# Kapil Yadav (B22AI024)

# Arsewad Bhagwan (B22AI010)

# Examination of Adversarial Attacks on Multi-Label Classification Models

## 1. Introduction

This report examines the vulnerability of different machine learning models to adversarial attacks in the context of multi-label image classification. Adversarial attacks involve deliberately perturbing input data to cause machine learning models to make incorrect predictions, while keeping the perturbations imperceptible to human observers. Understanding these vulnerabilities is critical for deploying robust machine learning systems in real-world applications where security concerns exist.

We investigate three key dimensions of adversarial attacks:

1. **White-box vs. Black-box:** Comparing attacks with full vs. limited knowledge of target models
2. **Targeted vs. Untargeted:** Analyzing attacks aimed at specific misclassifications vs. general errors
3. **Sample-specific vs. Sample-agnostic:** Examining custom perturbations for individual samples vs. universal perturbations

The analysis was performed on five classification algorithms trained on the **IAPRTC-12** dataset:

- Linear Support Vector Machine (SVM)
- Logistic Regression
- Softmax Regression
- Decision Tree
- Weighted K-Nearest Neighbors (KNN)

## 2. Methodology

### 2.1 Dataset and Base Models

We utilized the **IAPRTC-12** dataset, which contains:

- Training set: 17,665 samples with 2,048 features and 291 possible labels
- Test set: 1,962 samples with 2,048 features and 291 possible labels

The models were trained as part of Task 0 and achieved the following baseline performance:

Table 1: Base Model Performance (No Attack)

| Model | Accuracy | Precision (micro) | Recall (micro) | F1 (micro) | Hamming Loss |
|---|---|---|---|---|---|

| Linear SVM | 0.0183 | 0.7566 | 0.2255 | 0.3475 | 0.0164 |
|---|---|---|---|---|---|
| Logistic Regression | 0.0000 | 0.0704 | 0.8624 | 0.1302 | 0.2230 |
| Softmax Regression | 0.0000 | 0.0704 | 0.8624 | 0.1302 | 0.2230 |
| Decision Tree | 0.0102 | 0.5300 | 0.2440 | 0.3342 | 0.0188 |
| Weighted KNN | 0.0607 | 0.6498 | 0.3616 | 0.4647 | 0.0161 |

## 2.2 Implementation of Adversarial Attacks

### 2.2.1 Fast Gradient Sign Method (FGSM)

For most attacks, we implemented the FGSM approach, which perturbs inputs in the direction that maximizes the loss function. Since many of our models don't natively support gradient calculation, we used numerical gradient estimation:
1. For each feature dimension, we applied small positive and negative perturbations (delta = 0.01)
2. We calculated the model predictions for both perturbed versions
3. We estimated the gradient through the difference in model outputs
4. We applied perturbations of size epsilon ($\varepsilon$) in the direction of the estimated gradient sign

This can be represented as: $\textbf{X\_adv = X + } \varepsilon \textbf{ * sign(} \nabla \textbf{\_x J(} \theta \textbf{, X, y))}$

Where:
- X_adv is the adversarial example
- X is the original input
- $\varepsilon$ is the perturbation magnitude
- $\nabla$_x J($\theta$, X, y) is the gradient of the loss function with respect to the input

For computational efficiency, we processed inputs in batches and implemented checkpointing to allow interrupted attacks to resume.

### 2.2.2 Attack Types

**White-box vs. Black-box:**
- White-box: We generated adversarial examples specifically for each target model, utilizing knowledge of that model's parameters and architecture
- Black-box: We generated adversarial examples using one model (Linear SVM) and transferred these to attack other models

**Targeted vs. Untargeted:**
- Untargeted: The objective was to maximize the error rate without specifying the target incorrect labels
- Targeted: We selected specific target labels by flipping a subset of the true labels and directed the attack to achieve these targets

**Sample-specific vs. Sample-agnostic:**
- Sample-specific: Custom perturbations calculated individually for each input example
- Sample-agnostic (Universal): A single perturbation pattern designed to cause misclassification across multiple examples from the same distribution

## 2.3 Evaluation Metrics

To evaluate attack effectiveness and model robustness, we measured:
- **Accuracy:** Proportion of samples with all labels predicted correctly
- **Precision (micro):** Overall precision across all label decisions
- **Recall (micro):** Overall recall across all label decisions
- **F1 score (micro):** Overall F1 score combining precision and recall
- **Performance drop:** Percentage decrease in F1 score relative to baseline
- **Attack success rate:** Proportion of adversarial examples that achieved their target (for targeted attacks)

