

TSFool: Crafting Highly-Imperceptible Adversarial Time Series through Multi-Objective Attack

Yanyun Wang¹, Dehui Du^{2*}, Haibo Hu^{1*}, Zi Liang¹ and Yuanhao Liu²

¹ Department of Electrical and Electronic Engineering, The Hong Kong Polytechnic University, Hong Kong, China

² Software Engineering Institute, East China Normal University, Shanghai, China

* Corresponding Authors

Presenter: Yanyun Wang

□ Topic Introduction

- Neural network (NN) classifiers are vulnerable to **adversarial samples**, which means imperceptible perturbations added to the input can cause the output to change significantly.



x
“panda”
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

=



$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy.
"Explaining and harnessing adversarial examples." In ICLR, 2015.

□ Existing Knowledge

- **Adversarial attack** is to artificially craft adversarial samples to measure the robustness of NN models.

$$\vec{x}^* = \vec{x} + \delta_{\vec{x}} = \vec{x} + \min \|\vec{z}\| \text{ s.t. } f(\vec{x} + \vec{z}) \neq f(\vec{x})$$

□ Existing Knowledge

- **Adversarial attack** is to artificially craft adversarial samples to measure the robustness of NN models.

$$\vec{x}^* = \vec{x} + \delta_{\vec{x}} = \vec{x} + \min \|\vec{z}\| \text{ s.t. } f(\vec{x} + \vec{z}) \neq f(\vec{x})$$

- Gradient-based white-box adversarial attacks have achieved impressive performance on **feed-forward** NN classifiers and **image** data.

- FGSM $\delta_{\vec{x}} = \varepsilon \text{sign}(\nabla_{\vec{x}} \mathcal{L}(f, \vec{x}, \vec{y}))$

- PGD $\vec{x}^{t+1} = \Pi_{\epsilon} \{ \vec{x}^t + \varepsilon \cdot \text{sign}(\nabla_{\vec{x}} \mathcal{L}(f, \vec{x}^t, \vec{y})), \vec{x} \}$

□ Current Gap

- While recent years have witnessed the success of **recurrent** neural network (RNN) models in **time series** classification (TSC) tasks, the gradient-based white-box adversarial attacks cannot perform well on RNN-based TSC.

□ Current Gap

- While recent years have witnessed the success of **recurrent** neural network (RNN) models in **time series** classification (TSC) tasks, the gradient-based white-box adversarial attacks cannot perform well on RNN-based TSC.

□ Reasons

- From RNN: the unique cyclical computation in RNN architecture prevents **direct model differentiation**, which means the majority of gradient information is no longer directly available through the **chain rule**.

□ Current Gap

- While recent years have witnessed the success of **recurrent** neural network (RNN) models in **time series** classification (TSC) tasks, the gradient-based white-box adversarial attacks cannot perform well on RNN-based TSC.

□ Reasons

- From RNN: the unique cyclical computation in RNN architecture prevents **direct model differentiation**, which means the majority of gradient information is no longer directly available through the **chain rule**;
- From TSC: time series data are far more **visually sensitive** to perturbations than image data, which poses challenges to the conventional **local optimization objective** of adversarial attack to minimize the perturbation amount for every single sample.

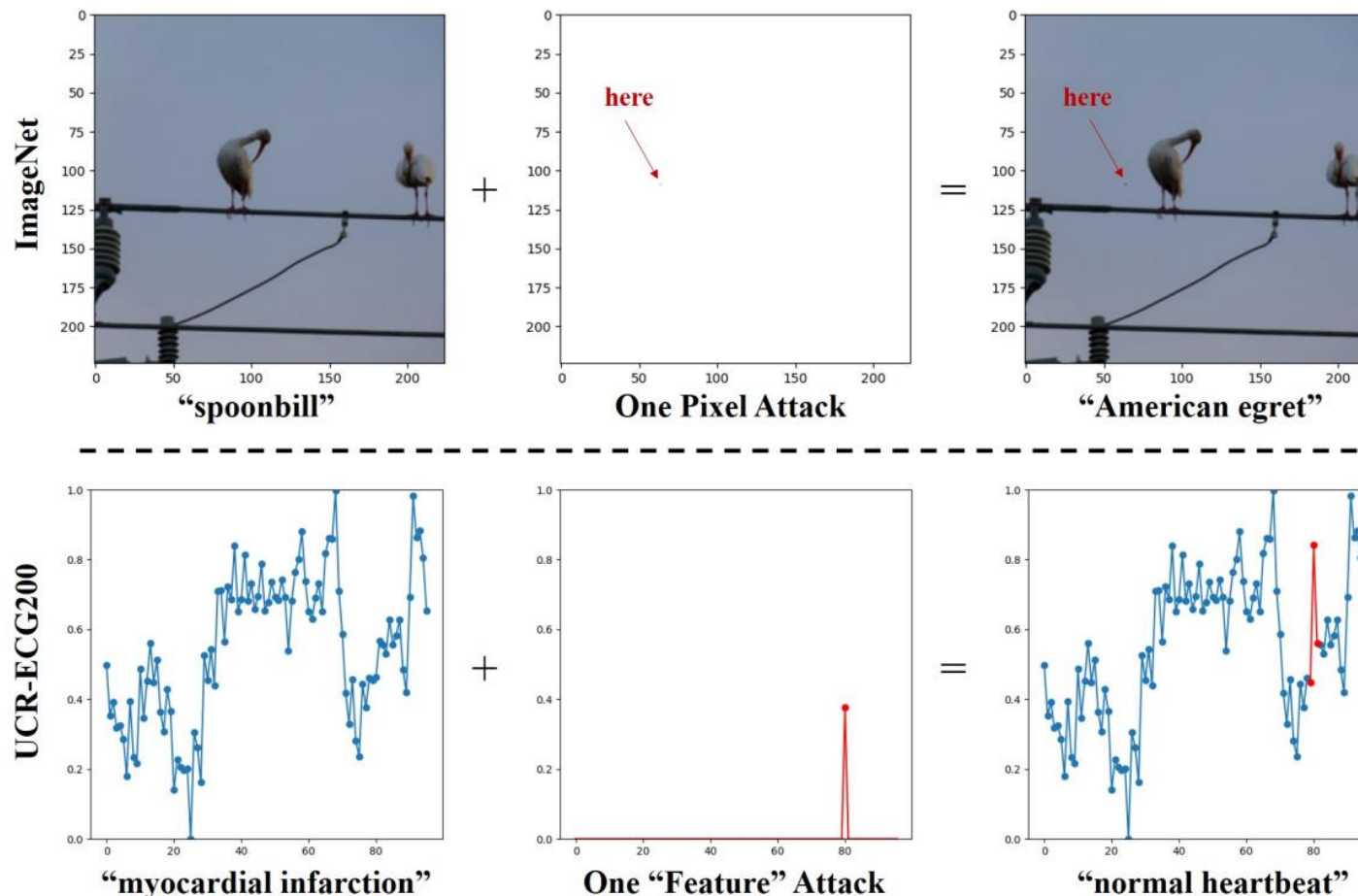
□ Current Gap

- While recent years have witnessed significant progress in **time series** classification (TSC), models cannot perform well on RNNs.

