

# Adversarial Machine Learning

Sanjay Seetharaman

TRDDC



# Acknowledgements

Rosni K V, Researcher <sup>1</sup>

Manish Shukla, Researcher <sup>1</sup>

Dr. Sachin Lodha, Head, Cybersecurity and Privacy Research <sup>1</sup>

Dr. Shina Sheen, Associate Professor <sup>2</sup>

Dr. N. Rajamanickam, Assistant Professor <sup>2</sup>

Dr. R. S. Lekshmi, Professor <sup>2</sup>

Dr. R. Nadarajan, Professor and Head <sup>2</sup>

<sup>1</sup> Tata Consultancy Services

<sup>2</sup> Department of Applied Mathematics and Computational Sciences,  
PSG College of Technology

# Outline

## 1 About TCS Research

## 2 Introduction

- Arms Race
- Attacker's Goal
- Attack Strategy

## 3 Attacks

- ALFA
- PGA

# About TCS Research

- TCS established its first lab in 1981 at Pune, India. It has invested in multiple research areas for over three decades.
- Every year, our researchers publish around 300 papers for Tier 1 conferences. TCS has also filed over 3,500 patents.
- Numerous tools and frameworks have been co-created for delivering complex projects for global enterprises.
- Stakeholders include industries, organizations and academia.

# Outline

## 1 About TCS Research

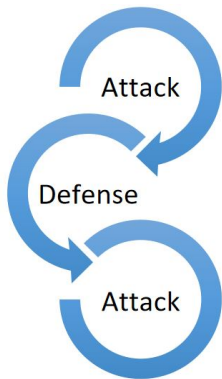
## 2 Introduction

- **Arms Race**
- Attacker's Goal
- Attack Strategy

## 3 Attacks

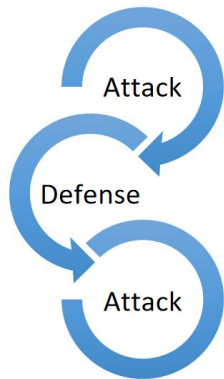
- ALFA
- PGA

# Security is an Arms Race



- Reactive Security
  - ▶ Unable to prevent the risk of never-seen-before attacks
- Proactive Security
  - ▶ Requires identification of relevant threats and development of corresponding countermeasures

# Security is an Arms Race



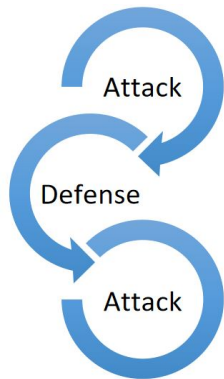
- Reactive Security

- ▶ Unable to prevent the risk of never-seen-before attacks

- Proactive Security

- ▶ Requires identification of relevant threats and development of corresponding countermeasures

# Security is an Arms Race



- Reactive Security
  - ▶ Unable to prevent the risk of `never-seen-before` attacks
- Proactive Security
  - ▶ Requires identification of `relevant` threats and development of corresponding countermeasures



# Outline

## 1 About TCS Research

## 2 Introduction

- Arms Race
- **Attacker's Goal**
- Attack Strategy

## 3 Attacks

- ALFA
- PGA

# Security Violation

- **Integrity violation:** Evade detection without compromising normal system operation
- **Availability violation:** Compromise the normal system functionalities available to legitimate users
- **Privacy violation:** Obtain private information about the system, its users or data by reverse-engineering the learning algorithm

# Attack and Error Specificity

## Attack Specificity

- **Targeted attack:** Cause misclassification of a specific set of samples/system users/services
- **Indiscriminate attack:** Target any system sample/system user/service

## Error Specificity

- **Generic:** Aim to have a sample misclassified
- **Specific:** Aim to have a sample misclassified as a specific class

# Attack and Error Specificity

## Attack Specificity

- **Targeted attack:** Cause misclassification of a specific set of samples/system users/services
- **Indiscriminate attack:** Target any system sample/system user/service

