# The GANfather: Controllable generation of malicious activity to improve defence systems

Ricardo Ribeiro Pereira\* Feedzai / University of Porto Portugal

> David Aparício University of Porto Portugal

Jacopo Bono Feedzai Portugal

Pedro Ribeiro University of Porto Portugal João Tiago Ascensão Feedzai Portugal

Pedro Bizarro<sup>†</sup>
Feedzai
Portugal

# **ABSTRACT**

Machine learning methods to aid defence systems in detecting malicious activity typically rely on labelled data. In some domains, such labelled data is unavailable or incomplete. In practice this can lead to low detection rates and high false positive rates, which characterise for example anti-money laundering systems. In fact, it is estimated that 1.7-4 trillion euros are laundered annually and go undetected. We propose *The GANfather*, a method to generate samples with properties of malicious activity, without label requirements. We propose to reward the generation of malicious samples by introducing an extra objective to the typical Generative Adversarial Networks (GANs) loss. Ultimately, our goal is to enhance the detection of illicit activity using the discriminator network as a novel and robust defence system. Optionally, we may encourage the generator to bypass pre-existing detection systems. This setup then reveals defensive weaknesses for the discriminator to correct. We evaluate our method in two real-world use cases, money laundering and recommendation systems. In the former, our method moves cumulative amounts close to 350 thousand dollars through a network of accounts without being detected by an existing system. In the latter, we recommend the target item to a broad user base with as few as 30 synthetic attackers. In both cases, we train a new defence system to capture the synthetic attacks.

#### **ACM Reference Format:**

Ricardo Ribeiro Pereira, Jacopo Bono, João Tiago Ascensão, David Aparício, Pedro Ribeiro, and Pedro Bizarro. 2023. The GANfather: Controllable generation of malicious activity to improve defence systems. In 4th ACM International Conference on AI in Finance (ICAIF '23), November 27–29, 2023, Brooklyn, NY, USA. ACM, New York, NY, USA, 8 pages. https://doi.org/10.1145/3604237.3626882

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICAIF '23, November 27-29, 2023, Brooklyn, NY, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 979-8-4007-0240-2/23/11...\$15.00 https://doi.org/10.1145/3604237.3626882

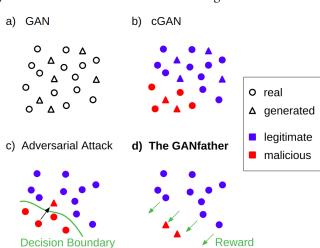


Figure 1: Comparison of our method to some widely used approaches. (a) GAN: a vanilla GAN setup does not require any labels, but one cannot choose the class of a generated sample since the distribution of the data is learned as a whole. (b) Conditional GAN (cGAN): using a cGAN, one learns the class-conditional distributions of the data, allowing the user to choose the class of a generated sample. However, labels are needed to train a cGAN. (c) Adversarial Attack (evasion): starting from a malicious example, perturbations are found such that a trained classifier is fooled and misclassifies the perturbed example. While labels are typically required to select the initial example as well as to train the classifier, eventually the adversarial attacks can be used to obtain a more robust classifier. (d) The GANfather: our method has some desirable properties from the three previous approaches: no labels are needed (as in a GAN), samples of a desired target class are generated (as in a cGAN) and a robust detection system can be trained (as in adversarial training). The combination of these properties in one framework is especially suitable in domains where no labelled data is available.

#### 1 INTRODUCTION

Digital systems' growing dominance in various aspects of our society opens up new opportunities for illicit actors. For example, digital banking enables clients to open bank accounts more easily but also facilitates complex money laundering schemes. It is estimated that

<sup>\*</sup>Corresponding author, ricardo.ribeiro@feedzai.com

<sup>&</sup>lt;sup>†</sup>Corresponding author, pedro.bizarro@feedzai.com

undetected money laundering activities worldwide accumulate to 1.7–4 trillion euros annually [16], while operational costs related to anti-money laundering (AML) compliance tasks incurred by financial institutions accumulate to \$37.1 billion [23]. Another example are recommender systems, which are often embedded in digital services to deliver personalised experiences. However, recommender systems may suffer from injection attacks whenever malicious actors fabricate signals (e.g., clicks, ratings, or reviews) to influence recommendations. These attacks have detrimental effects on the user experience. For example, a one-star decrease in restaurant ratings can lead to a 5 to 9 percent decrease in revenue [20].

