

Why Do Adversarial Attacks Transfer? Explaining Transferability of Evasion and Poisoning Attacks

Ambra Demontis[†], Marco Melis[†], Maura Pintor[†], Matthew Jagielski^{*}, Battista Biggio^{†,‡}, Alina Oprea^{*},
Cristina Nita-Rotaru^{*}, and Fabio Roli^{†,‡}

[†]Department of Electrical and Electronic Engineering, University of Cagliari, Italy

[‡]Pluribus One, Italy

^{*}Northeastern University, Boston, MA, USA

USENIX(2019)

给定攻击者对 $K \in \kappa$ 的了解, 攻击样本 $x' \in \Phi(x)$ 及其标签 y , 攻击者的目标可以用目标函数 $A(x', y, \kappa) \in R$ (例如损失函数) 来定义

$$\mathbf{x}^* \in \arg \max_{\mathbf{x}' \in \Phi(\mathbf{x})} \mathcal{A}(\mathbf{x}', y, \kappa).$$

Algorithm 1 Gradient-based Evasion and Poisoning Attacks

Input: \mathbf{x}, y : the input sample and its label; $\mathcal{A}(\mathbf{x}, y, \kappa)$: the attacker's objective; $\kappa = (\mathcal{D}, \mathcal{X}, f, \mathbf{w})$: the attacker's knowledge parameter vector; $\Phi(\mathbf{x})$: the feasible set of manipulations that can be made on \mathbf{x} ; $t > 0$: a small number.

Output: \mathbf{x}' : the adversarial example.

- 1: Initialize the attack sample: $\mathbf{x}' \leftarrow \mathbf{x}$
 - 2: **repeat**
 - 3: Store attack from previous iteration: $\mathbf{x} \leftarrow \mathbf{x}'$
 - 4: Update step: $\mathbf{x}' \leftarrow \Pi_{\Phi}(\mathbf{x} + \eta \nabla_{\mathbf{x}} \mathcal{A}(\mathbf{x}, y, \kappa))$, where the step size η is chosen with line search (bisection method), and Π_{Φ} ensures projection on the feasible domain Φ .
 - 5: **until** $|\mathcal{A}(\mathbf{x}', y, \kappa) - \mathcal{A}(\mathbf{x}, y, \kappa)| \leq t$
 - 6: **return** \mathbf{x}'
-

在规避攻击中，攻击者操纵测试样本以使其错误分类，即逃避学习算法的检测。

white-box:

$$\begin{aligned} \max_{\mathbf{x}'} \quad & \ell(y, \mathbf{x}', \mathbf{w}), \\ \text{s.t.} \quad & \|\mathbf{x}' - \mathbf{x}\|_p \leq \varepsilon, \\ & \mathbf{x}_{\text{lb}} \preceq \mathbf{x}' \preceq \mathbf{x}_{\text{ub}}, \end{aligned}$$

black-box:

对于黑盒案例，使用代理模型 \hat{f} 的参数 $\hat{\mathbf{w}}$ 就足够了。

假设有一个来自于替代模型的攻击样本 $x' = x + \hat{\delta}$ ，将这个攻击样本迁移到目标模型，迁移性采用损失函数进行刻画：

$$T = \ell(y, \mathbf{x} + \hat{\delta}, \mathbf{w})$$

进行线性化近似，可以把 $\hat{\delta}$ 拆出来：

$$T = \ell(y, \mathbf{x} + \hat{\delta}, \mathbf{w}) \cong \ell(y, \mathbf{x}, \mathbf{w}) + \hat{\delta}^\top \nabla_{\mathbf{x}} \ell(y, \mathbf{x}, \mathbf{w})$$

目标模型损失函数关于样本的梯度，梯度越大模型越容易受到迁移攻击

$$S(\mathbf{x}, y) = \|\nabla_{\mathbf{x}} \ell(y, \mathbf{x}, \mathbf{w})\|_q,$$

此处的q范数是扰动约束p范数的对偶范数

在黑盒情况下，攻击者掌握的替代模型和目标模型之间的区别：

$$R(\mathbf{x}, y) = \frac{\nabla_{\mathbf{x}} \hat{\ell}^{\top} \nabla_{\mathbf{x}} \ell}{\|\nabla_{\mathbf{x}} \hat{\ell}\|_2 \|\nabla_{\mathbf{x}} \ell\|_2}$$