## Error Specificity

- **Generic:** Aim to have a sample misclassified
- **Specific:** Aim to have a sample misclassified as a specific class

# Outline

## 1 About TCS Research

## 2 Introduction

- Arms Race
- Attacker's Goal
- **Attack Strategy**

## 3 Attacks

- ALFA
- PGA

# Notations

Input Data:  $\mathcal{D}$

Hypothesis space:  $\mathcal{H}$

Classification hypothesis:  $f \in \mathcal{H}$

Loss function:  $\mathcal{V}$

Attacker's knowledge space:  $\Theta$

Attacker's capability:  $\Phi$

# High-level Formulation

Given the attacker's knowledge  $\theta \in \Theta$  and a set of manipulated attack samples  $\mathcal{D}' \in \Phi(\mathcal{D})$ , the attacker's goal can be defined in terms of an objective function  $\mathcal{A}(\mathcal{D}', \theta) \in \mathbb{R}$  which measures how effective the attacks  $\mathcal{D}'$  are.

The optimal attack strategy can be thus given as:

$$\mathcal{D}^* \in \arg \max_{\mathcal{D}' \in \Phi(\mathcal{D})} \mathcal{A}(\mathcal{D}', \theta)$$

# Popular Attack Scenarios

## Test-time Evasion

Manipulation of input data to evade a trained classifier at test time

$$\begin{aligned} \max_{\mathbf{x}'} \quad & \mathcal{A}(\mathbf{x}', \theta) = \Omega(\mathbf{x}') = \max_{l \neq k} f_l(\mathbf{x}) - f_k(\mathbf{x}) \\ \text{s.t.} \quad & d(\mathbf{x}, \mathbf{x}') \leq d_{\max}, \mathbf{x}_{lb} \leq \mathbf{x}' \leq \mathbf{x}_{ub} \end{aligned}$$

## Train-time Poisoning

Manipulation of train data to increase the number of misclassified samples at test time

$$\begin{aligned} \mathcal{D}^* \in \arg \max_{\mathcal{D} \in \Phi(\mathcal{D})} \quad & \mathcal{A}(\mathcal{D}', \theta) = L(\mathcal{D}_{\text{val}}, \mathbf{w}^*) \\ \text{s.t.} \quad & \mathbf{w}^* \in \arg \min_{\mathbf{w}' \in \mathcal{W}} \mathcal{L}(\mathcal{D}_{\text{tr}} \cup \mathcal{D}', \mathbf{w}') \end{aligned}$$



# Popular Attack Scenarios

## Test-time Evasion

Manipulation of input data to evade a trained classifier at test time

$$\begin{aligned} \max_{\mathbf{x}'} \quad & \mathcal{A}(\mathbf{x}', \theta) = \Omega(\mathbf{x}') = \max_{l \neq k} f_l(\mathbf{x}) - f_k(\mathbf{x}) \\ \text{s.t.} \quad & d(\mathbf{x}, \mathbf{x}') \leq d_{\max}, \mathbf{x}_{lb} \leq \mathbf{x}' \leq \mathbf{x}_{ub} \end{aligned}$$

## Train-time Poisoning

Manipulation of train data to increase the number of misclassified samples at test time

$$\begin{aligned} \mathcal{D}^* \in \arg \max_{\mathcal{D} \in \Phi(\mathcal{D})} \quad & \mathcal{A}(\mathcal{D}', \theta) = L(\mathcal{D}_{\text{val}}, \mathbf{w}^*) \\ \text{s.t.} \quad & \mathbf{w}^* \in \arg \min_{\mathbf{w}' \in \mathcal{W}} \mathcal{L}(\mathcal{D}_{\text{tr}} \cup \mathcal{D}', \mathbf{w}') \end{aligned}$$