The detection of such malicious attacks is challenging in the following aspects. In many cases, these illicit activities are adversarial in nature, where an attacker and a defence system adapt to each other's behaviour over time. Additionally, labelled datasets are unavailable or incomplete in certain domains due to the absence of natural labels and the cost of manual feedback. For example, besides the large amount of undetected money laundering, the investigation of detected suspicious activity is often far from trivial, resulting in a feedback delay that can last months.

To address these issues, we propose *The GANfather*, a method to generate examples of illicit activity and train effective detection systems without any labelled examples. Starting from unlabelled data, which we assume to be predominantly legitimate, the proposed method leverages a GAN-like setup [12] to train a generator which learns to create malicious activity, as well as a detection model (the discriminator) learning to distinguish between real data and synthetic malicious data.

To be able to generate samples with malicious properties from legitimate data, we propose to include an additional optimisation objective in the training loss of the generator. This objective is a use-case-specific, user-defined differentiable formulation of the goal of the malicious agents. Furthermore, our method optionally allows to incorporate an existing defence system, as long as a differentiable formulation is possible. In that case, we penalise the generator when triggering existing detection mechanisms. Our method can then actively find liabilities in an existing system while simultaneously training a complementary detection system to protect against such attacks.

Our method has some desirable properties that make it particularly well-suited for adversarial domains where no labelled data is available:

- No labelled malicious samples are needed. Here, we assume that our unlabelled data is predominantly of legitimate nature.
- Samples with features of malicious activity are generated. The key to generate such samples from legitimate data is to introduce an extra objective function that nudges the generator to produce samples with the required properties. We implicitly assume that malicious activity shares many properties with legitimate behaviour. We justify this assumption since attackers often mimic legitimate activity to some degree, in order to avoid raising suspicion or triggering existing detection systems.
- A robust detection system is trained. By training a discriminator to distinguish between the synthetic malicious samples and real data, we conjecture that the defence against a variety of real malicious attacks can be strengthened.

While each of these properties can be found separately in other methods, we believe that the combination of all the properties in a single method is novel and useful in the discussed scenarios. In Figure 1, we illustrate visually how our method distinguishes itself from some well-known approaches. Finally, while we only perform experiments on two use-cases (anti-money laundering and recommender systems) in the following sections, we believe that the suggested approach is applicable in other domains facing similar constraints, i.e., no labelled data and adversarial attacks, subject to domain-specific adaptations.

## 2 METHODS

We provide a general description of our proposed framework in Section 2.1. We proceed to describe two use-cases: anti-money laundering (AML) (Section 2.2) and detection of injection attacks in recommendation systems (Section 2.3). In Section 2.4, we show theoretically, in a simplified setting, how our generator's loss function changes the learning dynamics compared to a typical GAN.

## 2.1 General description

Figure 2 depicts the general structure of our framework. It comprises a generator, a discriminator, an optimisation objective, and, optionally, an existing alert system. Each component is discussed in more detail below.

**Generator.** As in the classical GAN architecture, the generator G receives a random noise input vector and outputs an instance of data. However, unlike classical GANs, the loss of the generator  $\mathcal{L}(G)$  is a linear combination of three components: the optimisation objective for malicious activity  $\mathcal{L}_{Obj}(G)$ , the GAN loss  $\mathcal{L}_{GAN}(G,D)$  that additionally depends on the discriminator D, and the loss from an existing detection system A,  $\mathcal{L}_{Alert}(G,A)$ :

$$\mathcal{L}(G) = \alpha \mathcal{L}_{Obi}(G) + \beta \mathcal{L}_{GAN}(G, D) + \gamma \mathcal{L}_{Alert}(G, A)$$
 (1)

where  $\alpha$ ,  $\beta$  and  $\gamma$  are hyperparameters to tune the strength of each component. The last term is optional, and if no existing detection system is present we simply choose  $\gamma = 0$ . Note also that one of the parameters is redundant and we tune only two parameters in our experiments (or one if  $\gamma = 0$ ).

We show in Section 2.4 that the stable point of convergence for the generator in our theoretical example moves away from the data distribution for any  $\alpha > 0$ .

**Discriminator.** The discriminator setup is the same as in a classical GAN. It receives an example and produces a score indicating the likelihood that the example is real or synthetic. Importantly, as explained in Section 2.4, the generator subject to Equation 1 will generate data increasingly out of distribution for larger  $\alpha$ . Therefore, we do not require the discriminator accuracy to fall to chance level at training convergence, as is usual with GANs. Instead, the discriminator may converge to perfect classification and may be used as a detection system for illicit activity. In our experiments, we use the Wasserstein loss [2] as our GAN loss.

