

Knowledge Distillation-Based Model Extraction Attack using GAN-based Private Counterfactual Explanations

Fatima Ezzeddine¹, Omran Ayoub², Silvia Giordano²
¹ Università della Svizzera italiana, Lugano, Switzerland
² University of Applied Sciences and Arts of Southern Switzerland



Contribution

In this work, our contribution is:

- Novel Model Extraction Attack (MEA):** We propose a novel knowledge distillation (KD)-based MEA that exploit counterfactual explanations (CFs). We simulate an adversarial scenario where an attacker exploit CFs given by a Machine Learning as a Service (MLaaS).
- Privacy-Preserving CFs:** We introduce a novel technique to enhance the privacy of CFs generated by GAN-based models. We incorporate differential privacy (DP) into the GAN pipeline. We aim to mitigate the risk of privacy breaches while still providing meaningful explanations.

Specifically, we quantify:

- The effectiveness of KD-based MEA using agreement metrics.
- The quality of CFs generated using well-known metrics.

Results demonstrate that:

- Our proposed KD-based MEA outperforms the baseline.
- Our proposed private CFs method effectively preserves privacy against MEA while guaranteeing a specific level of privacy and CF quality.

Transparency in MLaaS

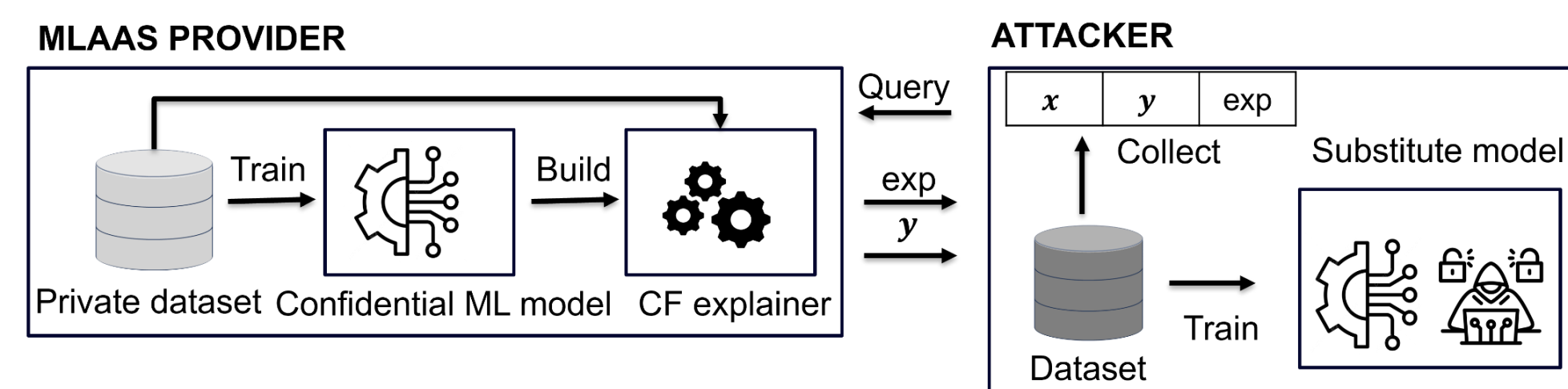
There is a notable increase in the deployment of MLaaS across production software applications.

ML models, demonstrably powerful, suffer from lack of **interpretability**. The absence of **transparency**, often referred to as the black box, undermines trust and urges the need for efforts to enhance their explainability.

MLaaS platforms now offer explanations alongside the ML prediction outputs and has elevated concerns regarding privacy, particularly in relation to **privacy leakage attacks** such as model extraction attacks (MEA).

Model Extraction Attack (MEA)

