

Hush-hush Gradients: A Review of Differential Privacy for Deep Learning

Søren Winkel Holm

June 10, 2022

Introduction

The field of Deep Learning (DL) is for many subfields moving towards a setup where large multi-purpose foundation models are developed and trained at major companies or research institutions, and then released for engineers to adapt to specific applications [Bom+21, pp. 3]. This application of the open-source principle to pre-trained models improves scientific reproduction ability [HO20, pp. 3] and technology accessibility [Bom+21, pp. 139]. One risk, however, is an adversarial actor exploiting a property of DL models: Parts of training data is generally recoverable from model weights [NSH19; Sho+17]. This might expose proprietary data or the private data of individuals as exemplified for Natural Language Processing (nlp) language models in Figure 1. As large-scale data sets are here to stay [Sun+17], algorithmic methods for improving the privacy of foundation models are required. The methods of Differential Privacy (DP) are suitable for this task and the relevant con-

cepts, algorithms and problems will here be reviewed.

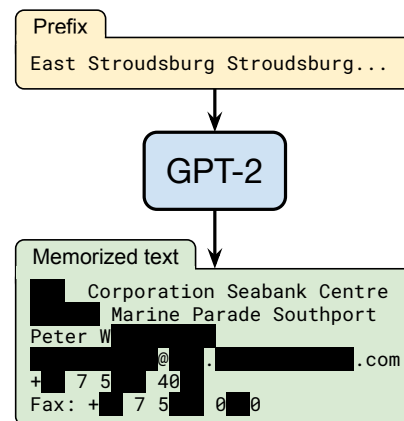


Figure 1: The extraction attack performed on GPT-2 [Car+21, Fig. 1] (private data redacted).

Fundamental Concepts

Achieving DP corresponds to making a promise of hiding information about individuals when publishing quantitative patterns about groups [DR14, pp. 5]. This

general problem is historically faced in releases of statistical data analyses by e.g. official organizations [Dal77; Wik22]. An algorithm is thus differentially private if a third party observer cannot extract individual information from its' computation. In this context, a Machine Learning (ML) model $f(x|w) = \hat{y} \approx y$ trained on a data set \mathcal{D} exposes data patterns when either its' parametrization w or predictions (x, \hat{y}) are released.

State of the Art

Open Problems

References

- [Bom+21] Rishi Bommasani et al. “On the Opportunities and Risks of Foundation Models”. In: *ArXiv* abs/2108.07258 (2021).
- [Car+21] Nicholas Carlini et al. “Extracting Training Data from Large Language Models”. In: *30th USENIX Security Symposium (USENIX Security 21)*. USENIX Association, Aug. 2021, pp. 2633–2650. ISBN: 978-1-939133-24-3. URL: <https://www.usenix.org/conference/usenixsecurity21/presentation/carlini-extracting>.
- [Dal77] Tore Dalenius. “Towards a methodology for statistical disclosure control”. In: *Statistisk Tidskrift* 15 (1977), pp. 429–444.
- [DR14] Cynthia Dwork and Aaron Roth. “The Algorithmic Foundations of Differential Privacy”. In: *Found. Trends Theor. Comput. Sci.* 9.3–4 (Aug. 2014), pp. 211–407. ISSN: 1551-305X. DOI: 10.1561/04000000042. URL: <https://doi.org/10.1561/04000000042>.
- [HO20] Matthew Hartley and Tjelvar S.G. Olsson. “dtoolAI: Reproducibility for Deep Learning”. In: *Patterns* 1.5 (2020), p. 100073. ISSN: 2666-3899. DOI: <https://doi.org/10.1016/j.patter.2020.100073>. URL: <https://www.sciencedirect.com/science/article/pii/S2666389920300933>.
- [NSH19] Milad Nasr, Reza Shokri, and Amir Houmansadr. “Comprehensive Privacy Analysis of Deep Learning: Passive and Active White-box Inference Attacks against Centralized and Federated Learning”. In: Mar. 2019. DOI: 10.1109/SP.2019.00065.
- [Sho+17] R. Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. “Membership Inference Attacks Against Machine Learning Models”. In: *2017 IEEE Symposium on Security and Privacy (SP)* (2017), pp. 3–18.
- [Sun+17] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. “Revisiting Unreasonable Effectiveness of Data in Deep Learning Era”. In: *2017 IEEE International Conference on Computer Vision (ICCV)*. 2017, pp. 843–852. DOI: 10.1109/ICCV.2017.97.
- [Wik22] Wikipedia contributors. *Differential privacy* — *Wikipedia, The Free Encyclopedia*. https://en.wikipedia.org/w/index.php?title=Differential_privacy&oldid=1091066967. [Online; accessed 2-June-2022]. 2022.