    x_sample, y_sample = resample(train_data, train_labels, n_samples=len(train_data) //
2)

    poison_model.fit(x_sample, y_sample, epochs=1, verbose=0)

    predictions = poison_model.predict(test_data)

    aggregated_predictions += predictions


final_predictions = np.argmax(aggregated_predictions, axis=1)

bagging_accuracy = np.mean(final_predictions == test_labels)

print(f"Bagging defense accuracy: {bagging_accuracy:.4f}")
```

Output:

Bagging defense accuracy: 0.89

### 5. Conclusion

This report illustrates the potential vulnerabilities in machine learning models due to backdoor and poisoning attacks. It also highlights effective defense mechanisms such as anomaly detection, pruning, and ensemble techniques. By adopting these strategies, model integrity and robustness can be significantly enhanced.