表征对抗攻击的可转移性，余弦相似度越高，迁移性越好

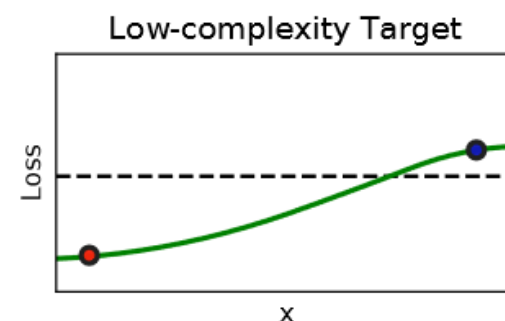
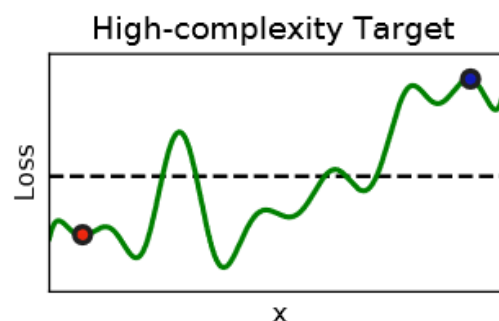
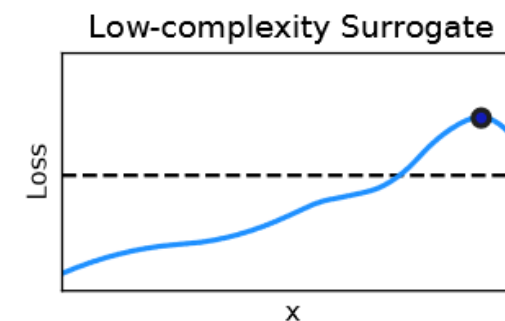
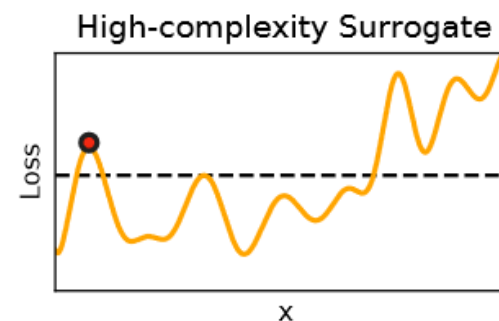
Variability of the Loss Landscape



更稳定且方差更低的替代模型的损失函数往往有利于基于梯度的优化攻击，以找到更好的局部最优

$$V(\mathbf{x}, y) = \mathbb{E}_{\mathcal{D}}\{\ell(y, \mathbf{x}, \hat{\mathbf{w}})^2\} - \mathbb{E}_{\mathcal{D}}\{\ell(y, \mathbf{x}, \hat{\mathbf{w}})\}^2,$$

E_D 是对不同替代采用的训练集的期望。



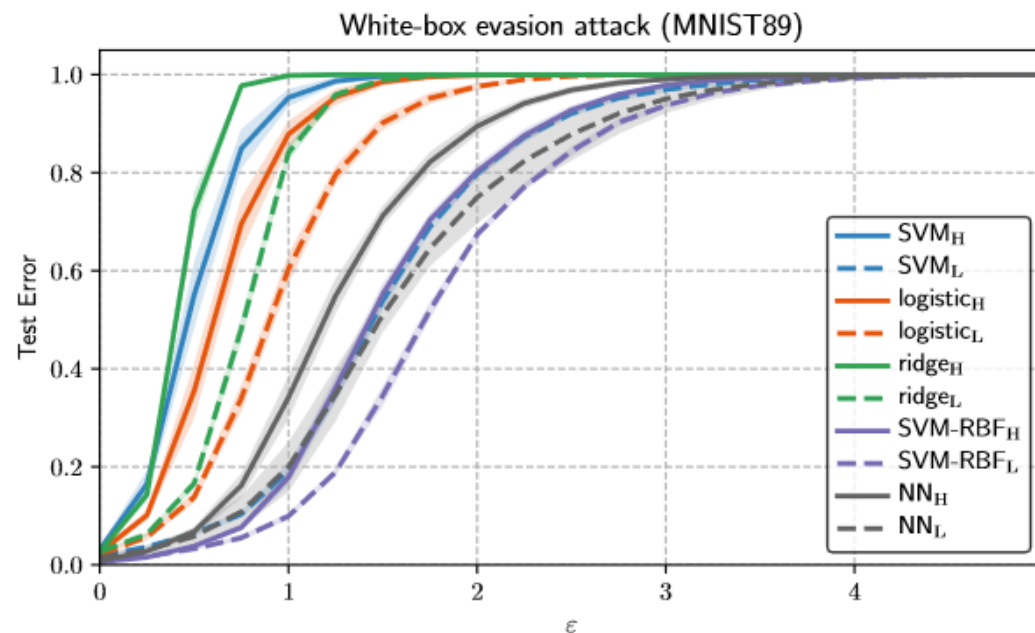


Figure 5: White-box evasion attacks on MNIST89. Test error against increasing maximum perturbation ϵ .