□ Reasons

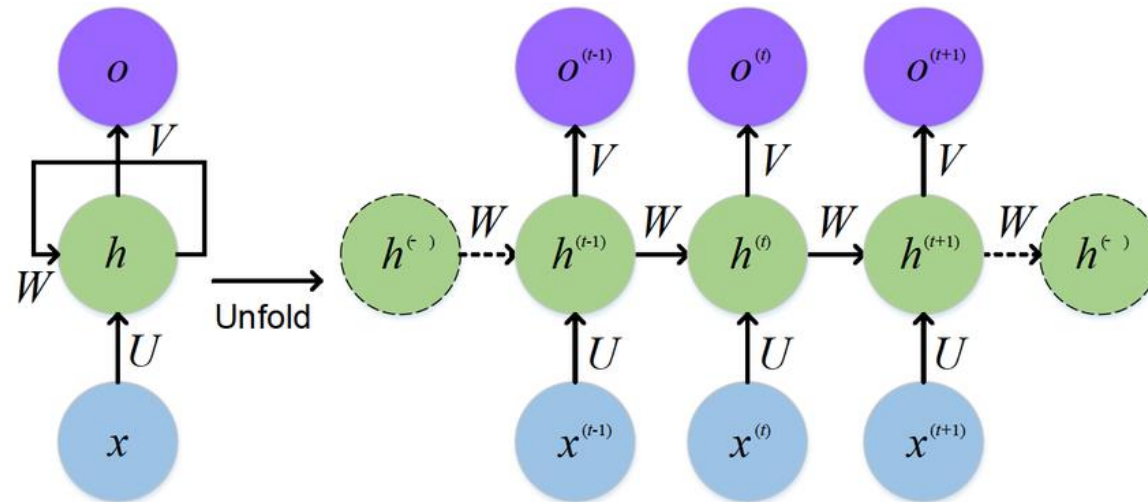
- From RNN: the unique cyclic **differentiation**, which means that the model is not available through the **chain rule**.
- From TSC: time series data are often noisy, which poses challenges for adversarial attack to minimize the perturbation.

Toy Perturbations with the Same Degree (ℓ_∞ norm: 37.72%) but Different Imperceptibility



□ Related Works

- Making RNN completely differentiable by **cyclical computational graph unfolding** to cater to the gradient-based methods.
 - Turns out to be **inefficient** and hard to stably **scale** in real-world practice.



□ Related Works

- Making RNN completely differentiable by **cyclical computational graph unfolding** to cater to the gradient-based methods.
 - Turns out to be **inefficient** and hard to stably **scale** in real-world practice.
- Just **ignoring** the given model knowledge and implementing **non-gradient-based** black-box attacks.

□ Related Works

- Making RNN completely differentiable by **cyclical computational graph unfolding** to cater to the gradient-based methods.
 - Turns out to be **inefficient** and hard to stably **scale** in real-world practice.
- Just **ignoring** the given model knowledge and implementing **non-gradient-based** black-box attacks.
 - By **model querying**:
 - A majority of these methods like One Pixel Attack and Square Attack can only work for **image** input.

□ Related Works

- Making RNN completely differentiable by **cyclical computational graph unfolding** to cater to the gradient-based methods.
 - Turns out to be **inefficient** and hard to stably **scale** in real-world practice.
- Just **ignoring** the given model knowledge and implementing **non-gradient-based** black-box attacks.
 - By **model querying**:
 - A majority of these methods like One Pixel Attack and Square Attack can only work for **image** input; and
 - The rest of these methods like Boundary Attack and HopSkipJump are extremely **time-consuming** because they rely on random-walking.

□ Related Works

- Making RNN completely differentiable by **cyclical computational graph unfolding** to cater to the gradient-based methods.
 - Turns out to be **inefficient** and hard to stably **scale** in real-world practice.
- Just **ignoring** the given model knowledge and implementing **non-gradient-based** black-box attacks.
 - By **model querying**:
 - A majority of these methods like One Pixel Attack and Square Attack can only work for **image** input; and
 - The rest of these methods like Boundary Attack and HopSkipJump are extremely **time-consuming** because they rely on random-walking.
 - By **adversarial transferability**:
 - Tends to achieve reasonable time and small perturbation, but **the worst** attack success rate.

□ Target

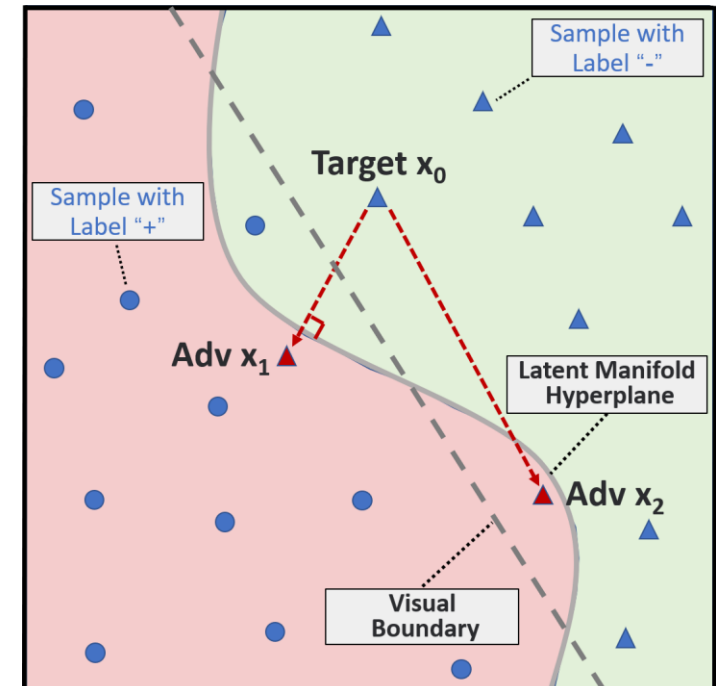
□ Target

- A non-gradient-based adversarial attack method (due to **RNN** model).