**Malicious optimisation objective.** The optimisation objective quantifies how well the synthetic example is fulfilling the goal of a malicious agent. It can be a mathematical formulation or a differentiable model of the goal. This objective allows the generator to find previously unseen strategies to meet malicious goals.

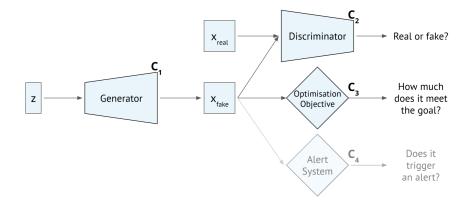


Figure 2: The GANfather framework. Its main components are a generator,  $C_1$ , which generates realistic attacks, a discriminator,  $C_2$ , which detects these attacks, and an optimisation objective,  $C_3$ , to incentivise the generation of malicious instances. Finally, our method optionally supports the inclusion of an existing alert system,  $C_4$ .

**Alert system.** If an existing, differentiable alert system is present, we can add it to our framework to teach the generator to create examples that do not trigger detection (see Equation 1). In that scenario, it is then beneficial for the discriminator to focus on the undetected illicit activity. Whenever the existing system is not differentiable, training a differentiable proxy may be possible.

Generator vs. Discriminator views. If required by the malicious optimisation objective, our generator can be adapted to generate samples which are only partially evaluated by the discriminator. For example, the layering stage of money laundering typically involves moving money through many financial institutions (FIs). However, each detection system operates within single institutions, limiting their view of the entire operation. Our method can be adapted to capture this situation, by generating samples containing various fictitious FIs, but only sending the partial samples corresponding to each FI to the discriminator. In recommender systems, the malicious objective can act on a group of synthetic illicit actors to generate coordinated attacks, while the detection of fraudulent users is typically performed on a single-user level.

**Architecture optimisations.** In the next sections, we provide more details about the specific architectures used in the two experiments. We note that the architecture details (layer types, widths and number of layers) were first optimised using a vanilla GAN setup (i.e. setting  $\alpha=0, \beta=1, \gamma=0$  in Equation 1). With the architecture fixed, the other hyperparameters were tuned as explained in the next sections.

**Code availability.** The Pytorch code for both models can be found on GitHub at https://github.com/feedzai/ganfather.

## 2.2 Anti-Money Laundering (AML)

We tackle the layering stage of money laundering, in which criminals attempt to conceal the origin of the money by moving large amounts across financial institutions through what are known as "mule accounts".

**Representing dynamic graphs as tensors.** To represent the dynamic graph of transactions, we can use a 3D tensor as depicted in Figure 3. We assume the nodes of the dynamic graph are accounts,

and the edges are transactions. The first two dimensions correspond to the weighted adjacency matrix of the accounts and the third dimension is time. We discretise the events into time windows of fixed length and group events that belong to the same entry in the tensor by summing their amounts. In other words, each entry  $A_{ijk}$  of the tensor corresponds to the cumulative amount sent between account i and account j on timestep k. Our representation covers any dynamic network with a 3D tensor whose size is fixed and prespecified, which allows us to avoid using recurrent models. While this approach may limit the size of generated data, domain experts reported that up to 95% of the money-laundering investigations involve cases containing up to 5 accounts.

**Architecture.** We implement the generator using a set of dense layers, followed by a set of transposed convolutions. Then, we create two branches: one to generate transaction amounts and the other to generate transaction probabilities. We use the probabilities to perform categorical sampling and generate sparse representations, similar to real transaction data. After the sampling step, the two branches are combined by element-wise multiplication, resulting in a final output tensor with the dimensions described above.

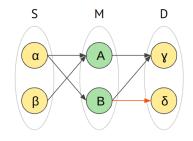
The discriminator receives two tensors with the same shape as inputs: one containing the total amount of money transferred per entry, and the other with the count information (mapping positive amounts to 1 and empty entries to 0). Each tensor passes through convolutional layers, followed by permutation-invariant operations over the internal and external accounts. Then, we concatenate both tensors. We reduce the dimensionality of the resulting vector to a classification outcome using dense layers.

We provide more details about both architectures on our GitHub repository.