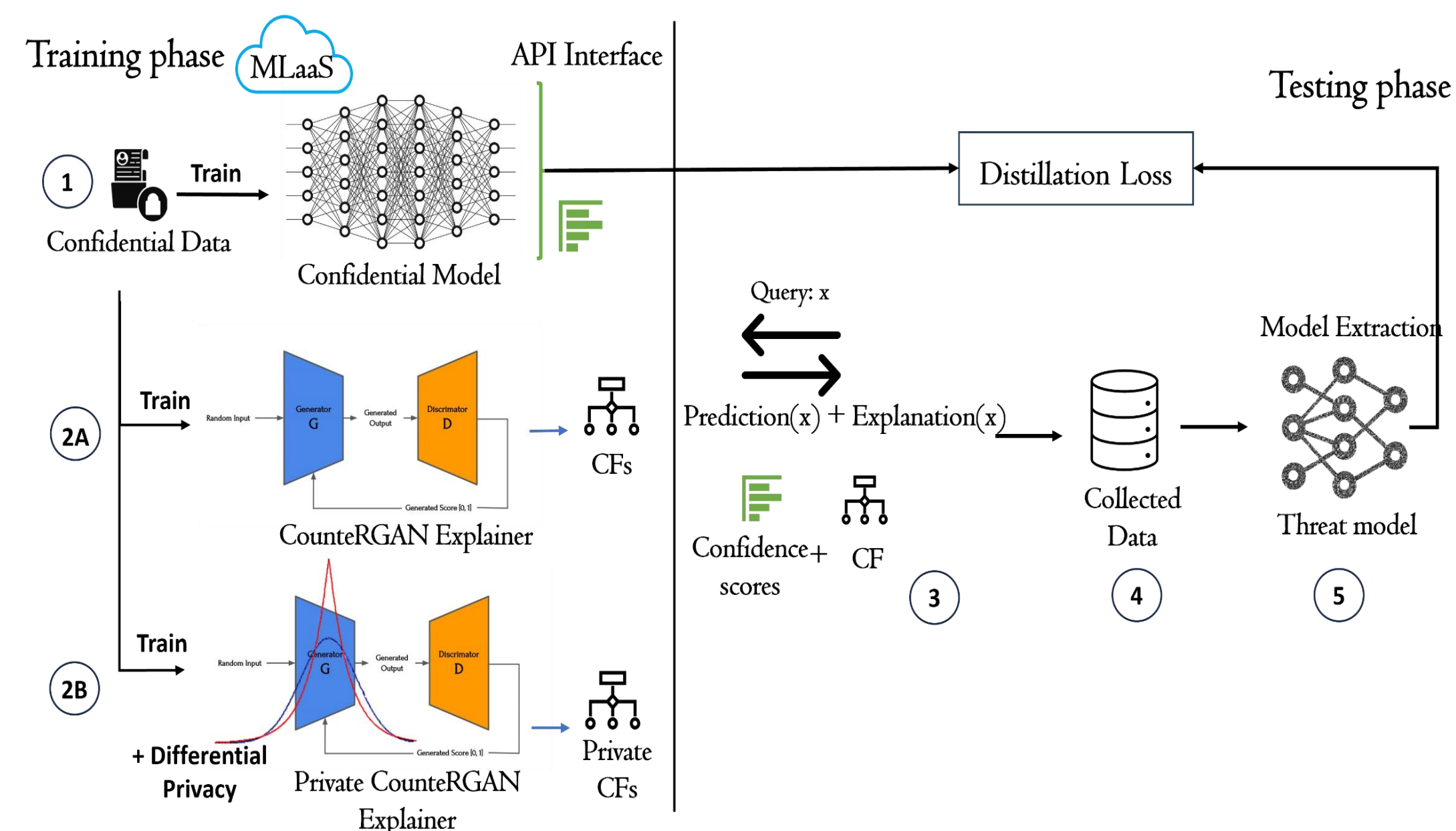
MEA derives a threat model t_x that closely resembles, in terms of functionality, the original model being targeted f_θ .



A user sends a **query**, which includes input data describing a data record x . Once the MLaaS API receives the user query, it performs the prediction $f_\theta(x)$ and generates a CF $c=E(x)$ where $f_\theta(c)$ has a different prediction and returns it along with its corresponding output to the user.

The attacker **extract (steal)** f_θ by training a **substitute (threat)** model t_x .

Knowledge Distillation-Based MEA



As an owner:

- Step 1:** train a DNN classifier f_θ on a private dataset.
- Step 2:** train **CounterGAN**¹ CF explainer and deploys as MLaaS.

As an attacker :

- Step 3:** query the model with **random queries** (with the assumption that the attacker does not have previous knowledge of the training set)
- Step 4:** collect a dataset to serve as input data to train t_x .
- Step 5:** apply **KD and** train t_x by minimizing the loss constituted by the threat model **classification loss** in addition to the **distillation loss**. We emphasize the importance of mimicking the output **probabilistic distribution** from the target model to the threat model.

$$loss = \alpha \cdot student_loss + (1 - \alpha) \cdot distillation_loss$$

$$JS(P||Q) = \frac{1}{2}KL\left(P\left\|\frac{P+Q}{2}\right.\right) + \frac{1}{2}KL\left(Q\left\|\frac{P+Q}{2}\right.\right)$$

