

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

Ambra Demontis¹, Marco Melis¹, Maura Pintor^{1,3}, Matthew Jagielski², Battista Biggio^{1,3}, Alina Oprea², Cristina Nita-Rotaru², Fabio Roli^{1,3}

1. University of Cagliari, Italy

2. Northeastern University, Boston, MA, USA

3. Pluribus One, Italy - <http://www.pluribus-one.it>

Pluribus One



UNIVERSITÀ
DEGLI STUDI
DI CAGLIARI



Northeastern
University

Introduction

Motivation

Previous works [1, 2, 3] have shown empirical evidence that adversarial attacks can transfer between models, however there is little understanding on the underlying reasons of this phenomenon.

Goals

- Formal definition of transferability
- Investigation on test-time (evasion) and training-time (poisoning) attacks
- Definition of new metrics for understanding when and why adversarial attacks transfer

Evasion and Poisoning Attacks

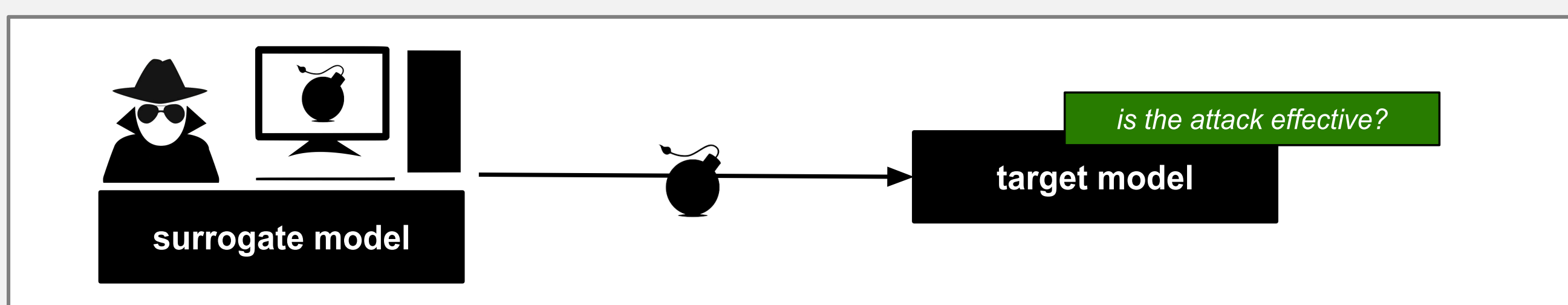
Evasion attacks [4]: small manipulations of testing data points that result in misprediction at testing time on those points

$$\begin{aligned} \max_{\mathbf{x}'} \quad & \ell(y, \mathbf{x}', \mathbf{w}) \\ \text{s.t.} \quad & \|\mathbf{x}' - \mathbf{x}\|_p \leq \varepsilon \\ & \mathbf{x}_{1b} \preceq \mathbf{x}' \preceq \mathbf{x}_{ub} \end{aligned}$$

Poisoning attacks [5]: the attacker controls a certain amount of training data, thus influencing the trained model and ultimately the predictions at test time on the testing set

$$\begin{aligned} \max_{\mathbf{x}'} \quad & L(\mathcal{D}_{\text{val}}, \mathbf{w}^*) = \sum_{j=1}^m \ell(y_j, \mathbf{x}_j, \mathbf{w}^*) \\ \text{s.t.} \quad & \mathbf{w}^* \in \arg \min_{\mathbf{w}} \mathcal{L}(\mathcal{D}_{\text{tr}} \cup (\mathbf{x}', y), \mathbf{w}) \end{aligned}$$

Transfer loss



perturbation computed with the surrogate model

$$T = \underbrace{\ell(y, \mathbf{x} + \hat{\delta}, \mathbf{w})}_{\text{loss attained by the target model on the adversarial point crafted against the surrogate}} \cong \underbrace{\ell(y, \mathbf{x}, \mathbf{w}) + \hat{\delta}^T \nabla_{\mathbf{x}} \ell(y, \mathbf{x}, \mathbf{w})}_{\text{loss increment in the target given by the attack crafted with the surrogate model}}$$

loss attained by the target model on the adversarial point crafted against the surrogate

loss increment in the target given by the attack crafted with the surrogate model

Our work

$$S(\mathbf{x}, y) = \|\nabla_{\mathbf{x}} \ell\|_2$$

S: size of input gradients
measures the vulnerability of the target model (white-box loss increment)