□ Target

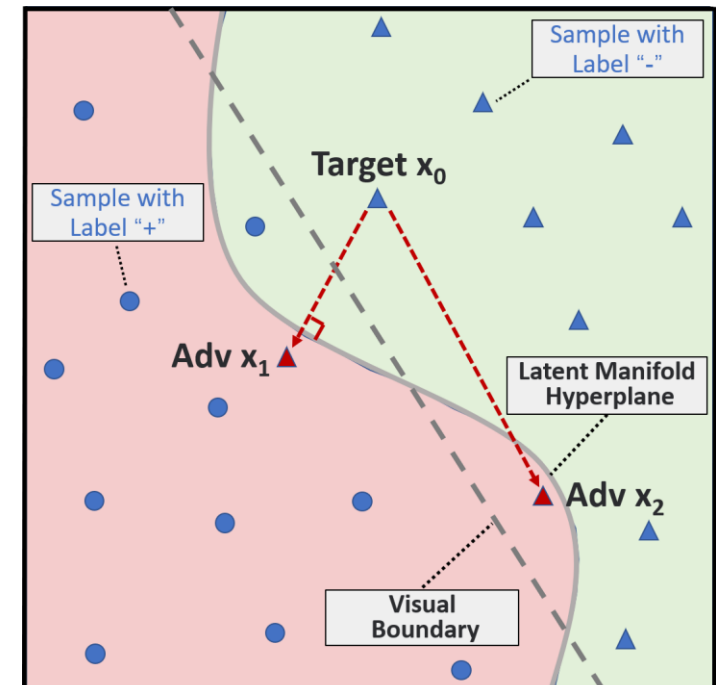
- A non-gradient-based adversarial attack method (due to **RNN** model),
- With additional consideration for the imperceptibility of perturbation (due to **time series** data).

□ Explanation for adversarial sample based on **Manifold Hypothesis**



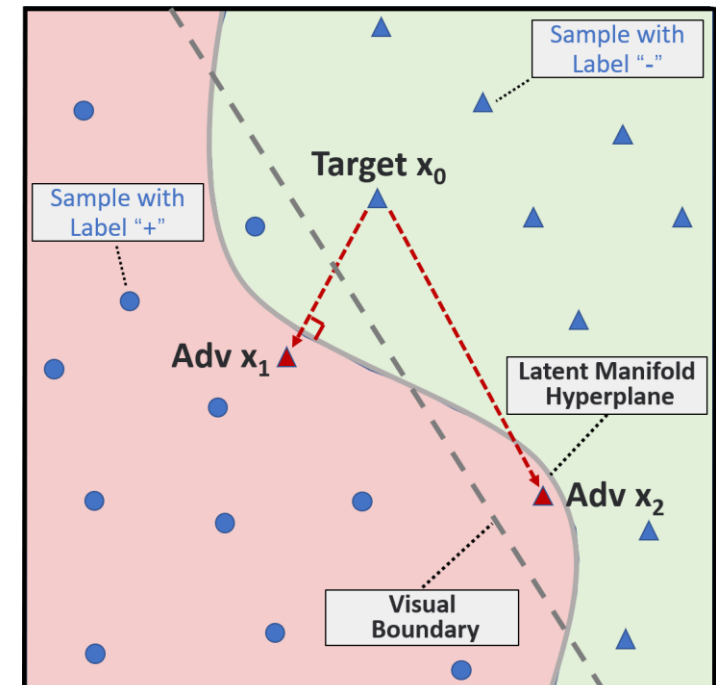
□ Explanation for adversarial sample based on **Manifold Hypothesis**

- NN-based classification relies on the latent manifold hyperplane (i.e., the classification boundary).



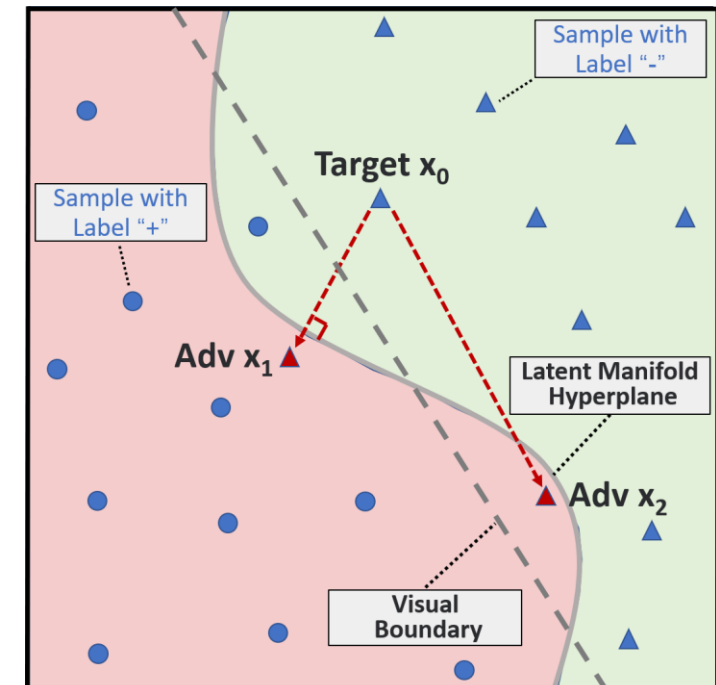
□ Explanation for adversarial sample based on **Manifold Hypothesis**

- NN-based classification relies on the latent manifold hyperplane (i.e., the classification boundary);
- The **features** of input data cannot always semantically and visually reflect the **latent manifold**.



□ Explanation for adversarial sample based on **Manifold Hypothesis**

- NN-based classification relies on the latent manifold hyperplane (i.e., the classification boundary);
- The **features** of input data cannot always semantically and visually reflect the **latent manifold**;
- A small perturbation in **human cognition** imposed on a sample may completely overturn the **perception of NN** to its latent manifold.

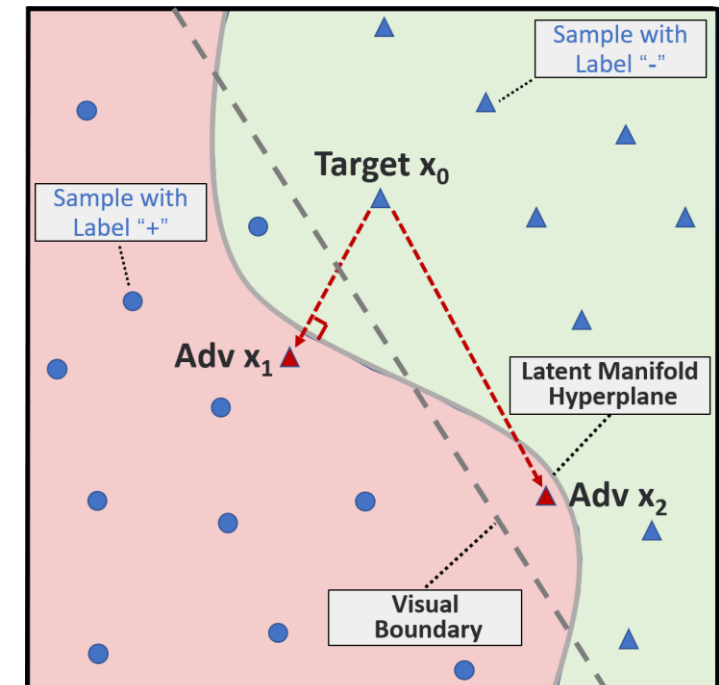


□ Explanation for adversarial sample based on **Manifold Hypothesis**

- NN-based classification relies on the latent manifold hyperplane (i.e., the classification boundary);
- The **features** of input data cannot always semantically and visually reflect the **latent manifold**;
- A small perturbation in **human cognition** imposed on a sample may completely overturn the **perception of NN** to its latent manifold.