**Money Mule objective.** To characterise the money flow behaviour of layering, where money is moved in and out of accounts while leaving little behind, we define the objective function as the geometric mean of the total amount of incoming  $(G(z)_{in})$  and outgoing  $(G(z)_{out})$  money per generated account (Equation 2).

$$\mathcal{L}_{Obj}(G) = -\int \sqrt{G(z)_{in} \times G(z)_{out}} \cdot p(z)dz$$
 (2)

Source	Target	Amount	Day
α	А	3.14	1
α	В	15.92	2
β	А	65.35	2
Α	γ	89.79	1
В	γ	32.38	1
В	δ	46.26	2
В	δ	43.38	2



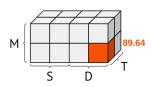


Figure 3: Data representation of transactional data. From the raw tabular data, we build the tripartite graph of the transactions, which is in turn represented as a 3D tensor. Here, S stands for Source accounts, M stands for Middle accounts, D stands for Destination accounts and T for the Time dimension.

Here z represents random noise input to the generator G and p(z) is its probability distribution. This objective encourages the generator to increase the amount of money sent and received per account and keep these two quantities similar, as observed in mule accounts.

Existing Alert System. In AML, it is common to have rule-based detection systems. In our case, the rules detection system contains five alert scenarios, capturing known suspicious patterns such as a sudden change in behaviour or rapid movements of funds. However, these rules are not differentiable, and our generator requires feedback in the form of a gradient. Hence, we construct a deep learning model as a proxy for the rules system. We hard-code a neural network mimicking the rules' logic operations by choosing the weights, biases and activation functions appropriately. This network gives the same feedback as the rules system would, but in a differentiable way.

## 2.3 Recommendation System

In this work, we consider collaborative filtering recommender systems. However, our method is compatible with any other differentiable recommender systems. The system receives a matrix of ratings R with shape  $(N_u, N_i)$ , where  $N_u$  is the number of users and  $N_i$  is the number of items. First, we compute cosine distances between users, resulting in the matrix D of shape  $(N_u, N_u)$ . Then, we compute the predicted ratings P as a matrix product between D and R. We decided to not represent time since most classical recommender systems do not account for it. However, it is possible to include temporal information using a similar setup to what we described in the AML use case. We also note that, unlike in the AML scenario, we do not have an existing detection system in this setup. We provide details about the architectures of both the generator and the discriminator on our GitHub repository.

**Injection Attack Objective.** We define the goal of malicious agents as increasing the frequency of recommendation of a specific item. The objective function in Equation 3 incentivizes the generator to increase the rating of the target item t for every user.

$$\mathcal{L}_{Obj}(G) = \int \sum_{i}^{N_u} \sum_{j}^{N_i} (P_{ij}(z) - P_{it}(z))_+ \cdot p(z) dz$$
 (3)

Here, the matrix of predicted ratings P depends on the random inputs z through the generator G and  $(\cdot)_+$  denotes a rectifier setting negative values to zero.

# 2.4 Theoretical justification

In this section, we provide a theoretical justification to enlighten certain aspects of our setup, in a simplified setting. We will assume no existing detection system is available ( $\gamma$  = 0 in Equation 1). In the case such a system would be available, we assume its effect is to limit how far the generated data distribution can be from the real data distribution. Furthermore, we assume that a malicious objective would promote a change in the distribution of at least one feature of the generated data compared to the real data.

In order to facilitate the analytical calculations, we make the following simplifying assumptions. Firstly, we assume that our data consists of only one feature, for which the regular (legitimate) activity is distributed following a normal distribution  $p_{\text{data}}$  with mean  $\mu_d$  and standard deviation  $\sigma_d$ :

$$p_{\text{data}} = \mathcal{N}\left(\mu_d, \sigma_d\right) \tag{4}$$

Secondly, we assume that we do not have any samples of malicious activity but that we know that it is characterised by larger values of this feature compared to the legitimate activity. Thirdly, we assume that the generated data follows a normal distribution  $p_{\rm gen}$  with mean  $\mu_g$  and standard deviation  $\sigma_g$ . Using  $\gamma=0$  and  $\beta=1-\alpha$  in Equation 1, assuming  $0\leq\alpha\leq1$ , we can write the training criterion of the generator as:

$$\mathcal{L}(G) = (1 - \alpha) \cdot (2 \cdot \text{JSD}(p_{\text{data}}|p_{\text{gen}}) - \log(4)) - \alpha \mu_q$$
 (5)

where the first term denotes the GAN loss [12] and the second term denotes our *malicious objective* rewarding the generator to produce samples with properties of the malicious data (i.e. increase the mean  $\mu_q$  as much as possible).

