

# 敵対的攻撃に対する堅牢性向上のための適応的マルチチャンネル選択法

Adaptive multi-channel selection method for enhancing robustness  
against adversarial attacks

松井清修<sup>1)</sup>

指導教員 青木輝勝<sup>2)</sup>

1) 東京工科大学大学院 バイオ・情報メディア研究科 コンピュータサイエンス専攻 青木研究室

2) 東京工科大学 コンピュータサイエンス学部 コンピュータサイエンス学科 青木研究室

本研究では、深層学習モデルの敵対的攻撃に対する新しい防御手法「適応的マルチチャンネル選択法 (AMCS)」を提案する。AMCSは、小さいカーネルと大きいカーネルの畳み込み層を組み合わせることでモデルは多様な特徴を抽出し、頑健性を向上させることに成功した。

敵対的攻撃、適応的マルチチャンネル選択法、深層学習モデル

## 1. はじめに

深層学習技術の発展により、自動運転車や医療画像診断システムなど、高度なセキュリティが要求されるアプリケーションの早期普及が期待されている。しかし、これらのシステムは一般に敵対的攻撃に対して脆弱であり、その影響は深刻である[1]。敵対的攻撃は、システムを悪用することを目的として行われる手法であり、特にAIモデルの認識を混乱させることを狙っている[2]。

本研究ではこの課題に対処するため、CNN (Convolutional Neural Network) の構造そのものに着目した新たなアプローチを提案する。

## 2. 予備実験

本章では、CNNの構造が敵対的攻撃に対する堅牢性に与える影響を分析する2つの実験について述べる。

### 2.1 チャンネル影響分析

本実験ではResNetを使用し、Animals-10データセット(10のカテゴリに属する約28,000枚)に対してFGSM[3]、PGD[4]、LAP[5]の3種類の攻撃を適用した。敵対的サンプルの特定チャンネルをクリーンサンプルのチャンネルに最大3チャンネル置換し、分類精度の評価を行った。

表1は、ResNetにおける入力層でのチャンネル置換の影響を示しており、3チャンネルの置換で最大

19.2%の精度向上が観察された。

表1 ResNet18におけるチャンネル置換の影響

指標	FGSM	PGD	LAP
クリーンサンプル精度	81.3%	81.3%	81.3%
敵対的サンプル精度	70.0%	54.2%	45.3%
置き換え後の精度	77.3%	73.4%	60.4%
置き換えたチャンネル	14, 18, 41	14, 18, 41	14, 18, 26

### 2.2 カーネルサイズの影響分析

本実験では、ResNetの入力層のカーネルサイズのみを変更し、 $3 \times 3$  (オリジナル)、 $5 \times 5$ 、 $7 \times 7$ 、 $9 \times 9$ の4種類のサイズで比較を行った。

表3に示している通り、カーネルサイズの増大に伴い敵対的サンプルに対する精度が向上する一方で、クリーンサンプルの精度が低下する傾向が観察され、 $7 \times 7$ のカーネルサイズが堅牢性と全体的な分類性能の最適なバランスを示した。

表2 ResNet18におけるカーネルサイズによる精度の変化

入力層のカーネルサイズ	クリーンサンプル精度	敵対的サンプル精度
$3 \times 3$	81.3%	54.2%
$5 \times 5$	78.9%	57.0%
$7 \times 7$	75.4%	61.2%
$9 \times 9$	70.1%	61.3%

### 2.3 考察

これらの実験結果から、以下の重要な知見が得られた。

第1に、入力層での特徴抽出が敵対的攻撃に対

する堅牢性に決定的な役割を果たしていることが明らかになった。第2に、特定のチャンネルが堅牢性に対して特に脆弱であり、これらのチャンネルを適切に処理することで大幅な改善が見込めることがわかった。第3に、カーネルサイズの適切な選択により、堅牢性を向上させられる可能性が示された。

