

Winter Term 2022

Application for privacy and security in Deep learning

Patrick NONKI¹

Contents

1	Motivation	4
2	The types of attacks on security	4
2.1	Model Extraction Attack	4
2.2	Model Inversion Attack	4
2.3	Adversarial Attack	4
2.4	Poisoning Attack	4
3	DL Algorithms of defence in security	4
4	The types of attacks on privacy	4
4.1	Model Extraction Attack	4
4.2	Model Inversion Attack	4
4.3	Adversarial Attack	4
4.4	Poisoning Attack	4
5	DL Algorithms of defence in privacy	4
6	Conclusion	4

¹ patrick.nonki@hshl.de

Abstract: Nowadays, a large amount of data is being generated for usage in big companies like Google, Siemens, Meta and it arose the challenge to deal with them in a context of making them as much as secure and private possible. Since the challenge of keeping them confidential with applications in a very wide range of applications like face recognition, farming, or path recognition, deep learning which is a part of the big family of Machine learning, could be a potential solution since it enables to deal with a lot of data at the same time and in a secure way since most of these companies deal with sensitive data of customers (for example private images, geographic positions, locations and passwords).

DL's privacy problems have recently come to light, and a number of attacks have been suggested. Model extraction attacks and model inversion attacks are two types of attacks that violate a model's privacy. In model extraction attacks, the adversary seeks to replicate the parameters/hyperparameters of the model that is used to deliver cloud-based ML services, jeopardizing the confidentiality of the ML algorithms and the service provider's intellectual property. Attackers use information already accessible to infer sensitive information in model inversion attacks.[1]

Several strategies have been put out to address privacy and security concerns in DL. In terms of DL privacy for privacy-preserving, there are now four primary technologies: differential privacy, homomorphic encryption, secure multi-party computation, and trusted execution environment. By using differential privacy, the adversary is prevented from deducing whether a certain instance was used to train the target model. The privacy of the training and testing data is prioritized by the homomorphic encryption and secure multi-party computation scheme. In order to safeguard training code and sensitive data, the trusted execution environment tries to employ hardware to create a secure and isolated environment.[1]

1 Motivation

2 The types of attacks on security

2.1 Model Extraction Attack

2.2 Model Inversion Attack

2.3 Adversarial Attack

2.4 Poisoning Attack

3 DL Algorithms of defence in security

4 The types of attacks on privacy

4.1 Model Extraction Attack

4.2 Model Inversion Attack

4.3 Adversarial Attack

4.4 Poisoning Attack

5 DL Algorithms of defence in privacy

6 Conclusion

Bibliography

- [1] Ximeng Liu; Lehui Xie; Yaopeng Wang; Jian Zou; Jinbo Xiong; Zuobin Ying; Athanasios V. Vasilakos, Privacy and Security Issues in Deep Learning: A Survey, 15th December 2020
- [2] Muhammad Imran Tariq ,Nisar Ahmed Memon, Shakeel Ahmed, Shahza di Tayyaba, Muhammad Tahir Mushtaq, Natash Ali Mian,5Muhammad Imran, and Muhammad W. Ashraf; A Review of Deep Learning Security and Privacy Defensive Techniques; 07 Apr 2020.
- [3] Georgios A. Kaissis, Marcus R. Makowski, Daniel Rückert Rickmer F. Braren, Secure, privacy-preserving and federated machine learning in medical imaging; 08 june 2020.