We can analytically solve the Jenson-Shannon Divergence (JSD) between the normal distributions, using  $\sigma_m^2 = \sigma_d^2 + \sigma_q^2$ ,

$$JSD (p_{\text{data}}|p_{\text{gen}}) = \frac{1}{2} KL (p_{\text{data}}|0.5 * (p_{\text{data}} + p_{\text{gen}}))$$

$$+ \frac{1}{2} KL (p_{\text{gen}}|0.5 * (p_{\text{data}} + p_{\text{gen}}))$$

$$= \frac{1}{2} \left[ \log \frac{\sigma_m}{\sigma_d} + \frac{\sigma_d^2 + (\mu_d - 0.5(\mu_d + \mu_g))^2}{2\sigma_m^2} - \frac{1}{2} + \log \frac{\sigma_m}{\sigma_g} + \frac{\sigma_g^2 + (\mu_g - 0.5(\mu_d + \mu_g))^2}{2\sigma_m^2} - \frac{1}{2} \right]$$

$$+ \log \frac{\sigma_m}{\sigma_g} + \frac{\sigma_g^2 + (\mu_g - 0.5(\mu_d + \mu_g))^2}{2\sigma_m^2} - \frac{1}{2}$$
(6)

From this, we can calculate the gradient w.r.t.  $\mu_q$ :

$$\frac{\partial \text{JSD}(p_{\text{data}}|p_{\text{gen}})}{\partial \mu_g} = \partial \left( \frac{1}{2} \left[ \log \frac{\sigma_m}{\sigma_d} + \frac{\sigma_d^2 + (\mu_d - 0.5(\mu_d + \mu_g))^2}{2\sigma_m^2} - \frac{1}{2} \right] + \log \frac{\sigma_m}{\sigma_g} + \frac{\sigma_g^2 + (\mu_g - 0.5(\mu_d + \mu_g))^2}{2\sigma_m^2} - \frac{1}{2} \right] \right) / \partial \mu_g$$

$$= \frac{\mu_g - \mu_d}{4\sigma_g^2 + 4\sigma_d^2} \tag{7}$$

Combining (5) and (7), we find that the gradient of the training objective of the generator w.r.t. the mean of the generated distribution  $\mu_q$  is

$$\frac{\partial \mathcal{L}(G)}{\partial \mu_g} = \frac{(1-\alpha)}{2} \frac{\mu_g - \mu_d}{\sigma_q^2 + \sigma_d^2} - \alpha \tag{8}$$

Without loss of generality, we set  $\sigma_q^2 + \sigma_{\text{data}}^2 = k/2$ , such that

$$\frac{\partial \mathcal{L}(G)}{\partial \mu_g} = (1 - \alpha) \frac{\mu_g - \mu_d}{k} - \alpha \tag{9}$$

Denoting  $\frac{\partial \mu_g}{\partial t}$  as the changes of  $\mu_g$  over time (i.e. a continuous version of the discrete gradient updates) and  $\eta$  as the learning rate, this leads to the following linear dynamical system which we can analyse in function of  $\mu_g$ ,  $\mu_d$  and the hyperparameter  $\alpha$ :

$$\begin{split} \frac{\partial \mu_g}{\partial t} &= -\eta \frac{\partial \mathcal{L}(G)}{\partial \mu_g} \\ &= -\eta (1 - \alpha) \frac{\mu_g - \mu_d}{k} + \eta \alpha \\ &= -\eta d\mu_g + \eta d\mu_d + \eta \alpha \end{split} \tag{10}$$

where we defined  $d = (1 - \alpha)/k$ . The stability of this linear system is defined by the sign of -d, which is always negative and hence the system has a stable fixed point. The stable fixed point for this dynamical system is easily found to be

$$\mu_g^{\star} = \mu_d + \frac{\alpha}{1 - \alpha} k \tag{11}$$

We plot the phase diagram of the dynamical system in Figure 4, showing the fixed point in function of the parameter  $\alpha$ .

