

# Time series adversarial attack

Xue Zhechang

<sup>1</sup> Karlsruhe Institute for Technology, Germany

<sup>2</sup> ugupg@student.kit.edu

**Abstract.** Nowadays, Time series data plays an important role in many fields. With the help of deep learning based model, people can classify and predict time series data, which is useful and efficient in data mining and recovering missing data. However, deep learning models are vulnerable to adversarial attacks. The researchers has transferred some adversarial attacks from image recognition domain to time series domain and prove their effectiveness. Thus, it's a big topic to detect these attacks and defense them to increasing the accuracy of classification and prediction. In this work, we will show the vulnerability of time series models, the method to attack models, to detect adversarial attacks and to prevent models from adversarial attacks.

**Keywords:** Time series · Adversarial attack · Deep Learning.

## 1 Introduction

Time series is a series of data points indexed in time order. Time series data are widely used in various fields ranging from mathematical statistics, signal processing and pattern recognition to quantitative finance and weather prediction.[1] People notice the big value of analysing time series data. Deep learning models have nowadays succeeded in many real life application such as speech recognition and computer visions. Thus, they are also applied to time series data. There's two main type of time series deep learning model. One is time series classification (TSC) model, which is mainly used to classify categories of time series data. Another is time series regression (TSR) model, which can predict the future data in a proper time point.

However, deep learning models are not robust. For example, researchers modify the original time series data, where the changes are too tiny to be detectable by human. Due to these changes, the accuracy of classification or prediction declines significantly. This is so-called adversarial attack.

In this paper, we will summarize some researches in time series domain. We will begin with adversarial attacks. On this basis, we will prove the vulnerability of time series model. Then we will show how to detect adversarial attacks from original data. Finally, we will introduce some methods to defense adversarial attacks.

## 2 Related Works

To the best of our knowledge, adversarial attacks were first studied by Fazle et al[2]. Their target models are 1-Nearest Neighbor Dynamic Time Warping (1-NN) DTW, a Fully Connected Network and a Fully Convolutional Network (FCN). They trained Adversarial transformation network (ATN) to attack target models and tested with University of California Riverside (UCR) datasets. Finally they proved that TSC models are vulnerable to adversarial attacks. They also proved the vulnerability of multivariate time series[4].

Gautam et al[7]. proved that TSR models are also vulnerable to adversarial attacks. They transferred existed attacks from computer vision domain to time series domain, which are called fast gradient sign method (FGSM) and basic iterative method (BIM). They also proved the transferability of adversarial attacks to different target models.

Hassan et al[3]. introduced how FGSM and BIM works in time series domain. They used multi-dimensional scaling as the measurement to evaluate the relation between perturbation and accuracy.

Pradeep et al[8]. defined 4 kinds of adversarial attacks: untargeted, targeted, individual universal attacks. Then they introduced the method of targeted attack and universal attack and proved their effectiveness on TSC models based on ResNet. Finally, they introduced that backpropagation algorithm can help increase the robustness of time series model.

Aidong et al[9]. introduced a new method to attack TSC models with higher efficiency. They measured the importance of adversarial samples and selected some of the most important adversarial samples to modify the original data. This method can decline the perturbation of original data but increase the effectiveness. Tiny perturbation of time series data can lead to big difference in classification and prediction. However, few researches are done about detection of adversarial attacks. Mubarak et al[?]. introduced a method to detect whether the time series data is adversarial generated by FGSM and BIM.

Shoaib et al[6]. transferred three defensive methods from computer vision domain to time series domain: Adversarial training, TRADES and feature denoising. To prove their effectiveness, they used FGSM and Projected Gradient Descent (PGD) as white-box attacks and noise attack, boundary attack and Simple Black-box Attack (SIMBA) as black-box attacks.

Zhongguo et al[10] introduced a new defend to adversarial attacks and proved its effectiveness by experiment. They used thermometer encoding to non-linear encode original time series data and trained a encode-decode model to decode the modified time series data. Then they trained the deep learning model on time series with these output, which can nearly completely ignore the affect of perturbation.

## 3 Adversarial attack methods

Adversarial attacks are normally divided into two categories: White-box attacks and black-box attacks.

In white-box attacks, the attacker has access to all the information about the targeted model. Thus, attacking with gradient is a common method in deep learning. Here we will introduce two attacks based on gradient: Fast gradient sign method (FGSM) and basic iterative method (BIM).

In black-box attacks, the attacker has no information or parameter about the targeted model. Thus, it's impossible to attack the targeted model with gradient. Attackers normally use the relation between input and output of the target to find the way to attack the model.

### 3.1 White-box attacks

*Definition 1* Time series data can be mathematically represented as set  $X = [x_1, x_2, x_3, \dots, x_T]$ , where  $T$  is the length of this set.

---

**Algorithm 1** Fast gradient sign method

---

Fast gradient sign method

Basic iterative method

### 3.2 Black-box attacks

Noise attack

Boundary attack

## 4 vulnerability

2009, 1903

### 4.1 Multi-Dimensional Scaling

### 4.2 Transferability

## 5 Detecting

2004

**5.1 Sample Entropy****5.2 Detrended Fluctuation Analysis****5.3 Class activation map**

2110

**5.4 The importance of adversarial sample****6 Defense**

2008, 2101, 2110

**6.1 Adversarial training****6.2 TRADES****6.3 Feature Denoising****6.4 Backpropagation****6.5 Non-linear transfer**

thermometer encoding

encode-decode model

**Method****7 Conclusion**

*Sample Heading (Fourth Level)* The contribution should contain no more than four levels of headings. Table 1 gives a summary of all heading levels.