□ Arguments

- Even the **minimal local perturbation** is not necessarily the most imperceptible one from the **global** perspective.

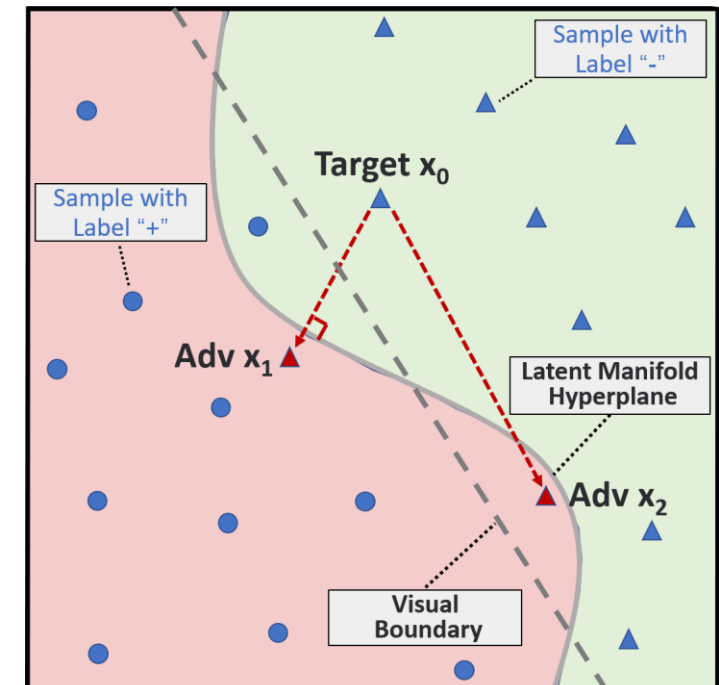


□ Explanation for adversarial sample based on **Manifold Hypothesis**

- NN-based classification relies on the latent manifold hyperplane (i.e., the classification boundary);
- The **features** of input data cannot always semantically and visually reflect the **latent manifold**;
- A small perturbation in **human cognition** imposed on a sample may completely overturn the **perception of NN** to its latent manifold.

□ Arguments

- Even the **minimal local perturbation** is not necessarily the most imperceptible one from the **global** perspective;
- The conventional approach to approximate the local optimization objective does **not** always lead to a highly-imperceptible adversarial attack.



□ Camouflage Coefficient

- A novel **global optimization objective** that takes the relative position between adversarial samples and class clusters into consideration, to measure the imperceptibility of adversarial samples from the perspective of **class distribution**.

□ Camouflage Coefficient

- A novel **global optimization objective** that takes the relative position between adversarial samples and class clusters into consideration, to measure the imperceptibility of adversarial samples from the perspective of **class distribution**.

$$\mathcal{C}(\vec{x}^*) = \frac{\|\vec{x}^* - \vec{m}_i\|/d_i}{\|\vec{x}^* - \vec{m}_j\|/d_j}$$

$$\vec{m}_i = \frac{1}{|\mathcal{X}_i|} \sum_{\vec{x}' \in \mathcal{X}_i} \vec{x}'$$

$$d_i = \frac{1}{|\mathcal{X}_i|} \sum_{\vec{x}' \in \mathcal{X}_i} \|\vec{x}' - \vec{m}_i\|$$

- It is the relative proportion of the norm distance between the adversarial sample and the **original class** to the distance between it and the **misclassified class** regarding the different cluster ranges of the classes.

□ Multi-objective Optimization Problem

- By adding the Camouflage Coefficient, we refine the adversarial attack task to a **multi-objective optimization** problem.

□ Multi-objective Optimization Problem

- By adding the Camouflage Coefficient, we refine the adversarial attack task to a **multi-objective optimization** problem.

□ Approximative Solution

□ Multi-objective Optimization Problem

- By adding the Camouflage Coefficient, we refine the adversarial attack task to a **multi-objective optimization** problem.

□ Approximative Solution

- To cross the manifold hyperplane in the position that is:
 - **The closest** to the benign sample.

□ Multi-objective Optimization Problem

- By adding the Camouflage Coefficient, we refine the adversarial attack task to a **multi-objective optimization** problem.

□ Approximative Solution

- To cross the manifold hyperplane in the position that is:
 - **Sufficiently close** to the benign sample.

□ Multi-objective Optimization Problem

- By adding the Camouflage Coefficient, we refine the adversarial attack task to a **multi-objective optimization** problem.

