

ADVANCE ALL MARCH NEVER LEARN NG

ATTACK MODELS: FGS M

- ▶ Proposed by Goodfellow et al.
- ▶ Uses L_∞ distance metric
- ▶ General Attack Model
- ▶ Targeted version
- ▶ Fast but does not guarantee optimal or closest adversary

18

Goodfellow, J., Shlens, J., & Szegedy, C. (2014). Explaining adversarial examples.

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

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

ATTACK MODELS: FGSM^K OR PGD

- ▶ This is an iterative version of FGSM:

$$x_{adv}^0 = x; \quad x_{adv}^{t+1} = Clip_{x,\epsilon} \{ x_{adv}^t + \alpha \cdot sign(\nabla_x \mathcal{L}(x, \theta, y)) \}$$

- ▶ Targeted version:

$$x_{adv}^0 = x; \quad x_{adv}^{t+1} = Clip_{x,\epsilon} \{ x_{adv}^t - \alpha \cdot sign(\nabla_x \mathcal{L}(x, \theta, \ell)) \}$$

- ▶ Madry et al. proposes PGD as universal first order attack method
- ▶ This means defence against this attack would guarantee defence against all gradient based method

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