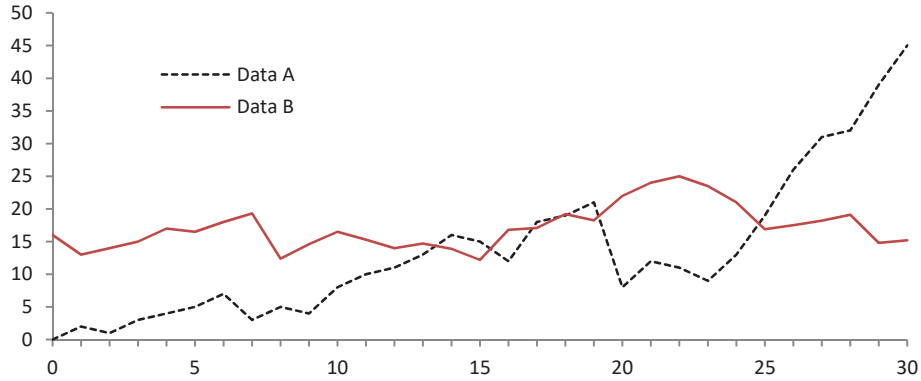
**Table 1.** Table captions should be placed above the tables.

Heading level	Example	Font size and style
Title (centered)	<b>Lecture Notes</b>	14 point, bold
1st-level heading	<b>1 Introduction</b>	12 point, bold
2nd-level heading	<b>2.1 Printing Area</b>	10 point, bold
3rd-level heading	<b>Run-in Heading in Bold.</b> Text follows	10 point, bold
4th-level heading	<i>Lowest Level Heading.</i> Text follows	10 point, italic

Displayed equations are centered and set on a separate line.

$$x + y = z \tag{1}$$

Please try to avoid rasterized images for line-art diagrams and schemas. Whenever possible, use vector graphics instead (see Fig. 1).



**Fig. 1.** A figure caption is always placed below the illustration. Please note that short captions are centered, while long ones are justified by the macro package automatically.

**Theorem 1.** *This is a sample theorem. The run-in heading is set in bold, while the following text appears in italics. Definitions, lemmas, propositions, and corollaries are styled the same way.*

*Proof.* Proofs, examples, and remarks have the initial word in italics, while the following text appears in normal font.

For citations of references, we prefer the use of square brackets and consecutive numbers. Citations using labels or the author/year convention are also acceptable. The following bibliography provides a sample reference list with entries for journal articles [?], an LNCS chapter [?], a book [?], proceedings without editors [?], and a homepage [?]. Multiple citations are grouped [?, ?, ?], [?, ?, ?, ?].

## References

1. Wikipedia contributors. Time series, 2022. URL: [https://en.wikipedia.org/wiki/Time\\_series](https://en.wikipedia.org/wiki/Time_series)
2. Fazle Karim and Houshang Darabi, Adversarial Attacks on Time Series, 2019 URL: <https://ieeexplore.ieee.org/abstract/document/9063523>
3. Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar and Pierre-Alain Muller IRIMAS, Universite Haute-Alsace, Mulhouse, France, Adversarial Attacks on Deep Neural Networks for Time Series Classification, URL: <https://ieeexplore.ieee.org/abstract/document/8851936>
4. Samuel Harford, Fazle Karim, and Houshang Darabi, Adversarial Attacks on Multivariate Time Series, URL: <https://arxiv.org/abs/2004.00410>
5. Mubarak G. Abdu-Aguye, Walid Gomaa, Yasushi Makihara, Yasushi Yagi, Cyber Physical Systems Lab, Egypt Japan University of Science and Technology, Egypt, Faculty of Engineering, Alexandria University, Egypt, The Institute of Scientific and Industrial Research, Osaka University, Japan, Detecting adversarial attacks in time-series data, URL: <https://ieeexplore.ieee.org/abstract/document/9053311>

6. Shoaib Ahmed Siddiqui, Andreas Dengel, and Sheraz Ahmed, German Research Center for Artificial Intelligence (DFKI), 67663 Kaiserslautern, Germany, Benchmarking adversarial attacks and defenses for time-series data, URL: [https://link.springer.com/chapter/10.1007/978-3-030-63836-8\\_45](https://link.springer.com/chapter/10.1007/978-3-030-63836-8_45)
7. Gautam Raj Mode and Khaza Anuarul Hoque, Department of Electrical Engineering Computer Science, University of Missouri, Columbia, MO, USA, Adversarial Examples in Deep Learning for Multivariate Time Series Regression, URL: <https://ieeexplore.ieee.org/abstract/document/9425190>
8. Pradeep Rathore, Arghya Basak, Sri Harsha Nistala, Venkataramana Runkana, TCS Research, Pune, India, Untargeted, Targeted and Universal Adversarial Attacks and Defenses on Time Series, URL: <https://ieeexplore.ieee.org/abstract/document/9207272>
9. Aidong Xu, Xuechun Wang, Yunan Zhang, Tao Wu, Xingping Xian. Adversarial Attacks on Deep Neural Networks for Time Series Prediction, URL: <https://dl.acm.org/doi/10.1145/3485314.3485316>
10. Zhongguo Yang, Irshad Ahmed Abbasi, Fahad Algarni, Sikan-dar Ali, and Mingzhu Zhang, An IoT Time Series Data Security Model for Adversarial Attack Based on Thermometer Encoding URL: <https://www.hindawi.com/journals/scn/2021/5537041/>