□ Approximative Solution

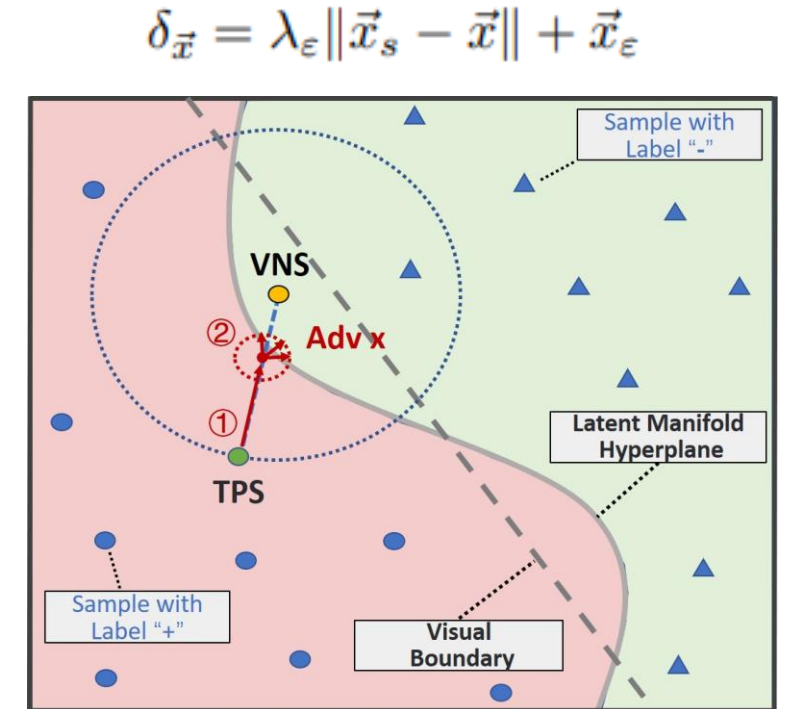
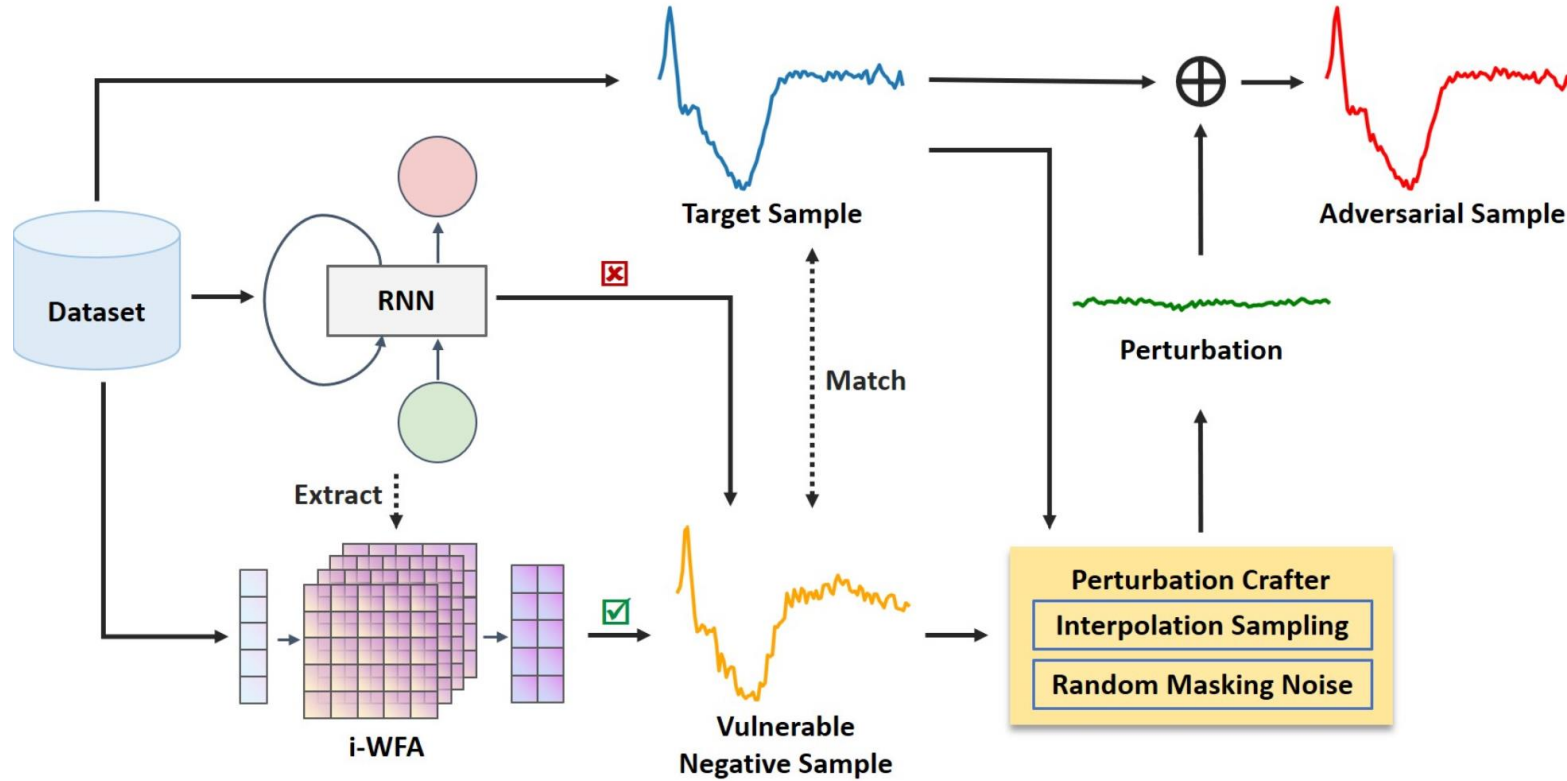
- To cross the manifold hyperplane in the position that is:
 - **Sufficiently close** to the benign sample; and also
 - **Sufficiently close** to the center of mass of the benign class (i.e., m_i).

□ Multi-objective Optimization Problem

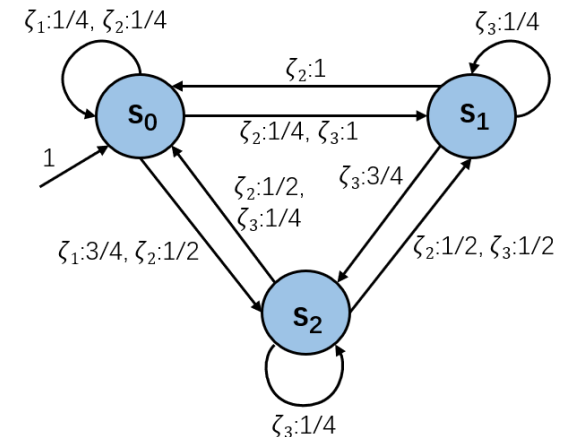
- By adding the Camouflage Coefficient, we refine the adversarial attack task to a **multi-objective optimization** problem.