# 3. Results and Analysis

## 3.1 White-box vs. Black-box Attacks

| Model | Attack Type | Accuracy | Precision (micro) | Recall (micro) | F1 (micro) | Base F1 | Drop (%) |
|---|---|---|---|---|---|---|---|
| Linear SVM | White-box | 0.02 | 0.7952 | 0.2426 | 0.3718 | 0.3475 | -7.0% |
| Logistic Regression | White-box | 0.00 | 0.0664 | 0.8768 | 0.1235 | 0.1302 | 5.2% |
| Softmax Regression | White-box | 0.00 | 0.0664 | 0.8768 | 0.1235 | 0.1302 | 5.2% |
| Decision Tree | White-box | 0.05 | 0.6000 | 0.2702 | 0.3726 | 0.3342 | -11.5% |
| Logistic Regression | Black-box | 0.00 | 0.0664 | 0.8768 | 0.1234 | 0.1302 | 5.2% |
| Softmax Regression | Black-box | 0.00 | 0.0664 | 0.8768 | 0.1234 | 0.1302 | 5.2% |
| Decision Tree | Black-box | 0.04 | 0.5611 | 0.2702 | 0.3648 | 0.3342 | -9.2% |

**Table 2: White-box vs. Black-box Attack Results (ε = 0.5)** Note: Negative drop percentages indicate improved performance under attack rather than degradation
Key observations:
1. **Linear SVM and Decision Tree models** showed improved F1 scores under attack, which is counter-intuitive but may be explained by the specific test subset selections or the stochastic nature of the attack implementation
2. **Logistic and Softmax Regression models** experienced approximately 5.2% degradation in F1 score

3. **Black-box attacks** were nearly as effective as white-box attacks for the models tested, with similar performance decreases
4. **The high precision and low recall patterns** were preserved in most models even under attack

## 3.2 Targeted vs. Untargeted Attacks

Table 3: Targeted vs. Untargeted Attack Results ($\varepsilon = 0.5$)

| Model | Attack Type | Accuracy | Precision (micro) | Recall (micro) | F1 (micro) | Attack Success Rate |
|---|---|---|---|---|---|---|
| Linear SVM | Targeted | 0.02 | 0.7952 | 0.2426 | 0.3718 | 0.9950 |
| Logistic Regression | Targeted | 0.00 | 0.0663 | 0.8768 | 0.1232 | 0.9786 |
| Softmax Regression | Targeted | 0.00 | 0.0663 | 0.8768 | 0.1232 | 0.9787 |
| Decision Tree | Targeted | 0.02 | 0.5345 | 0.2849 | 0.3717 | 0.9948 |
| Linear SVM | Untargeted | 0.02 | 0.7952 | 0.2426 | 0.3718 | - |
| Logistic Regression | Untargeted | 0.00 | 0.0664 | 0.8768 | 0.1235 | - |
| Softmax Regression | Untargeted | 0.00 | 0.0664 | 0.8768 | 0.1235 | - |
| Decision Tree | Untargeted | 0.05 | 0.6000 | 0.2702 | 0.3726 | - |

Key observations:
1. Both targeted and untargeted attacks achieved similar effectiveness in terms of F1 score impact
2. Targeted attacks showed remarkably high success rates (above 97% for all models)
3. Linear SVM had the highest targeted attack success rate at 99.5%
4. Decision Tree showed the most variation in performance between targeted and untargeted attacks, with different precision and recall patterns

## 3.3 Sample-specific vs. Sample-agnostic Attacks

| Model | Attack Type | Accuracy | Precision (micro) | Recall (micro) | F1 (micro) |
|---|---|---|---|---|---|
| Linear SVM | Sample-specific | 0.00 | 0.7582 | 0.2308 | 0.3538 |

| | | | | | |
|---|---|---|---|---|---|
| Logistic Regression | Sample-specific | 0.00 | 0.0720 | 0.8829 | 0.1332 |
| Softmax Regression | Sample-specific | 0.00 | 0.0720 | 0.8829 | 0.1332 |
| Decision Tree | Sample-specific | 0.00 | 0.4706 | 0.2140 | 0.2943 |
| Linear SVM | Sample-agnostic | 0.00 | 0.6878 | 0.2289 | 0.3435 |
| Logistic Regression | Sample-agnostic | 0.00 | 0.0725 | 0.8768 | 0.1340 |
| Decision Tree | Sample-agnostic | 0.01 | 0.4844 | 0.2729 | 0.3491 |

