# A3: Training robust neural networks

Alexandros Kouvatseas | Luka Lafaye de Micheaux | Matthieu Neau

# Attack Techniques

- PGD - Projected Gradient Descent

- FGSM - Fast Gradient Sign Method

- MIM - Momentum Iterative Method

# PGD – Projected Gradient Descent

**Algorithm 1** PGD Attack

1: **Input:** Model $model$, images $X$, labels $Y$
2: **Parameters:** Perturbation $\epsilon$, step size $\alpha$, iterations $N$, device $device$
3: **Output:** Perturbed images $X'$
4:
5: $X' \leftarrow X$
6: Enable gradient computation for $X'$
7: **for** $i \leftarrow 1$ to $N$ **do**
8: $\quad outputs \leftarrow model(X')$
9: $\quad loss \leftarrow F.nll\_loss(outputs, Y)$
10: $\quad$ Reset gradients: $model.zero\_grad()$
11: $\quad$ Compute gradients: $loss.backward()$
12: $\quad X' \leftarrow X' + \alpha \cdot \text{sign}(X'.grad)$
13: $\quad$ Clip $X'$ within $[X - \epsilon, X + \epsilon]$ and $[0, 1]$
14: $\quad$ Detach $X'$ from the current graph
15: $\quad$ Re-enable gradient computation for $X'$
16: **end for**
17: **return** $X'$

**Strengths**

- Theoretically grounded in constrained optimization
- Iterative Refinement

**Weaknesses**

- Computationally intensive
- Sensitive to hyperparameters
- Overfitting risk

# FGSM

**Algorithm 2** FGSM attack

1: **Input:** Neural network $model$, images $X$, labels $Y$
2: **Parameters:** Perturbation $\epsilon$, computation device $device$
3: **Output:** Perturbed images $X'$
4:
5: Enable gradient computation for $X$
6: $outputs \leftarrow model(X)$
7: $loss \leftarrow F.nll\_loss(outputs, Y)$
8: $model.zero\_grad()$
9: Compute gradients: $loss.backward()$
10: $X' \leftarrow X + \epsilon \cdot \text{sign}(X.grad)$
11: Clip $X'$ to be within valid pixel range $[0, 1]$
12: **return** $X'$

**Strengths**

- Fast with a single step
- Simple to implement

**Weaknesses**

- Sensitive to ε
- Mainly designed for $\ell\infty$ bounded perturbations

# MIM

**Algorithm 3** MIM Attack

1: **Input:** Model $model$, images $X$, labels $Y$
2: **Parameters:** Max perturbation $\epsilon$, step size $\alpha$, iterations $N$, momentum $\mu$, device $device$
3: **Output:** Adversarially perturbed images $X'$
4:
5: Initialize $g \leftarrow \mathbf{0}$ (same shape as $X$)
6: **for** $i \leftarrow 1$ to $N$ **do**
7:      Enable gradient computation for $X'$
8:      $outputs \leftarrow model(X')$
9:      $loss \leftarrow F.nll\_loss(outputs, Y)$
10:      Compute gradients: $grad \leftarrow$ autograd.grad$(loss, X')[0]$
11:      Normalize gradients: $grad\_norm \leftarrow$ torch.norm$(grad.view(grad.shape[0], -1), p = 1, dim = 1)$
12:      $grad\_normalized \leftarrow grad/(grad\_norm.view(-1, 1, 1, 1) + 1e - 8)$
13:      Update momentum: $g \leftarrow \mu \cdot g + grad\_normalized$
14:      $X' \leftarrow X' + \alpha \cdot \text{sign}(g)$
15:      Clip change: $delta \leftarrow$ torch.clamp$(X' - X, min = -\epsilon, max = \epsilon)$
16:      Clip $X'$: $X' \leftarrow$ torch.clamp$(X + delta, min = 0, max = 1)$
17:      Detach $X'$ from computation graph
18: **end for**
19: **return** $X'$