□ Approximative Solution

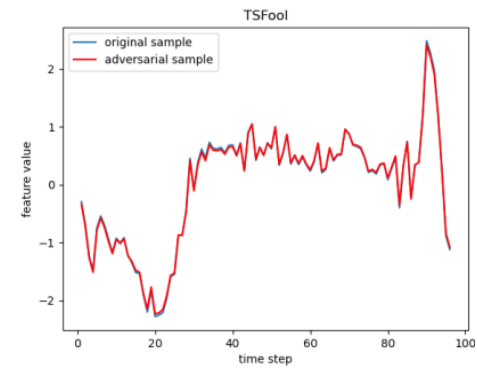
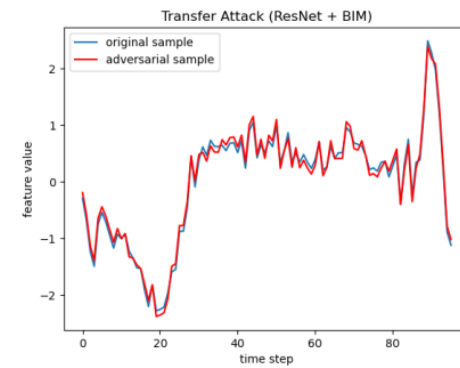
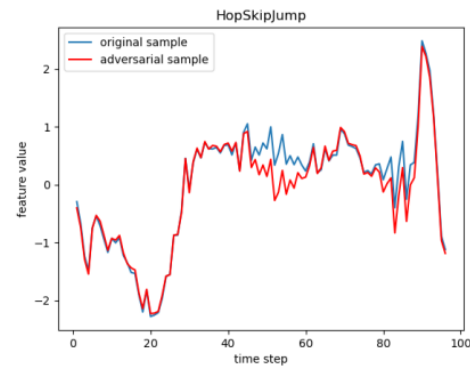
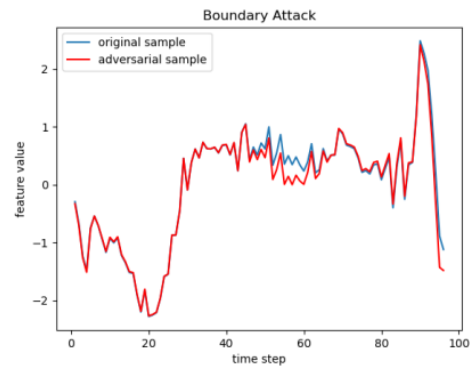
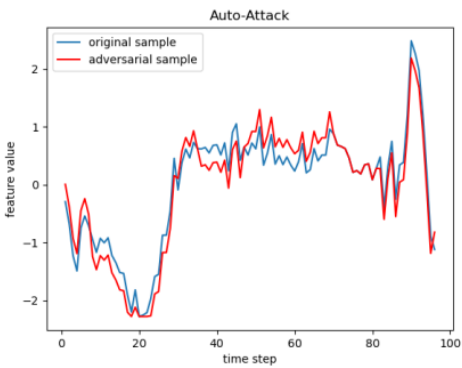
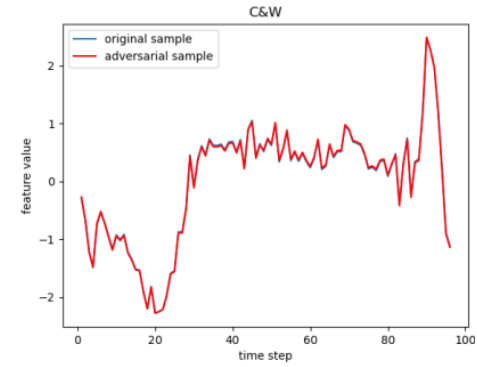
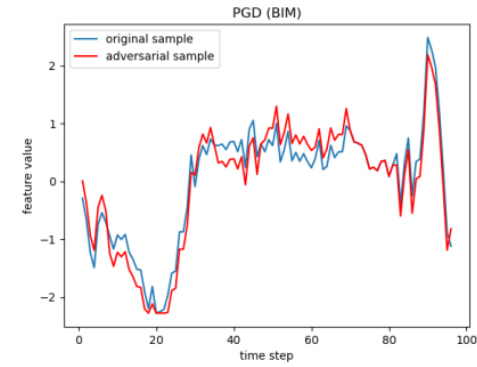
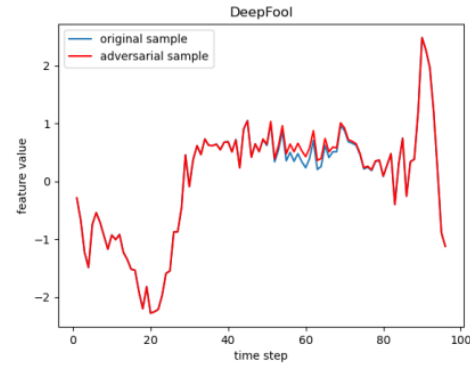
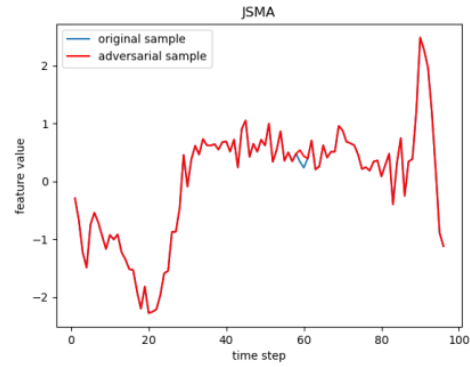
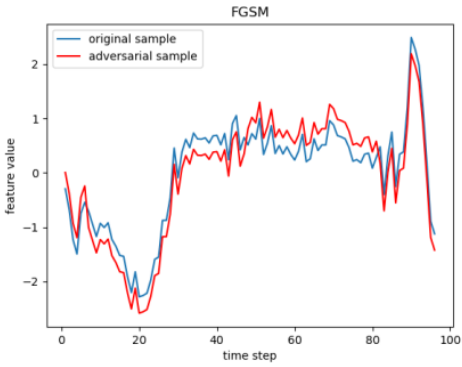
- To cross the manifold hyperplane in the position that is:
 - **Sufficiently close** to the benign sample; and also
 - **Sufficiently close** to the center of mass of the benign class (i.e., m_i),
- In a **non-gradient-based** way.

