

Draft research question

April 22, 2020

1 Roadmap

1. introduction and literature study (1.5 weeks) introduction + related works
2. **draft methodology to answer research question (1 week)**
3. experiment with simple attacks (1-2 weeks)
4. implement methodology and start experiment (3-4 weeks)
5. analyse results (2-3 weeks)
6. write paper (2 weeks)
7. presentation(1-2 weeks)

2 Research question

Machine learning applications have grown in popularity over the last decade. One of the most used kinds of models are neural networks which grew in popularity with the advances in hardware and deep learning algorithms, and because they excel in image recognition and natural language processing applications. Cloud providers have included Machine Learning as a Service to their offer, it is now easier than ever for businesses to deploy machine learning models. MLaaS platforms such as [AWS](https://aws.amazon.com/machine-learning/)¹, [Microsoft Azure](https://azure.microsoft.com/en-us/services/machine-learning/)², [Google Cloud AI](https://cloud.google.com/ai-platform)³ or [IBM Watson](https://www.ibm.com/cloud/machine-learning)⁴ offer their customers to train their model given

¹<https://aws.amazon.com/machine-learning/>

²<https://azure.microsoft.com/en-us/services/machine-learning/>

³<https://cloud.google.com/ai-platform>

⁴<https://www.ibm.com/cloud/machine-learning>

they provide the dataset. ~~The final model can be queried as a black-box or not, it is nevertheless sensible to privacy attacks aiming at inferring information from the data it was trained on.~~ I will add a description of the different kinds of attacks on deep learning (higher level class of attacks compared to the next paragraph)

Privacy attacks can go from determining whether a data point was used for training (*membership inference*) (Shokri, Stronati, Song, & Shmatikov, 2017; Salem et al., 2019), reconstructing class representatives (*model inversion*) (Fredrikson, Jha, & Ristenpart, 2015), or inferring hidden properties of the training set (*properties inference*) (Ganju, Wang, Yang, Gunter, & Borisov, 2018).

The implementation of the GDPR in 2018 has put a question on whether those kinds of models are compliant (*question: can it be argued that if the model learns more information than it needs, then Art.5.1.c - data minimization applies?* I still need to write this part).

Membership inference attacks work the following way: given a data point, determine whether it was used to train the given model. The issue with such attacks is they are not always applicable in real-life settings as the attacker needs to get a copy of one of the training data before determining whether it belonged to the training set. Inversion model attacks operate such that, for a classifier, for example, a class representative can be created. While they are impressive for face recognition models, they only build the average representation of the class and do not work well if the data instances are not similar. Property inference attacks (PIA) aim at inferring properties learned by the model that are independent of the characteristics of the class the model is trained to recognize. There have been fewer studies on this kind of attack, ~~according to (He, Meng, Chen, Hu, & He, 2019) they only account for 28% of the published papers on model inversion attacks against 72% for membership inference attacks.~~ In (Ganju et al., 2018) the authors identify the invariance of fully connected neural networks to node permutation as a limiting factor of PIA efficiency and show that using a Deep Set architecture (Zaheer et al., 2017) is an effective way to tackle this limitation. In (Melis, Song, De Cristofaro, & Shmatikov, 2019) the authors successfully perform PIA in a collaborative learning setting where the attacker ~~is~~ one of the members.

The kind of model used has an impact on the kind of attacks that can be used on it (I still need to do that. *remark: I probably need to include a*

source here, it is a bit of a strong claim.). As computer vision and natural language processing popularity increases, more complex neural network architectures are created even if their vulnerability to privacy attacks is still not fully determined (He et al., 2019). None of the current work studied the relationship between the depth and width of a deep neural network and its sensibility to PIA. ~~Geiping, Bauermeister, Dröge, and Moeller~~ Geiping et al. analyzed the effect of model architecture on sensibility to model inversion attacks on non-sensible data (cars, animals).

In case of a trained Neural Network, are sensitive properties of the underlying dataset at risk due to property inference attacks, and if so, how does than risk relate to the architecture of the model such as depth and width?

3 Methodology

3.1 Threat model

We assume the target model is a deep neural network trained for a classification task on a sensitive dataset containing personal data. The attacker has access to the full architecture and parameters of the model (white-box setting). While a white-box model facilitates the job of the attacker, it is reasonable to make as it is common in some machine learning paradigms such as Federated Learning (Shokri & Shmatikov, 2015) where the clients, and sometimes the server, have full knowledge of the model. There also exists model extraction attacks with close to perfect performance (Papernot et al., 2016). The goal of the attacker is to infer general information about the training dataset. Such information could be the proportion of the training data having a property unrelated to the main classification task of the model.

3.2 Property Inference Attack

The main idea of a PIA is to train an attack model M_a which takes as input a trained model and outputs the probability that P is true for the training set of the input model. Once the attack model is trained, the attacker can give it the target model M_t as input and know whether the P is true for the dataset used to train the target model. To train M_a , the attacker uses several shadow models ($M_{s1} \dots M_{sk}$) which are respectively trained on datasets $D_1 \dots D_k$ specifically crafted by the attacker to contain or not the

target property P . Once the shadow models are trained, they are used as the training set for M_a . The general overview of the attack is described in Figure 1.

