





ADVANCE ALL MARCH EVERLASTING

# PRACTICAL BLACKBOX ATTACK

- ▶ Papernot et al. proposed a practical black box attack on CNN.
- ▶ They first train a substitute DNN on target classifier (oracle  $\tilde{o}$ )
- ▶ They used Adversarial Examples crafted on substitute DNN to attack oracle.



2

4

**Algorithm 1 - Substitute DNN Training:** for oracle  $\tilde{O}$ , a maximum number  $max_\rho$  of substitute training epochs, a substitute architecture  $F$ , and an initial training set  $S_0$ .

---

**Require:**  $\tilde{O}$ ,  $max_\rho$ ,  $S_0$ ,  $\lambda$

- 1: Define architecture  $F$
  - 2: **for**  $\rho \in 0 \dots max_\rho - 1$  **do**
  - 3:     *// Label the substitute training set*
  - 4:      $D \leftarrow \{(\vec{x}, \tilde{O}(\vec{x})) : \vec{x} \in S_\rho\}$
  - 5:     *// Train  $F$  on  $D$  to evaluate parameters  $\theta_F$*
  - 6:      $\theta_F \leftarrow \text{train}(F, D)$
  - 7:     *// Perform Jacobian-based dataset augmentation*
  - 8:      $S_{\rho+1} \leftarrow \{\vec{x} + \lambda \cdot \text{sign}(J_F[\tilde{O}(\vec{x})]) : \vec{x} \in S_\rho\} \cup S_\rho$
  - 9: **end for**
  - 10: **return**  $\theta_F$
-

Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z. B., & Swami, A. (2017, April). Practical black-box attacks against machine learning.



## ATTACK MODELS : TRANSFERABILITY

- ▶ Transferability : Adversarial Examples crafted using one type of ML models can be used to attack other types of models
- ▶ This phenomena is observed because of similarity in decision boundaries of various ML models
- ▶ In The Space of Transferable Adversarial Examples F. Tramèr et al. analyse adversarial subspace.
- ▶ This is very important in devising Black Box Attacks

# PRACTICAL BLACK BOX ATTACK

- ▶ Papernot et al. proposed a practical black box attack on CNN.
- ▶ They first train a substitute DNN on target classifier (oracle  $\tilde{O}$ )
- ▶ They used Adversarial Examples crafted on substitute DNN to attack oracle.

---

**Algorithm 1 - Substitute DNN Training:** for oracle  $\tilde{O}$ , a maximum number  $max_\rho$  of substitute training epochs, a substitute architecture  $F$ , and an initial training set  $S_0$ .

---

**Require:**  $\tilde{O}$ ,  $max_\rho$ ,  $S_0$ ,  $\lambda$

- 1: Define architecture  $F$
- 2: **for**  $\rho \in 0 \dots max_\rho - 1$  **do**
- 3:     *// Label the substitute training set*
- 4:      $D \leftarrow \{(\vec{x}, \tilde{O}(\vec{x})) : \vec{x} \in S_\rho\}$
- 5:     *// Train  $F$  on  $D$  to evaluate parameters  $\theta_F$*
- 6:      $\theta_F \leftarrow \text{train}(F, D)$
- 7:     *// Perform Jacobian-based dataset augmentation*
- 8:      $S_{\rho+1} \leftarrow \{\vec{x} + \lambda \cdot \text{sign}(J_F[\tilde{O}(\vec{x})]) : \vec{x} \in S_\rho\} \cup S_\rho$
- 9: **end for**
- 10: **return**  $\theta_F$

Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z. B., & Swami, A. (2017, April). Practical black-box attacks against machine learning.