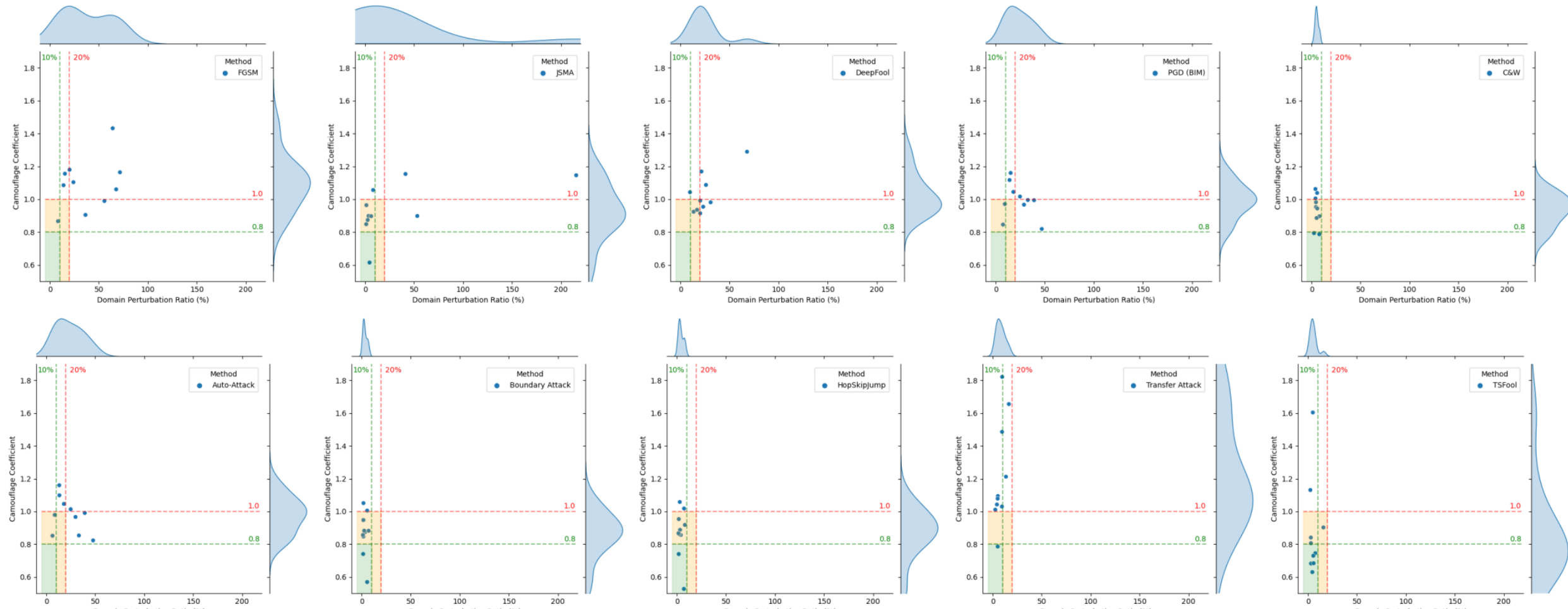


The **representation model** named i-WFA is built only upon the RNN's outputs. It can fit the manifold hyperplane of an RNN classifier but distinguish samples by their original features like humans. As a result, it can capture deeply embedded vulnerable samples whose **features** deviate from the **latent manifold**.

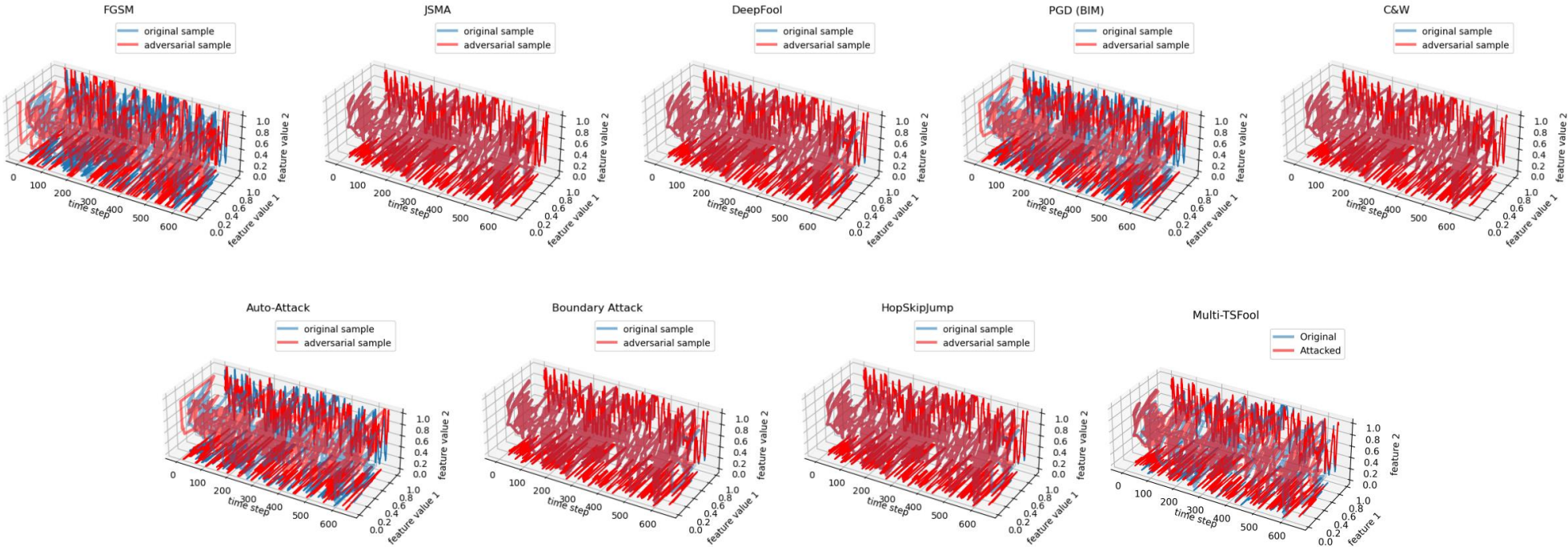


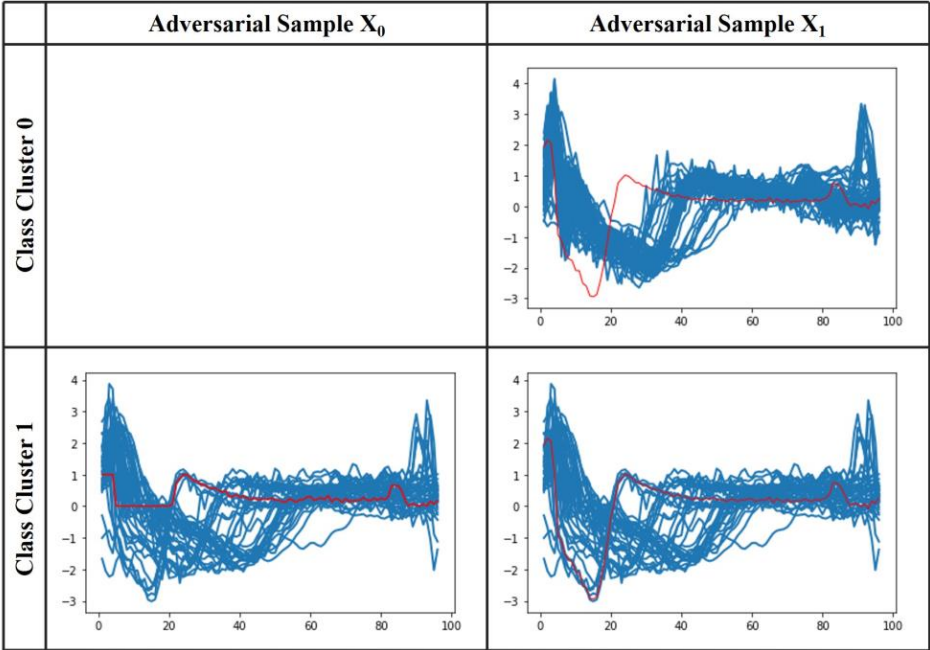
Method	Attack Success Rate	Generation Number	Average Time Cost (s)	Perturbation Ratio (ρ^*)	Camouflage Coefficient
FGSM	72.12%	300.75	<u>0.0018</u>	37.13%	1.0804
JSMA	83.53%		1.0287	15.06%	0.9476
DeepFool	81.58%		0.0276	21.45%	1.0107
PGD (BIM)	76.84%		0.1327	22.71%	0.9938
C&W	69.90%		3.2016	5.16%	0.9372
Auto-Attack	80.11%		0.1824	22.55%	0.9745
Boundary Attack	79.01%		9.0399	<u>3.04%</u>	0.8788
HopSkipJump	83.17%	250	12.3068	<u>3.86%</u>	0.8872
Transfer Attack	19.54%		-	7.68%	1.2010
TSFool	<u>87.76%</u>	305	0.0230	4.63%	<u>0.8147</u>