# Outline

## 1 About TCS Research

## 2 Introduction

- Arms Race
- Attacker's Goal
- Attack Strategy

## 3 Attacks

- **ALFA**
- PGA

# Supervised Classification

## Problem

Given training instances  $S := \{(\mathbf{x}_i, y_i) \mid \mathbf{x}_i \in \mathcal{X}, y_i \in \mathcal{Y}\}_{i=1}^n$ , the goal is to find a classification hypothesis  $f_S \in \mathcal{H}$  by solving the Tikhonov regularization problem

$$f_S := \arg \min_f \gamma \sum_{i=1}^n V(y_i, f(\mathbf{x}_i)) + \|f\|_{\mathcal{H}}^2$$

# Attack Formulation

- Introduce a set of variables  $z_i \in \{0, 1\}, i = 1, \dots, n$ .
- Replace  $y_i$  with  $y'_i := y_i (1 - 2z_i)$
- Denote  $S' := \{(\mathbf{x}_i, y'_i)\}_{i=1}^n$  the *tainted* training set

$$\max_{\mathbf{z}} \sum_{(\mathbf{x}, y) \in T} V(y, f_{S'}(\mathbf{x}))$$

$$\text{s.t.} \quad f_{S'} \in \arg \min_f \gamma \sum_{i=1}^n V(y'_i, f(\mathbf{x}_i)) + \|f\|_{\mathcal{H}}^2$$

$$\sum_{i=1}^n c_i z_i \leq C, z_i \in \{0, 1\}$$

where  $c_i \in \mathbb{R}_{0+}$  is the cost (or risk) of flipping label  $y_i$  and  $C$  is the total adversarial cost.

# Attack Formulation

- Set  $U := \{(\mathbf{x}_i, y_i)\}_{i=1}^{2n}$  is constructed as follows

$$\begin{aligned}(\mathbf{x}_i, y_i) &\in \mathcal{S}, & i &= 1, \dots, n \\ \mathbf{x}_i &:= \mathbf{x}_{i-n}, & i &= n+1, \dots, 2n \\ y_i &:= -y_{i-n} & i &= n+1, \dots, 2n\end{aligned}$$

- The near-optimal label flips problem is rewritten as

$$\min_{\mathbf{q}, f} \quad \gamma \sum_{i=1}^{2n} q_i [V(y_i, f(\mathbf{x}_i)) - V(y_i, f_S(\mathbf{x}_i))] + \|f\|_{\mathcal{H}}^2$$

$$\text{s.t.} \quad \sum_{i=n+1}^{2n} c_i q_i \leq C$$

$$\begin{aligned}q_i + q_{i+n} &= 1, & i &= 1, \dots, n \\ q_i &\in \{0, 1\}, & i &= 1, \dots, 2n\end{aligned}$$

# Attack Formulation

$$\min_{\mathbf{q}, \mathbf{w}, \epsilon, b} \quad \gamma \sum_{i=1} q_i (\epsilon_i - \xi_i) + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t.} \quad y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0, \quad i = 1, \dots, 2n$$

$$\sum_{i=n+1}^{2n} c_i q_i \leq C$$

$$q_i + q_{i+n} = 1, \quad i = 1, \dots, n$$

$$q_i \in \{0, 1\}, \quad i = 1, \dots, 2n$$

# Attack Formulation - Iterative Approach

QP

$$\min_{\mathbf{w}, \epsilon, b} \quad \gamma \sum_{i=1}^{2n} q_i \epsilon_i + \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t.} \quad y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \epsilon_i, \quad \epsilon_i \geq 0, \quad i = 1, \dots, 2n$$

LP

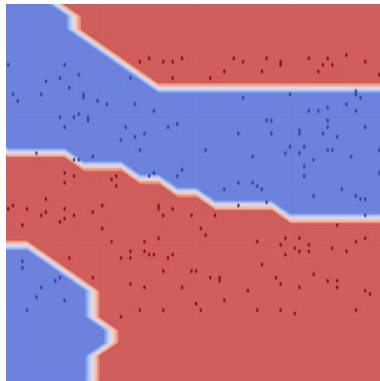
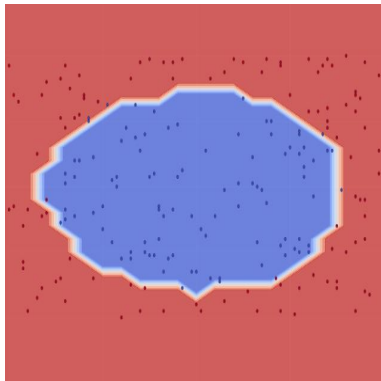
$$\min_{\mathbf{q}} \quad \gamma \sum_{i=1}^{2n} q_i (\epsilon_i - \xi_i)$$