これらの知見は、CNNの構造自体を最適化することで、従来の防御手法を補完し、より効果的な敵対的攻撃対策を実現できる可能性を示唆している。

### 3. 提案手法

本研究で提案する適応的マルチチャンネル選択法 (AMCS) は、2.3で得られた知見に基づき、異なるサイズのカーネルを組み合わせた特徴抽出と、バッチサイズの動的調整を行うことにある。

AMCSは、デュアルカーネル入力層 ( $3 \times 3$  と  $7 \times 7$  の畳み込み層)、バッチサイズ調整機構、特徴統合層、そして標準的なCNNアーキテクチャから構成される。大きいバッチサイズは堅牢性の向上に、小さいカーネルサイズは精度向上に寄与する。

図1に示すAMCSの構造はこれらの要素を統合することで、従来のトレードオフを克服し、実用的な防御手法の確立が期待される。

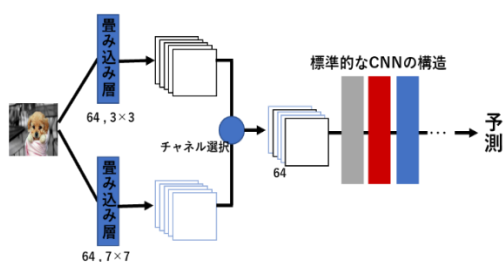


図1 提案手法

## 4. 実験

本章では、提案手法であるAMCSの有効性を検証するために行った実験の結果について述べる。

### 4.1 実験条件

実験ではResNet, EfficientNet, WideResNeの3つのモデルを使用し、従来のモデルとAMCSを加えた提案手法の比較を行う。攻撃にはFGSM, PGD, LAPの3種類を使用し、Animals-10データセット

を使用して実験を行う。

### 4.2 実験結果

表4に各モデルと攻撃手法の組み合わせにおける分類精度を示す。結果としてAMCSを適用したモデルは、全ての攻撃手法に対して従来モデルよりも高い堅牢性を示した。クリーンサンプルに対する精度の変化は2%程度に抑えられた。

表3 チャンネル交換による各攻撃手法の精度の変化と交換したチャンネル

モデル	クリーンサンプル精度	FGSM精度	PGD精度	LAP精度
ResNet18	81.3%	53.1%	48.3%	49.5%
AMCS-ResNet18	81.0%	69.8%	70.1%	65.2%
EfficientNet-B0	75.2%	48.3%	40.4%	39.9%
AMCS-EfficientNet-B0	77.8%	60.5%	51.4%	49.8%
WideResNet-50-2	79.8%	59.7%	57.1%	61.4%
AMCS-WideResNet-50-2	78.4%	66.2%	60.3%	68.9%

## 5. 結論

実験結果から、AMCSを適用したモデルは、従来のモデルと比較して敵対的攻撃に対する堅牢性が大幅に向上することが確認された。今後は、チャンネル選択アルゴリズムの改良やカーネルサイズの最適化など、AMCSの性能を更に高めることが重要な課題である。

## 参考文献

- 田籠 照博. "AIセキュリティと敵対的サンプルの脅威." NRIセキュアブログ, NRIセキュア, 2021年8月25日, <https://www.nri-secure.co.jp/blog/hostile-sample>. (2024年7月26日参照)
- 松永, スキルアップ AI. "機械学習モデルへの敵対的攻撃とは【スキルアップ AI キャンプ勉強記録】." スキルアップ AI Journal, スキルアップ AI, 2023年3月28日, <https://www.skillupai.com/blog/tech/adversarial-attacks/>. (2024年7月26日参照)
- Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and Harnessing Adversarial Examples." arXiv, 2014, <https://arxiv.org/abs/1412.6572>. (2024年7月26日参照)
- Madry, Aleksander, et al. "Towards Deep Learning Models Resistant to Adversarial Attacks." ICLR, 2018, <https://openreview.net/forum?id=rJzIBfZAb>. (2024年7月26日参照)
- Laidlaw, Cassie, and Soheil Feizi. "Laplacian Pyramid Adversarial Examples." (NeurIPS), 2019, <https://arxiv.org/abs/1905.03345>. (2024年7月26日参照)