





ADVANCE ALL MARCH EVERLASTING

DEFENSE STRATEGIES

31

- ▶ Adversarial Training
- ▶ Distillation
- ▶ Regularisation (Dropout, Weight Decay etc., Label Smoothing)
- ▶ Ensemble
- ▶ Virtual Adversarial Training

## DEFENSE: ADVERSARIAL TRAINING

- ▶ Modified loss function:

$$\tilde{J}(\theta, x, y) = \alpha J(\theta, x, y) + (1 - \alpha) J(\theta, x + \epsilon \text{sign}(\nabla_x J(\theta, x, y)))$$

- ▶ Augment the adversarial examples into training dataset.

---

**Algorithm 1** Adversarial training of network  $N$ .

Size of the training minibatch is  $m$ . Number of adversarial images in the minibatch is  $k$ .

---

- 1: Randomly initialize network  $N$
  - 2: **repeat**
  - 3:   Read minibatch  $B = \{X^1, \dots, X^m\}$  from training set
  - 4:   Generate  $k$  adversarial examples  $\{X_{adv}^1, \dots, X_{adv}^k\}$  from corresponding clean examples  $\{X^1, \dots, X^k\}$  using current state of the network  $N$
  - 5:   Make new minibatch  $B' = \{X_{adv}^1, \dots, X_{adv}^k, X^{k+1}, \dots, X^m\}$
  - 6:   Do one training step of network  $N$  using minibatch  $B'$
  - 7: **until** training converged
-

## DEFENCE STRATEGIES

- ▶ Adversarial Training
- ▶ Distillation
- ▶ Regularisation (Dropout, Weight Decay etc., Label Smoothing)
- ▶ Ensemble
- ▶ Virtual Adversarial Training