## Strengths

- Incorporates momentum to stabilize updates
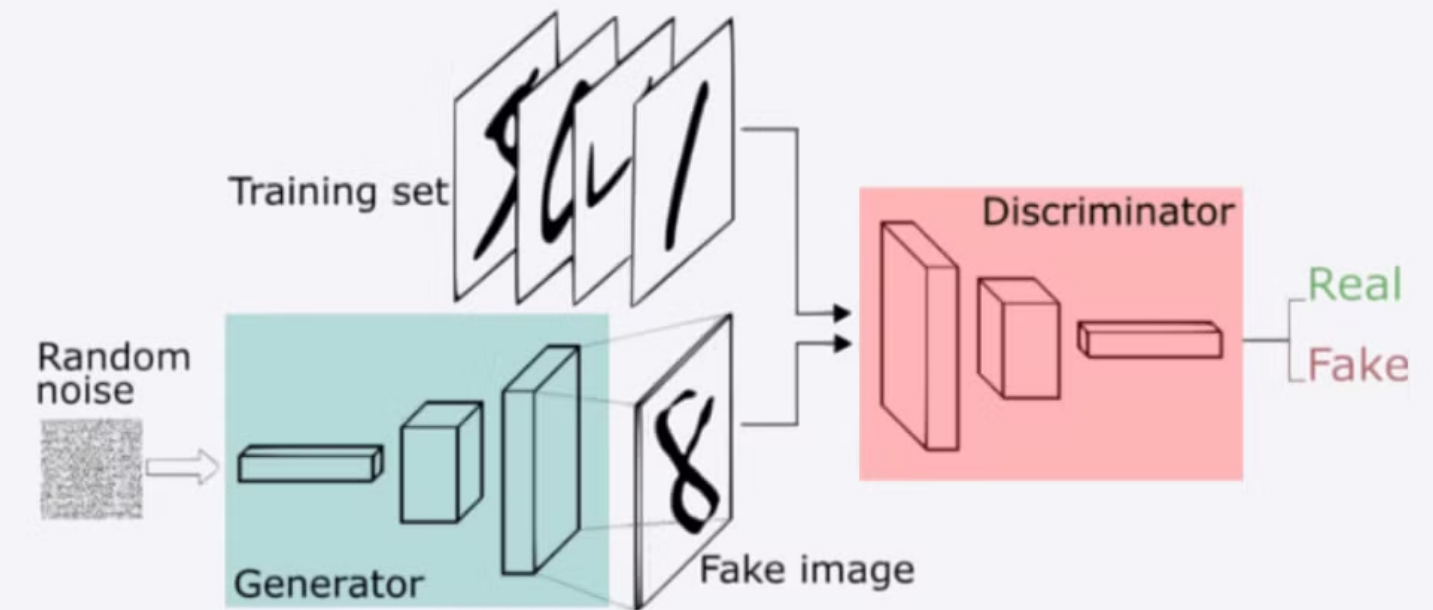- Reduces oscillations in gradient-based optimization

## Weaknesses

- Computationally intensive
- Sensitive to momentum decay

# Defense Techniques

- Adversarial Training - enhances resistance to attacks but reduces clean accuracy

- Regularization Techniques - maintains stable predictions under small perturbations

- Multi-attack technique (adversarial training using all attacks)



$$\mathcal{L}_{\text{total}} = \mathcal{L}(f(x), y) + \lambda \|\nabla_x \mathcal{L}(f(x), y)\|_p^2$$

$\mathcal{L}(f(x), y)$: Original loss function.

$\|\nabla_x \mathcal{L}(f(x), y)\|_p^2$: Regularization term penalizing large gradients

# Adversarial Training

- Trains the model on modified data points known as adversarial examples.

- **Objective of Adversarial Examples**: Introduce errors into the model's predictions, aiming to maximize the prediction error during training to make the model more robust.

- **Balanced Training Approach**: Adjusts the training process to include a mix of both natural and adversarially altered inputs, with the aim of minimizing overall prediction errors and enhancing model resilience.

