

ADVANCE ALL MARCH EVERLASTING

ATTACK MODELS: CIPHER ATTACKS

- ▶ Nicholas Carlini and David Wagner proposer three attack models with L_0 , L_2 and L_∞ distance
- ▶ L_2 attack is most optimal and broke all the existing defences
- ▶ They defined the problem as following:

2

2

Cardini, N., & Wagner, D. (2017, May). Towards evaluating the robustness of neural networks.

minimize $\mathcal{D}(x, x + \delta)$

such that $\mathcal{C}(x + \delta) = \ell$

$x + \delta \in [0, 1]^n$

\mathcal{D} can be L_0, L_2, L_∞ distance

ATTACK MODELS: CARLINI WAGNER L2 ATTACK

- ▶ This is a hard problem. So they defined f such that:

$$\mathcal{C}(x + \delta) = \ell \text{ if and only if } f(x + \delta) \leq 0$$

- ▶ So the optimization problem becomes:

$$\text{minimize } \mathcal{D}(x, x + \delta) + c \cdot f(x + \delta) \quad \text{s.t. } x + \delta \in [0, 1]^n$$

- ▶ Variable replacement

$$\text{They replaced } \delta_i \text{ as } \delta_i = \frac{1}{2}(\tanh(w_i) + 1) - x_i$$

- ▶ So the optimization problem becomes

$$\text{minimize } \left\| \frac{1}{2}(\tanh(w) + 1) - x \right\|_2^2 + c \cdot f\left(\frac{1}{2}(\tanh(w) + 1)\right)$$

with f defined as

$$f(x') = \max(\max\{\mathcal{H}(x')_i : i \neq l\} - \mathcal{H}(x')_l, -\kappa).$$

- ▶ They solve optimization using multiple starting point gradient descent

ATTACK MODELS: CARLINI WAGNER ATTACKS

- ▶ Nicholas Carlini and David Wagner proposer three attack models with L_0 , L_2 and L_∞ distance
- ▶ L_2 attack is most optimal and broke all the existing defences
- ▶ They defined the problem as following:

$$\begin{aligned} &\text{minimize } \mathcal{D}(x, x + \delta) \\ &\text{such that } \mathcal{C}(x + \delta) = \ell \\ &\quad x + \delta \in [0, 1]^n \end{aligned}$$

\mathcal{D} can be L_0, L_2, L_∞ distance