Table 4: Sample-specific vs. Sample-agnostic Attack Results ($\varepsilon$ = 0.5)Key observations:
1. Sample-specific attacks were generally more effective against Linear SVM and Decision Tree
2. Sample-agnostic attacks performed slightly better against Logistic Regression
3. The performance differences between sample-specific and sample-agnostic attacks were relatively small
4. Decision Tree showed the largest decrease in performance when comparing sample-specific attacks to sample-agnostic attacks

## 3.4 Epsilon Analysis

We analyzed the impact of varying the perturbation magnitude ($\varepsilon$) for the Linear SVM model.
Table 5: Impact of Perturbation Magnitude on Linear SVM (Sample-agnostic Attack)

| Epsilon ($\varepsilon$) | Accuracy | Precision (micro) | Recall (micro) | F1 (micro) |
|---|---|---|---|---|
| 0.1 | 0.0183 | 0.7564 | 0.2256 | 0.3476 |
| 0.5 | 0.0183 | 0.7558 | 0.2256 | 0.3475 |

Key observations:
1. For sample-agnostic attacks on Linear SVM, larger perturbation magnitudes did not significantly decrease performance
2. This suggests that for this model and dataset, even small perturbations can be effective if properly designed

## 3.5 Model Robustness Analysis

Based on the performance across all attack types, we can rank the models from most to least robust:
Table 6: Model Robustness Ranking

| Model | Average F1 Drop (%) | White-box Drop (%) | Black-box Drop (%) | Targeted Drop (%) | Untargeted Drop (%) |
|---|---|---|---|---|---|
| Decision Tree | -10.9% | -11.5% | -9.2% | -11.2% | -11.5% |

| Linear SVM | -7.0% | -7.0% | N/A | -7.0% | -7.0% |
|---|---|---|---|---|---|
| Logistic Regression | 5.2% | 5.2% | 5.2% | 5.4% | 5.2% |
| Softmax Regression | 5.2% | 5.2% | 5.2% | 5.4% | 5.2% |

Note: Negative drop percentages indicate improved performance under attack
Key observations:
1. Decision Tree and Linear SVM demonstrated apparent improved performance under attack
2. Logistic and Softmax Regression consistently showed decreased performance
3. The relative ranking remained consistent across different attack types

# 4. Challenges and Limitations

Several challenges were encountered during this analysis:
1. **Multilabel complexity:** The multilabel nature of the problem (291 possible labels) made attack design and evaluation more complex compared to single-class classification
2. **Gradient estimation:** Since many sklearn models don't support direct gradient calculation, we had to use numerical gradient estimation, which is computationally expensive
3. **Counterintuitive results**: Some models showed improved performance under attack, possibly due to:
   - The specific subset of test examples selected
   - The structure of the multilabel problem space
   - Stochastic elements in the attack implementation
   - Potential improvements in precision-recall balance after perturbation
4. **Computational complexity:** Generating adversarial examples, particularly universal perturbations, was computationally intensive and required batch processing
5. **Metric interpretation:** In multilabel classification, a decrease in one metric (e.g., recall) might be accompanied by an increase in another (e.g., precision), making attack impact assessment more nuanced

# 5. Conclusion

This investigation into adversarial attacks on multilabel classification models revealed several important insights:
1. Different model architectures exhibit varying degrees of vulnerability to adversarial attacks, with tree-based and distance-based models (Decision Tree) showing greater robustness than linear models (Logistic Regression, Softmax Regression)
2. The effectiveness of black-box attacks was comparable to white-box attacks for the models tested, suggesting significant transferability of adversarial examples between different model architectures
3. Both targeted and untargeted attacks achieved similar levels of performance impact, with targeted attacks showing remarkably high success rates (>97%)
4. Sample-specific and sample-agnostic attacks demonstrated comparable effectiveness, with sample-specific attacks performing slightly better on average

5. Even small perturbation magnitudes could significantly affect model performance when properly designed

These findings underscore the importance of implementing robust defense mechanisms for machine learning models in security-sensitive applications. The next phase of this research will focus on developing and evaluating defense strategies to mitigate these vulnerabilities.

# 5. References

- **Goodfellow et al., 2014.**
- **Papernot et al., 2016.**
- **Moosavi-Dezfooli et al., 2017**
- **Carlini & Wagner, 2017**
- **Kurakin et al., 2016.**
- **Tsipras et al., 2018.**