# Regularization Techniques (1)

- **Spectral Normalization**: This technique adjusts each layer's weights by dividing them by their largest singular value, effectively moderating the layer's sensitivity to input perturbations.
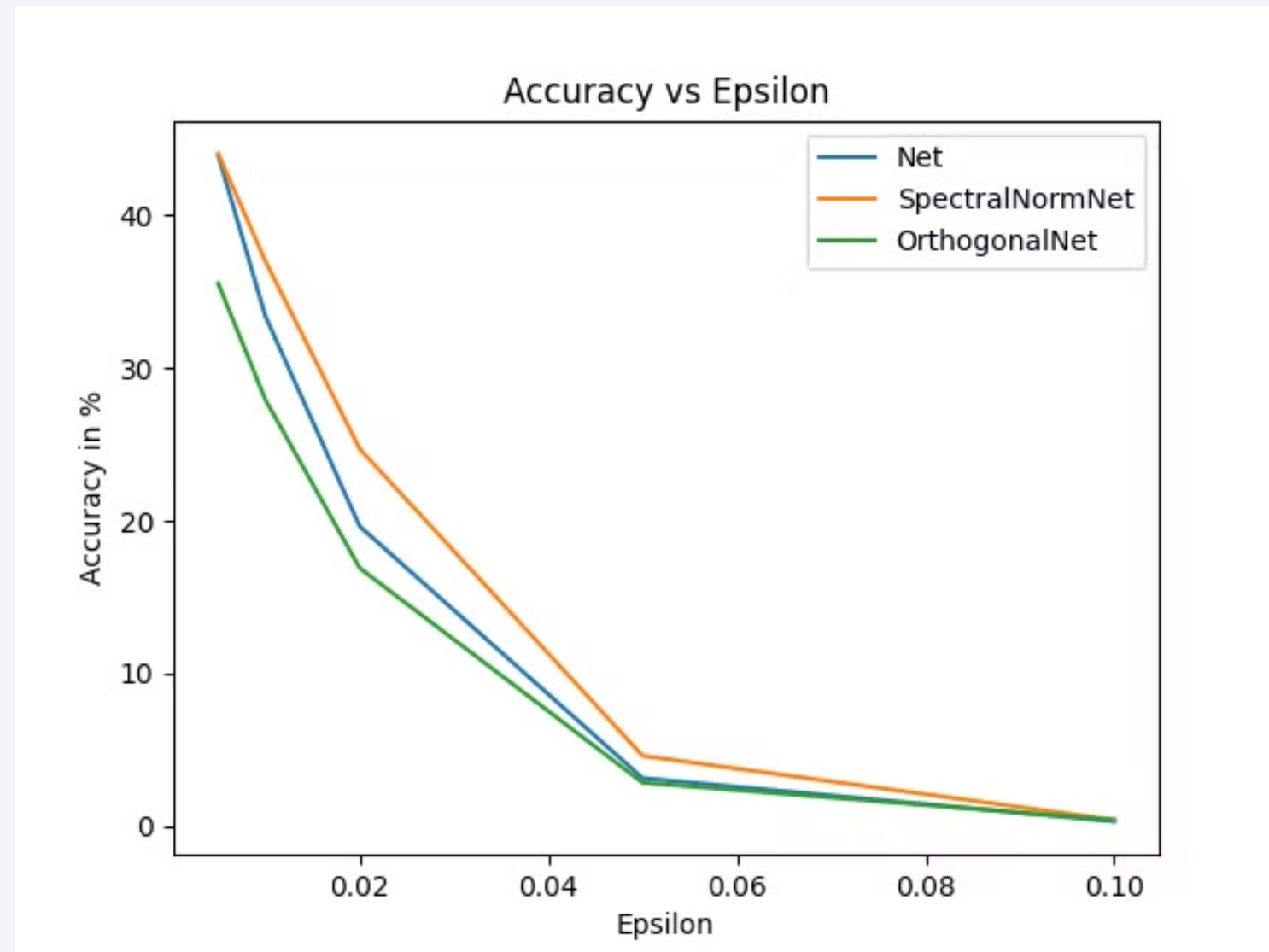
$$\sigma(a) = \max_{\mathbf{h}:\mathbf{h}\neq 0} \frac{\|A\mathbf{h}\|_2}{\|\mathbf{h}\|_2} = \max_{\|\mathbf{h}\|_2 \leq 1} \|A\mathbf{h}\|_2 \qquad \bar{W}_{\mathrm{SN}}(W) = W/\sigma(W)$$