From these simplified setting calculations, we can conclude that:

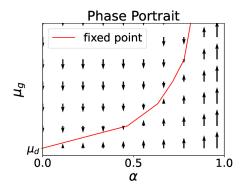


Figure 4: Phase portrait of our toy system. The fixed point of  $\mu_g$  depends on hyperparameter  $\alpha$ . For  $\alpha \to 1$ , the fixed point approaches infinity. For  $\alpha \to 0$ , the fixed point converges to  $\mu_d$ . Arrows denote the direction of the gradient  $\frac{\partial \mu_g}{\partial t}$ .

- For α > 0, our generated data will move away from the real data distribution and increasingly comply with the malicious objective.
- Different values of α will result in varying levels of deviation from the real data. In the absence of ground truth to evaluate the system, hyperparameter tuning and empirical testing are necessary.
- When generated data deviates from real data, the discriminator will increasingly achieve a perfect performance even at training completion. This is a major difference to standard GAN training.

#### 3 RESULTS

We evaluate the efficiency of *The GANfather* to generate and detect attacks in two use-cases: money laundering (Section 3.1) and recommendation system (Section 3.2).

## 3.1 Money Laundering

**Setup.** We use a real-world dataset of financial transactions, containing approximately 200,000 transactions, between 100,000 unique accounts, over 10 months<sup>1</sup>. Some of these accounts are labelled as suspicious of money laundering. We build a real test set of 5000 accounts, 184 of which are label positive (suspicious). We implement *The GANfather*'s generator and discriminator following the architectures presented in Section 2.2.

**Results.** We conduct a hyperparameter random search over the learning rate ( $[10^{-4}, 3 \times 10^{-3}]$ ) and the weights  $\alpha$  (set to 1),  $\beta$  ( $[10^2, 10^5]$ ) and  $\gamma$  ( $[10^3, 4 \times 10^3]$ ) mentioned in Equation 1. These ranges were empirically selected to allow convergence during training.

In Figure 5, we compare the distribution of money flows from such a generator compared to the real data distribution. We can observe that the generated samples successfully move more money than real data (up to 350,000 dollars vs. up to 9,000 dollars respectively). Interestingly, the distribution of amounts used is similar to real data, and the main difference is the number of transactions.

 $<sup>^{1}\</sup>mathrm{Due}$  to the confidential nature we cannot disclose the actual dataset.

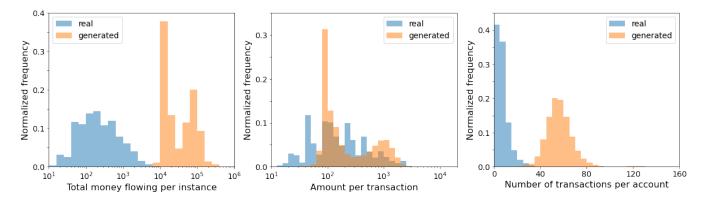


Figure 5: Comparing distributions of total money flow, amounts and counts between generated and real data.

Next, we test the detection performance of the trained discriminators on generated data. To detect potential bias in a discriminator trained solely on samples of the corresponding generator, we first build a *mixed* dataset, where synthetic malicious data is sampled from various generators. We combine this synthetic dataset with real data, and use it to evaluate the trained discriminators. Importantly, no retraining on this mixed dataset is performed. We observe that most discriminators can distinguish between real and generated examples with 100% accuracy, especially those trained with higher values of the  $\beta$  hyperparameter (see Equation 1, and note that in this experiment  $\alpha$  was fixed to a value of 1).

Then, we evaluate the detection performance on the real test set. We train a model DM with the same architecture as the discriminator using the mixed dataset mentioned in the previous paragraph. This training does not require real labels, since we use generated data as positive examples (suspicious) and assume that all real examples are negative (legitimate). After training, we evaluate three detection scenarios: the set of rules mentioned in Section 2.2; the model DM, with the threshold tuned to match the alert rate of the rules<sup>2</sup>; a combination of both (alert if either of them triggers). The results are shown in Table 1. We see that, even though the model DM was trained using only generated data as positive examples, it achieves better detection performance than the rules. Furthermore, only 10 of the 128 alerts of the Rules+Model scenario were alerted by both detection systems, and the true positives had little overlap as well (5 out of 54). This means that, by including the rules' feedback in the loss of the generator, it learns to create synthetic examples that are not captured by the rules but are similar to real examples of suspicious activity. As such, a model trained with those synthetic examples can be used to complement the rules, with the advantage of being easy to tune to a desired alert rate.

	Alert Rate %	Recall %	Precision %
Rules	1.4	13.6	36.2
Model	1.4	18.5	49.3
Rules + Model	2.6	29.3	42.2

Table 1: Detection of real labels.