$$\text{s.t.} \quad \sum_{i=n+1}^{2n} c_i q_i \leq C$$

$$q_i + q_{i+n} = 1, \quad i = 1, \dots, n$$

$$0 \leq q_i \leq 1, \quad i = 1, \dots, 2n$$

# Experimental Results





# Outline

## 1 About TCS Research

## 2 Introduction

- Arms Race
- Attacker's Goal
- Attack Strategy

## 3 Attacks

- ALFA
- **PGA**

# Objective Function

The attacker wants the learner's learned weight vector  $\mathbf{w}$  as close to  $\mathbf{w}^*$  as possible

$$\max_{\mathbf{z}} \quad \frac{\mathbf{w}^\top \mathbf{w}^*}{\|\mathbf{w}\| \|\mathbf{w}^*\|}$$

$$\text{s.t.} \quad f \in \arg \min_{g \in \mathcal{H}} C \sum_{i=1}^n L(y'_i, g(\mathbf{x}_i)) + \frac{1}{2} \|g\|^2$$

$$\sum_{i=1}^n z_i \leq B$$

$$y'_i = y_i (1 - 2z_i), \forall i \in [n]$$

$$z_i \in \{0, 1\}, \forall i \in [n]$$

# Objective Function

The attacker wants the learner's learned weight vector  $\mathbf{w}$  as close to  $\mathbf{w}^*$  as possible

$$\max_{\mathbf{z}} \quad \frac{\mathbf{w}^\top \mathbf{w}^*}{\|\mathbf{w}\| \|\mathbf{w}^*\|}$$

$$\text{s.t.} \quad f \in \arg \min_{g \in \mathcal{H}} C \sum_{i=1}^n L(y'_i, g(\mathbf{x}_i)) + \frac{1}{2} \|g\|^2$$

$$\sum_{i=1}^n z_i \leq B$$

$$y'_i = y_i (1 - 2z_i), \forall i \in [n]$$

$$z_i \in \{0, 1\}, \forall i \in [n]$$

- Relax binary variables  $z_i$  to interval  $[0, 1]$
- Update  $\mathbf{z}$  along its approximate gradients until convergence or the iteration limit is reached
- Project  $\mathbf{z}$  onto a  $l_{\text{inf}}$  norm ball by truncating each  $z_i$  into range  $[0, 1]$
- Further project  $\mathbf{z}$  onto a  $l_1$  norm ball with diameter  $B$

- Relax binary variables  $z_i$  to interval  $[0, 1]$
- Update  $\mathbf{z}$  along its approximate gradients until convergence or the iteration limit is reached
- Project  $\mathbf{z}$  onto a  $l_{\text{inf}}$  norm ball by truncating each  $z_i$  into range  $[0, 1]$
- Further project  $\mathbf{z}$  onto a  $l_1$  norm ball with diameter  $B$

- Relax binary variables  $z_i$  to interval  $[0, 1]$
- Update  $\mathbf{z}$  along its approximate gradients until convergence or the iteration limit is reached
- Project  $\mathbf{z}$  onto a  $l_{\text{inf}}$  norm ball by truncating each  $z_i$  into range  $[0, 1]$
- Further project  $\mathbf{z}$  onto a  $l_1$  norm ball with diameter  $B$

- Relax binary variables  $z_i$  to interval  $[0, 1]$
- Update  $\mathbf{z}$  along its approximate gradients until convergence or the iteration limit is reached
- Project  $\mathbf{z}$  onto a  $l_{\text{inf}}$  norm ball by truncating each  $z_i$  into range  $[0, 1]$
- Further project  $\mathbf{z}$  onto a  $l_1$  norm ball with diameter  $B$

# Gradient Computation

$$\nabla_{\mathbf{z}} U = \nabla_{\mathbf{w}} U \cdot \nabla_{y'} \mathbf{w} \cdot \nabla_{\mathbf{z}} y'$$

$$\frac{\partial U}{\partial w_j} = \frac{\|\mathbf{w}\|^2 w_j^* - \mathbf{w}^\top \mathbf{w}^* w_j}{\|\mathbf{w}\|^3 \|\mathbf{w}^*\|}$$

$$\frac{\partial y'_i}{\partial z_j} = -\mathbb{1}(i = j) 2y_i$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y'_i \mathbf{x}_i$$

$$\frac{\partial w_j}{\partial y'_i} = \alpha_i x_{ij}$$



# Gradient Computation

$$\nabla_{\mathbf{z}} U = \nabla_{\mathbf{w}} U \cdot \nabla_{y'} \mathbf{w} \cdot \nabla_{\mathbf{z}} y'$$

$$\frac{\partial U}{\partial w_j} = \frac{\|\mathbf{w}\|^2 w_j^* - \mathbf{w}^\top \mathbf{w}^* w_j}{\|\mathbf{w}\|^3 \|\mathbf{w}^*\|}$$

$$\frac{\partial y'_i}{\partial z_j} = -\mathbb{1}(i = j) 2y_i$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y'_i \mathbf{x}_i$$

$$\frac{\partial w_j}{\partial y'_i} = \alpha_i x_{ij}$$

# Gradient Computation

$$\nabla_{\mathbf{z}} U = \nabla_{\mathbf{w}} U \cdot \nabla_{y'} \mathbf{w} \cdot \nabla_{\mathbf{z}} y'$$

$$\frac{\partial U}{\partial w_j} = \frac{\|\mathbf{w}\|^2 w_j^* - \mathbf{w}^\top \mathbf{w}^* w_j}{\|\mathbf{w}\|^3 \|\mathbf{w}^*\|}$$

$$\frac{\partial y'_i}{\partial z_j} = -\mathbb{1}(i = j) 2y_i$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y'_i \mathbf{x}_i$$

$$\frac{\partial w_j}{\partial y'_i} = \alpha_i x_{ij}$$

# Gradient Computation

$$\nabla_{\mathbf{z}} U = \nabla_{\mathbf{w}} U \cdot \nabla_{y'} \mathbf{w} \cdot \nabla_{\mathbf{z}} y'$$

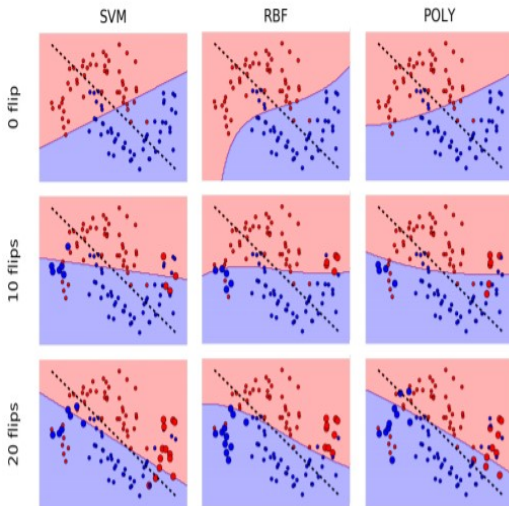
$$\frac{\partial U}{\partial w_j} = \frac{\|\mathbf{w}\|^2 w_j^* - \mathbf{w}^\top \mathbf{w}^* w_j}{\|\mathbf{w}\|^3 \|\mathbf{w}^*\|}$$

$$\frac{\partial y'_i}{\partial z_j} = -\mathbb{1}(i = j) 2y_i$$

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y'_i \mathbf{x}_i$$

$$\frac{\partial w_j}{\partial y'_i} = \alpha_i x_{ij}$$

# Experimental Results



*Thank You*