Counterfactual Explanations generation with Differential Privacy

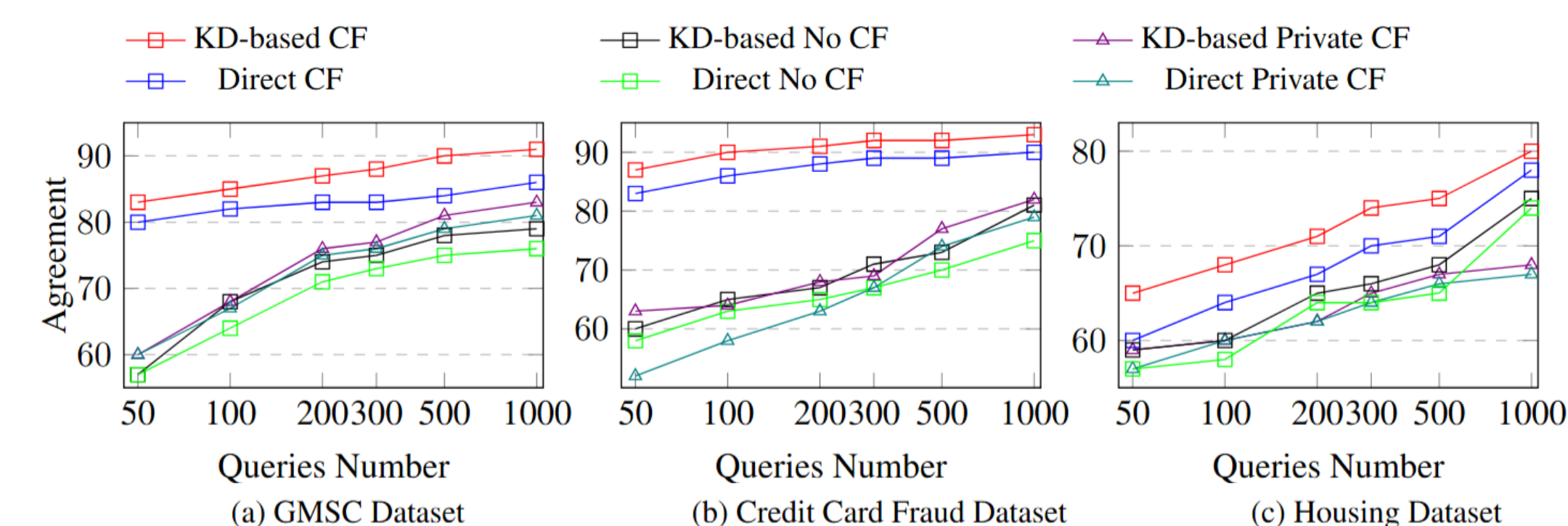
Objective: Prevent the generation of CFs that closely resemble the private training data, reducing the resemblance of CFs with the training data and mitigate the risks of privacy attacks.

Method: Inject DP into the generator during the optimization process. We employ the Adam DP optimizer. The process of DP Adam often involves multiple iterations of **adding noise** repeatedly over several rounds into the **gradients** of the **generator** in addition to a step of gradient **clipping**.

Results and Evaluation

Effectiveness of KD-based MEA

Agreement Metric: Predictions similarity between the ML models t_x and f_θ .



Main takeaways:

- Our proposed KD-based MEA **outperform** the **baseline** by achieving superior agreement levels.
- When CFs are employed, MEA requires significantly **few instances** to reach a specific high agreement level.
- The incorporation of **DP** can play a crucial role in **maintaining agreement** levels comparable to scenarios without CFs.

Impact of Incorporating DP in CF generator on quality of explanations

Metrics:

- Actionability:** Amount of modification of CFs.
- Realism:** If a data instance fits a known data distribution (reconst error of AE).
- Prediction Gain:** the probability changes of the CF explanation for a target class

Data	Prediction Gain		Actionability	
	CFs	Private CFs	CFs	Private CFs
GMSC	0.243 ± 0.011	0.121 ± 0.01	24.567 ± 0.364	16.981 ± 0.158
Credit Fraud	0.700 ± 0.084	0.445 ± 0.06	35.269 ± 0.328	10.507 ± 0.238
Housing	0.633 ± 0.052	0.678 ± 0.024	3.852 ± 0.053	1.004 ± 0.016

Data	Realism		
	random points	CFs	Private CFs
GMSC	15.649 ± 0.033	8.56 ± 0.142	15.723 ± 0.033
Credit Fraud	3.104 ± 0.019	3.072 ± 0.1647	4.73 ± 0.027
Housing	2.070 ± 0.04	1.356 ± 0.031	2.0 ± 0.01

Main takeaways:

- The initial random query points are inherently unrealistic (do not exhibit high realism).
- The private CF generation approach ensures this unrealistic nature is preserved when queried with a random data point.
- The integration of DP has an impact on prediction gain and actionability, constraining the explainer's progress toward the desired class.

References

Daniel Nemirovsky, Nicolas Thiebaud, Ye Xu, Abhishek Gupta. Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence, PMLR 180:1488-1497, 2022.