## 3.2 Recommender System

**Setup.** We use the MovieLens 1M dataset<sup>3</sup>, comprised of a matrix of 6040 users and 3706 movies, with ratings ranging from 1 to 5 [13]. We implement the generator and discriminator and collaborative filtering recommender system as described in Section 2.3. To compute the predicted ratings, during training we take a weighted average of ratings considering all users in the dataset. We consider all users during training because the initially generated ratings are random, and only providing feedback from the top-N closest users limits the strategies that the generator can learn. In contrast, we consider the top-400 closest neighbours to compute predicted ratings at inference since we observed empirically that this value produces the lowest recommendation loss.

In this scenario, we do not use an existing detection component, corresponding to  $\gamma=0$  in Equation 1. We train our networks with 300 synthetic attackers but evaluate the generator's ability to influence the recommender system with injection attacks of various sizes. We also define four baseline attacks: (1) a rating of 5 for the target movie and 0 otherwise, (2) a rating of 5 for the target movie and ~90 random ratings for randomly chosen movies, (3) a rating of 5 for the target movie and ~90 random ratings for the top 10% highest rated movies, (4) a rating of 5 for the target movie and ~90 random ratings for the top 10% most rated movies.

**Results.** We choose  $\beta = 1 - \alpha$  in Equation 1, with  $0 \le \alpha \le 1$  and perform a hyperparameter search over  $\alpha$ . We observe that increasing  $\alpha$  leads to generators whose attacks increasingly recommend the target movie, at the cost of moving further away from the rating distributions of real profiles.

In Table 2, we show how many real users have the target movie in their top-10 recommendations, depending on the number of generated users that we inject and how they were generated (through *The GANfather* or the described baselines). We observe that even with a very limited proportion of generated users (30 among 6040 real users, 0.5%), they are able to greatly influence many real users (3.7%). In contrast, the baselines have very small impact on the recommendations of real users. Lastly, as expected, increasing the number of injected users increases the target movie's recommendation frequency to real users.

<sup>&</sup>lt;sup>2</sup>We assume that the rules are fixed, so we cannot tune the number of their alerts.

<sup>3</sup>https://www.kaggle.com/datasets/odedgolden/movielens-1m-dataset

Generation strategy	30 users	60 users	120 users
The GANfather	225	290	428
Baseline 1	0	0	0
Baseline 2	0	0	0
Baseline 3	1	3	7
Baseline 4	0	0	0

Table 2: Number of real users with the target movie in their top-10 recommendations, after injecting 30, 60, or 120 generated users.

Finally, we analyse the detection of synthetic attacks. As in the AML scenario we build a test set containing real and synthetic data, where the synthetic data contains a mixture of samples from various trained generators to identify the possible bias of a discriminator to attacks by the corresponding generator. We then quantify the AUC of the trained discriminators. We observe that most discriminators trained in a GAN setting (taking turns with a generator to update their weights) achieve around 0.75 AUC. Unlike the AML scenario, this suggests that the discriminators are tuned to detect synthetic data from their respective generators, but less so from other generators. If instead we build a *mixed* training set combining real samples with synthetic data from various generators and use it to retrain a discriminator, it achieves near-perfect classification (above 0.99 AUC).

## 4 RELATED WORK

Controllable data generation. Wang et al. [31] review controllable data generation with deep learning. Among the presented works, we highlight [8]. It leverages a GAN trained with reinforcement learning to generate small molecular graphs with desired properties. Their work is similar to ours in that we both (1) extend a GAN with an extra objective and (2) use similar data representations, namely sparse tensors. However, whereas [8] uses a labelled dataset of molecules and their chemical properties, our method does not rely on any labelled data.

Adversarial Attacks. A vast amount of literature exists on the generation of adversarial attacks (see [36] for a recent review). Such attacks have been studied in various domains and using various setups (e.g. cybersecurity evasion using reinforcement learning [1], intrusion detection evasion using GANs [30], sentence sentiment misclassification using BERT [11]). In all cases, a requirement is that labelled examples of malicious attacks exist.