- **Orthogonal Normalization**: Attempts to maintain weight matrices close to orthogonality to stabilize learning, though found less effective for tasks unrelated to disentangling latent spaces.
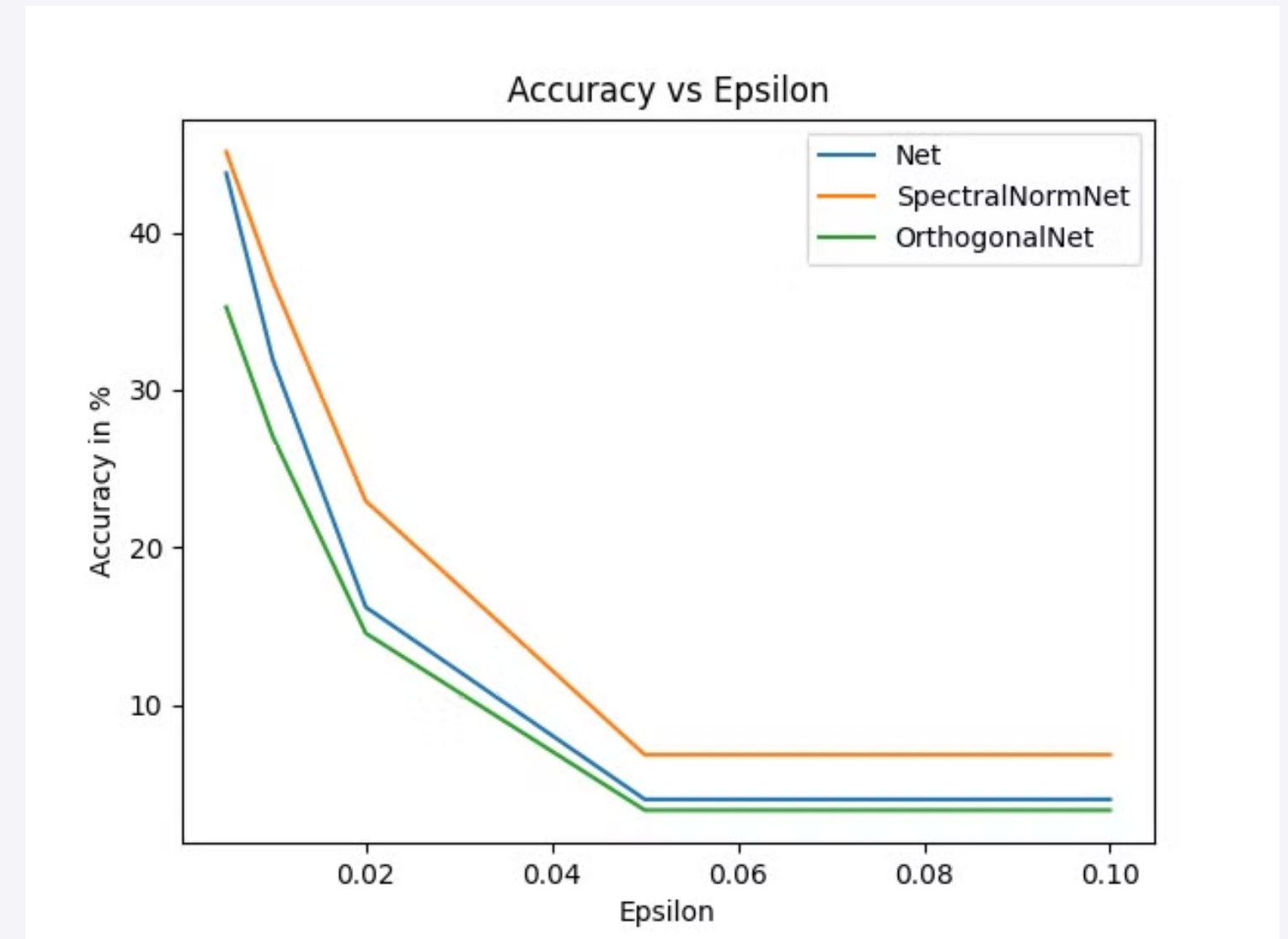
$$\mathcal{L}_{ortho} = \sum(|WW^T - I|)$$

# Regularization Results (1)



FGSM attacks



PGD attacks, alpha = 0.001, n_iter = 40

# Regularization Techniques (2)

- **Gradient Regularization:** Combines adversarial training with a gradient penalty to counteract the effects of input perturbations, improving model robustness against adversarial attacks

$$P = \frac{1}{N} \sum_{i=1}^{N} \left( \left\| \nabla_{\mathbf{x}'_i} L \right\|_2 \right)^2$$

# Multi-Attack

- **Multi-Attack Technique**: Multiple adversarial attacks (FGSM, MIM, PGD) are applied in a cyclical manner during training, enhancing model robustness by exposing it to varied perturbations.

- **Cyclical Application**: Each training batch applies a different attack based on the batch index, ensuring even exposure to all attack types throughout the training epochs.

- **Diverse Input Utilization**: Utilizing adversarial examples from different attacks as inputs for training prevents model overfitting and contributes to superior performance by increasing input diversity.

# Results

| Attack | Defense | Nat Accuracy | PGD linf | PGD l2 |
|--------|---------|--------------|----------|--------|
| PGD | Adversarial Training | 56.25% | 17.26% | 26.4% |
| FGSM | Adversarial Training | 50% | 0.46% | 8.48% |
| MIM | Adversarial Training | 37.5% | 17.23% | 26.36% |
| PGD | Gradient Regularization | 25% | 27.43% | 27.01% |
| PGD/MIM/FGSM | Multi-Attack | 43.75% | 28.39% | 34.87% |

# Thank you for listening!