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

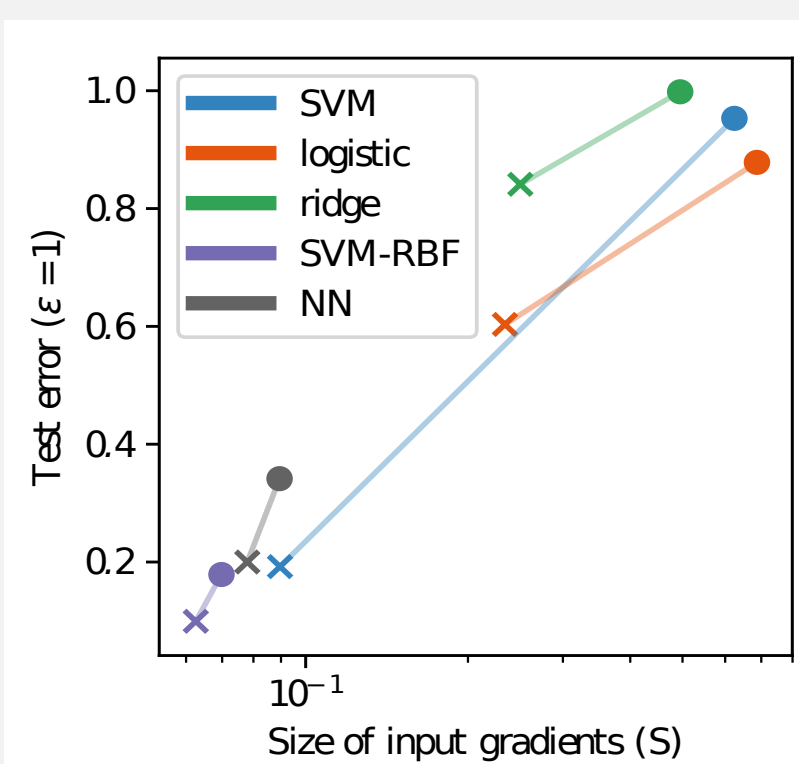
R: gradient alignment
measures black-box to white-box loss increment ratio

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

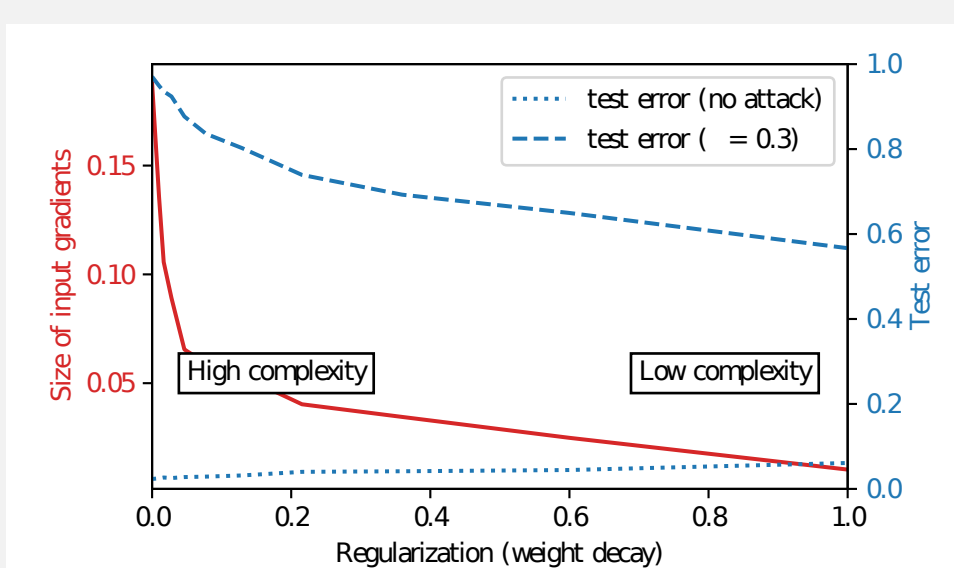
V: variability of the loss landscape

Evaluates the variability of the surrogate classifier under training data resampling

Experiments



Test error under attack vs. average size of input gradients (S) for low- (denoted with 'x') and high-complexity (denoted with 'o') classifiers