Anti-Money Laundering. Typical anti-money laundering solutions are rule-based [25, 33, 34]. However, rules suffer from high false-positive rates, may fail to detect complex schemes, and are costly to maintain. Machine learning-based solutions tackle these problems [7]. Given the lack of labelled data, most solutions employ unsupervised methods like clustering [26, 32], and anomaly detection [5, 10]. These assume that illicit behaviours are rare and distinguishable, which may not hold whenever money launderers mimic legitimate behaviour. Various supervised methods have been explored [14, 21, 22, 24, 29], but most of these works use synthetic positive examples or incompletely labelled datasets. To avoid this, Lorenz et al. [19] propose efficient label collection with active learning. Deng et al. [9] and [6] explore data augmentation using

conditional GANs. Lastly, Li et al. [17] and Sun et al. [27] propose a metric to detect dense money flows in large transaction graphs, resulting in an anomaly score. Their method does not involve training of a classifier, and instead relies on generating many subsets of nodes and iteratively calculating the anomaly score.

Recommender systems (RS) injection attacks. Most injection attacks on RS are hand-crafted according to simple heuristics. Examples include random and average attacks [15], bandwagon attacks [3] and segmented attacks [4]. However, these strategies are less effective and easily detectable as most generated rating profiles differ significantly from real data and correlate with each other. Tang et al. [28] address the optimisation problem of finding the generated profiles that maximise their goals directly through gradient descent and a surrogate RS. Some studies apply GANs to RS to generate attacks and defend the system. Wu et al. [35] combines a graph neural network (GNN) with a GAN to generate their attack. The former selects which items to rate, and the latter decides the ratings. Zhang et al. [37] and Lin et al. [18] propose a similar setup to ours in which they train a GAN to generate data and add a loss function to guide the generation of rating profiles. The main differences to our work are the usage of template rating profiles to achieve the desired sparsity, the chosen architecture and loss functions. In our work, sparsity is learned by the generator through the categorical sampling branch (see Section 2). Moreover, our method allows the generation of coordinated group attacks by generating multiple attackers from a single noise vector.

#### 5 CONCLUSION

In this work, we propose *The GANfather* to generate data of a class for which no labelled examples are available (malicious activity), while simultaneously training a detection network to classify this class correctly.

We performed experiments in two domains. In the anti-money laundering setting, the generated attacks are able to move up to 350,000 dollars using just five internal accounts, and without triggering an existing detection system. Then we show that for a real-world labelled dataset, a model trained with these generated attacks can be used to complement the rules, alerting previously undetected suspicious activity. In the recommender system setting, we generate attacks that are substantially more successful at recommending the target item than naive baselines. Then, we train a near-perfect classifier to detect the synthetic malicious activity. While a real test in a deployment scenario is lacking and should be addressed in future work, we believe our current experiments provide a proof of value of the method. In these experiments, our method generates a variety of successful attacks, and we therefore believe it can be a valuable method to improve the robustness of defence systems.

The limitations of our method lie in its assumptions. Firstly, we assume that the unlabelled data is dominated by legitimate events, and our method would not work in settings where this is not the case. Secondly, we assume that we can quantify the malicious objective in terms of available features. In this case, one could argue we can just use the malicious objective as a detection score. However, the detection system often has a (much) smaller view than the malicious objective. For example, anti-money laundering systems only view incoming and outgoing transactions for *one* financial

institution. However, our objective can be adapted to generate malicious activity mimicking flows across *multiple* synthetic financial institutions, while keeping the view of the discriminator on an individual institution level. Thirdly, while our method does not prevent generated data to be very different from real data, we argue that the strength of our method is in generating more subtle attacks that are not immediately distinguishable from real data. Finally, while we chose the malicious objectives to be as simple as possible in our proof of concept experiments, there is no restriction to make them more complex as long as they are differentiable.

To conclude, our method fits the adversarial game between criminals and security systems by simulating various meaningful attacks. If existing defences are in place, our method may learn to avoid them and, eventually, train a complementary model. Incorporating machine learning models into the detection system typically enhances the detection of illicit activity by triggering more precise alerts, while being easier to fine-tune and maintain. We believe our work contributes to increase the robustness of detection methods of illicit activity.

## REFERENCES

- Giovanni Apruzzese, Mauro Andreolini, Mirco Marchetti, Andrea Venturi, and Michele Colajanni. 2020. Deep Reinforcement Adversarial Learning Against Botnet Evasion Attacks. *IEEE Transactions on Network and Service Management* 17, 4 (2020), 1975–1987. https://doi.org/10.1109/TNSM.2020.3031843
- [2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein generative adversarial networks. In *International conference on machine learning*. PMLR, 214–223.
- [3] