

ADVANCE ALL MARCH EVERLASTING

ATTACK MODELS: L-BFGS

- ▶ Targeted attack proposed by Christian Szegedy et al.
- ▶ This is a hard problem. So instead they solves following problem using box constrained L-BFGS Method



Szegedy, C, Zaremba, M, Sutskever, I, Bruna, J, Erhan, D, Goodfellow, I, & Fergus, R. (2013). Intriguing properties of neural networks.

$$\arg \min_{||\delta||_2} \text{ s.t. } \mathcal{F}(x+\delta) = \ell \text{ \& } x+\delta \in [0,1]^m$$

$$\text{minimize } c\|\delta\|_2 + \mathcal{L}(x+\delta, \ell) \text{ subject to } x+\delta \in [0, 1]^n$$

ATTACK MODELS: FGSM

- ▶ Proposed by Goodfellow et al.
- ▶ Uses L_∞ distance metric
- ▶ General Attack Model

$$x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x, y, \theta))$$

- ▶ Targeted version

$$x_{adv} = x - \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(x, \ell, \theta))$$

- ▶ Fast but does not guarantee optimal or closest adversary

ATTACK MODELS: L-BFGS

- ▶ Targeted attack proposed by Christian Szegedy et al.

$$\arg \min_{||\delta||_2} \text{ s.t. } \mathcal{F}(x + \delta) = \ell \ \& \ x + \delta \in [0, 1]^m$$

- ▶ This is a hard problem. So instead they solve following problem using box constrained L-BFGS Method

$$\text{minimize } c||\delta||_2 + \mathcal{L}(x + \delta, \ell) \text{ subject to } x + \delta \in [0, 1]^m$$