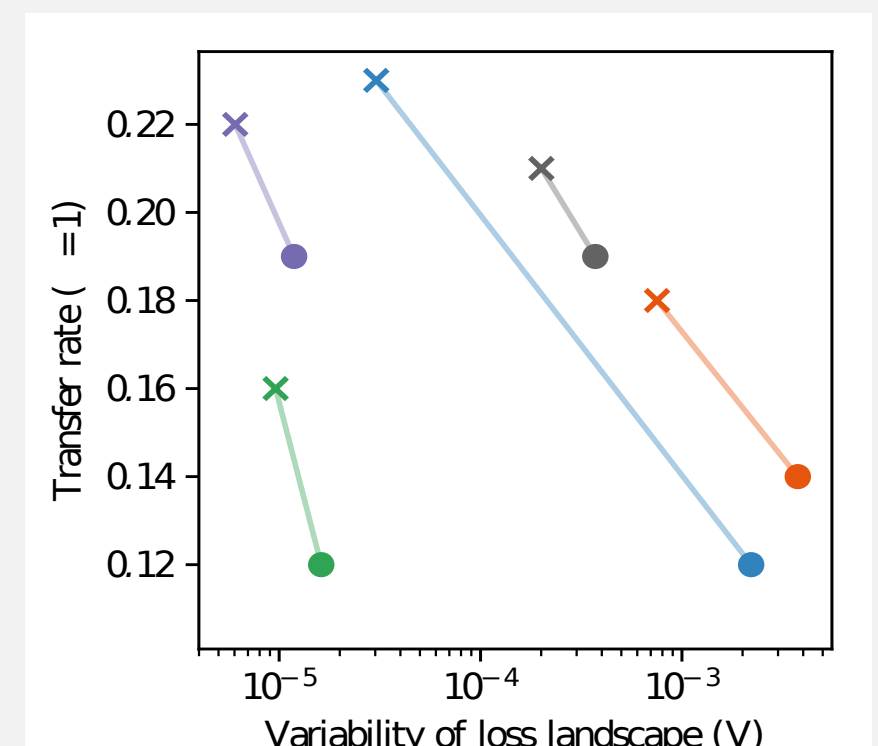


The size of the input gradients (S) is an indication of how much the target model is implicitly vulnerable against adversarial examples, including the ones crafted with the surrogate model.

SVM _H	0.14	0.35	0.19	0.29	0.13	0.25	0.26	0.32	0.28	0.32
SVM _L	0.32	0.88	0.42	0.63	0.26	0.50	0.68	0.83	0.67	0.79
logistic _H	0.18	0.45	0.25	0.37	0.18	0.32	0.35	0.42	0.36	0.41
logistic _L	0.26	0.64	0.35	0.51	0.24	0.43	0.49	0.59	0.51	0.58
ridge _H	0.12	0.26	0.16	0.23	0.18	0.28	0.21	0.25	0.21	0.24
ridge _L	0.22	0.49	0.29	0.41	0.27	0.47	0.39	0.46	0.40	0.44
SVM-RBF _H	0.25	0.69	0.33	0.50	0.21	0.40	0.67	0.75	0.58	0.66
SVM-RBF _L	0.30	0.83	0.39	0.58	0.25	0.47	0.75	0.87	0.66	0.78
NN _H	0.26	0.68	0.34	0.51	0.22	0.41	0.57	0.67	0.65	0.68
NN _L	0.30	0.81	0.39	0.58	0.24	0.46	0.67	0.79	0.68	0.80

Gradient alignment (R) and perturbation correlation (Pearson correlation coefficient between white-box and black-box perturbation) for evasion attacks on MNIST89. Rows indicate the surrogate classifier and columns indicate the target

Average transfer rate vs variability of the loss landscape (V). Low-complexity models are denoted with 'x' and high-complexity models are denoted with 'o'.



The main factors contributing to transferability are the intrinsic vulnerability of the target model and the complexity of the surrogate. The variability of the loss landscape of the surrogate under different training set resamplings can influence the stability of the solution of the optimization problem.

Key insights and future work

Evasion: decreasing complexity of the surrogate model by properly adjusting the hyperparameters of its learning algorithm provides adversarial examples that transfer better to a range of models

Poisoning: the best surrogates are generally models with similar levels of regularization as the target model

Future work: framework for understanding causes of success and failure of adversarial attacks against a model, comparing attacks with quantitative metrics in a competitive but fair environment

References and useful links

[1] Goodfellow et al., *Explaining and harnessing adversarial examples*, ICLR 2015.

[2] Szegedy et al., *Intriguing properties of neural networks*, ICLR 2014.

[3] Papernot et al., *Practical black-box attacks against machine learning*, ASIACCS 2017

Code available at: <https://gitlab.com/secml/secml>

Paper: Demontis et al., *Why do Adversarial Attacks Transfer? Explaining Transferability of Evasion and Poisoning Attacks*, USENIX 2019

[4] Biggio et al., *Evasion attacks against machine learning at test time*, ECML PKDD 2013.

[5] Biggio et al., *Poisoning attacks against support vector machines*, ICML 2012.

Download full text