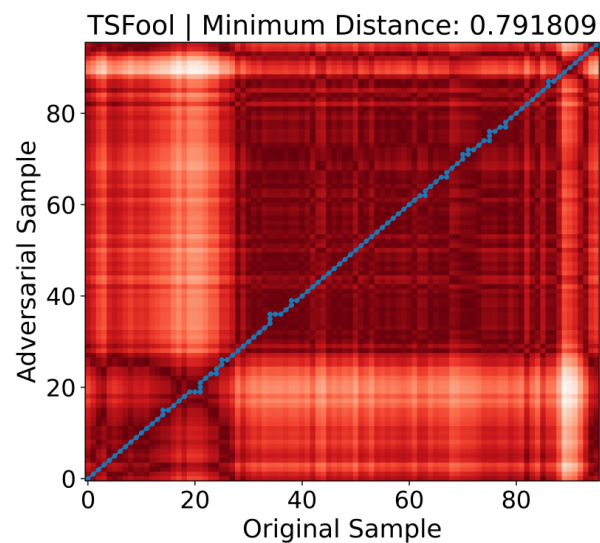
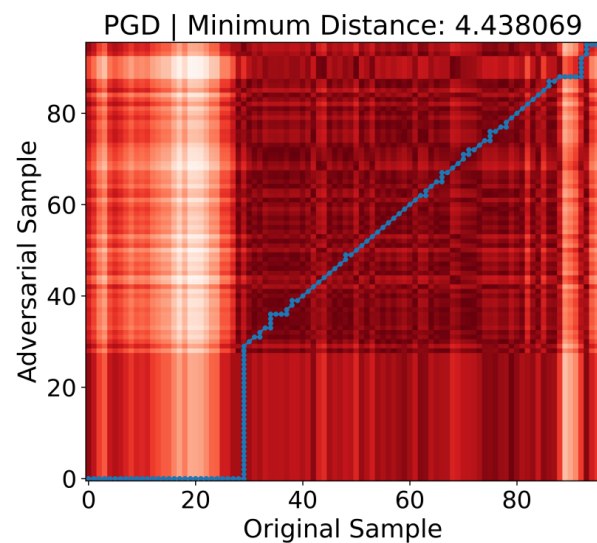
Target Model		Method	Attack Success Rate	Generation Number	Average Time Cost (s)	Perturbation Ratio (ρ^*)			Camouflage Coefficient
Dataset	Accuracy					ℓ_1	ℓ_2	ℓ_∞	
AF	0.8000	FGSM	80.00%	15	0.0127	24.52%	25.87%	29.04%	1.0211
		JSMA	86.67%		7.4301	3.04%	3.74%	4.99%	0.8126
		DeepFool	80.00%		0.9900	0.54%	0.64%	0.85%	0.7551
		PGD (BIM)	80.00%		0.9406	16.79%	18.22%	21.72%	0.9918
		C&W	40.00%		13.2858	0.48%	0.52%	0.63%	0.7958
		Auto-Attack	80.00%		2.1487	15.80%	17.11%	20.35%	0.9918
		Boundary Attack	66.67%		418.2511	0.43%	0.51%	0.67%	0.8045
		HopSkipJump	86.67%	20	78.3258	0.84%	0.97%	1.27%	0.8066
		TSFool	100.00%		0.0960	5.89%	6.69%	8.65%	0.6047





Human Study 2

Study	Option	Question							Sum	Count
		1	2	3	4	5	6	7		
Study 1	The Original Class Cluster	20	13	55	37	47	55	31	258	5
	The Misclassified Class Cluster	41	48	4	22	13	8	28	164	2
	Neutral	4	4	6	6	5	2	6	33	0
Study 2	The Adversarial Sample from TSFool	54	58	51	57	58	60	58	396	7
	The Adversarial Sample from PGD	10	5	12	6	4	3	5	45	0
	Neutral	1	2	2	2	3	2	2	14	0



Method	Metric	Time Series Anomaly Detection			
		OCSVM	IF	LOF	LSTMOD
PGD (BIM)	<i>Pre</i>	0.1755	0.1899	0.2794	0.2107
	<i>Re</i>	0.4890	0.5377	0.7018	0.7091
	<i>F1</i>	0.2454	0.2692	0.3782	0.3092
C&W	<i>Pre</i>	0.0637	0.0534	0.0463	0.0968
	<i>Re</i>	0.1304	0.1201	0.0745	0.3377
	<i>F1</i>	0.0798	0.0696	0.0534	0.1432
HopSkipJump	<i>Pre</i>	0.0693	0.0473	0.0801	0.1381
	<i>Re</i>	0.1561	0.1115	0.2282	0.5127
	<i>F1</i>	0.0897	0.0640	0.1146	0.2129
TSFool	<i>Pre</i>	0.0505	0.0346	0.0460	0.0741
	<i>Re</i>	0.1012	0.0829	0.1274	0.3218
	<i>F1</i>	0.0622	0.0469	0.0639	0.1175

□ **Insights**

□ **Contributions**

□ Insights

- **General** consideration beyond **image** data and **feed-forward** models is still lacking in the current knowledge of adversarial attack.

□ Contributions

□ Insights

- **General** consideration beyond **image** data and **feed-forward** models is still lacking in the current knowledge of adversarial attack;
- **Imperceptibility measures** of adversarial samples have not received sufficient attention, without which it would be hard to fairly define “adversarial”.

□ Contributions

□ Insights

- **General** consideration beyond **image** data and **feed-forward** models is still lacking in the current knowledge of adversarial attack;
- **Imperceptibility measures** of adversarial samples have not received sufficient attention, without which it would be hard to fairly define “adversarial”.

□ Contributions

- A novel global **optimization objective** "Camouflage Coefficient" to refine the adversarial attack as a multi-objective optimization problem.

□ Insights

- **General** consideration beyond **image** data and **feed-forward** models is still lacking in the current knowledge of adversarial attack;
- **Imperceptibility measures** of adversarial samples have not received sufficient attention, without which it would be hard to fairly define “adversarial”.

□ Contributions

- A novel global **optimization objective** "Camouflage Coefficient" to refine the adversarial attack as a multi-objective optimization problem;
- A new latent manifold-based **methodology** to heuristically approximate the solution of the suggested optimization problem, which opens a **new feasible path** to craft imperceptible adversarial samples.

Thank You!

Acknowledgements

This work was supported by National Natural Science Foundation of China (Grant No: 92270123, 62072390), and the Research Grants Council, Hong Kong SAR, China (Grant No: PolyU 15203120, 15226221, 15209922, and 15210023).

Contact Information

Haibo Hu: haibo.hu@polyu.edu.hk & Dehui Du: dhdu@sei.ecnu.edu.cn