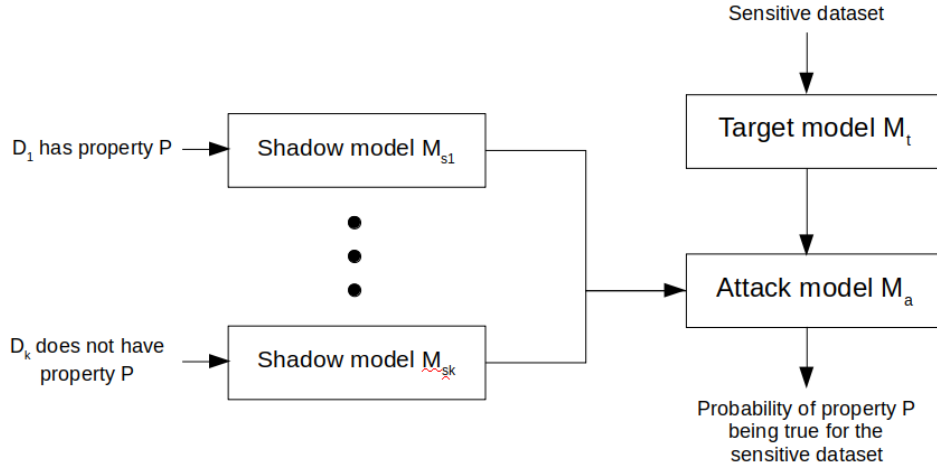


Figure 1: Property Inference Attack using shadow models.

3.3 Datasets

CelebFaces Attributes (CelebA) (Liu, Luo, Wang, & Tang, 2015) is a face attributes dataset containing more than 200 000 images of more than 10 000 celebrities taken from online. The images are labeled using 40 physical attributes such as hair color, smiling, wearing a hat. We use the dataset to train model M_t to detect whether the person is wearing glasses. To add other datasets if we end up using more.

25 different architectures of target models are trained using 5 width and 5 depth values. The width values are 8, 16, 32, 64, 128 (*to add: I dont know what reasonable width values to use at the moment*). The depth values are 10, 20, 40, 80, 160 (*to add: same, I am not sure what values yet. I saw 18, 20, 110 in (Li, Li, & Ribeiro, 2020) but it was not for face attributes detection*). The list of the architectures and their test accuracy is presented in the table

		depth				
		10	20	40	80	160
width	8					
	16					
	32					
	64					
	128					

Table 1: Target models accuracies on CelebA.

3.4 Evaluation

The two properties we try to infer are both sensitive according to the GDPR and are listed as follow:

- P_1 : whether the training set contains an imbalanced number of male/female (*question: maybe infer the exact proportion is a more impressive feat, it is probably more difficult*)
- P_2 : whether the training set contains an imbalanced number of light skin celebrities

to add: I want to describe which metrics I'm going to use to assess the attack.

References

- Fredrikson, M., Jha, S., & Ristenpart, T. (2015). Model inversion attacks that exploit confidence information and basic countermeasures. Proceedings of the ACM Conference on Computer and Communications Security, 2015-Octob, 1322–1333. doi: 10.1145/2810103.2813677
- Ganju, K., Wang, Q., Yang, W., Gunter, C. A., & Borisov, N. (2018). Property inference attacks on fully connected neural networks using permutation invariant representations. Proceedings of the ACM Conference on Computer and Communications Security, 619–633. doi: 10.1145/3243734.3243834
- Geiping, J., Bauermeister, H., Dröge, H., & Moeller, M. (2020). Inverting Gradients – How easy is it to break privacy in federated learning? (1). Retrieved from <http://arxiv.org/abs/2003.14053>
- He, Y., Meng, G., Chen, K., Hu, X., & He, J. (2019). Towards Privacy and Security of Deep Learning Systems: A Survey. , 1–23. Retrieved from <http://arxiv.org/abs/1911.12562>
- Li, J., Li, N., & Ribeiro, B. (2020). Membership Inference Attacks and Defenses in Supervised Learning via Generalization Gap. Retrieved from <http://arxiv.org/abs/2002.12062>
- Liu, Z., Luo, P., Wang, X., & Tang, X. (2015, December). Deep learning face attributes in the wild. In Proceedings of international conference on computer vision (iccv).
- Melis, L., Song, C., De Cristofaro, E., & Shmatikov, V. (2019). Exploiting unintended feature leakage in collaborative learning. Proceedings - IEEE Symposium on Security and Privacy, 2019-May, 691–706. doi: 10.1109/SP.2019.00029
- Papernot, N., McDaniel, P. D., Goodfellow, I. J., Jha, S., Celik, Z. B., & Swami, A. (2016). Practical black-box attacks against deep learning systems using adversarial examples. CoRR, abs/1602.02697. Retrieved from <http://arxiv.org/abs/1602.02697>
- Salem, A., Zhang, Y., Humbert, M., Berrang, P., Fritz, M., & Backes, M. (2019). ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models. (February). doi: 10.14722/ndss.2019.23119
- Shokri, R., & Shmatikov, V. (2015). Privacy-preserving deep learning. In Proceedings of the 22nd acm sigsac conference on computer and communications security (pp. 1310–1321).
- Shokri, R., Stronati, M., Song, C., & Shmatikov, V. (2017).

Membership Inference Attacks Against Machine Learning Models.
Proceedings - IEEE Symposium on Security and Privacy, 3–18. doi:
10.1109/SP.2017.41

Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R. R.,
& Smola, A. J. (2017). Deep sets. In Advances in neural information
processing systems (pp. 3391–3401).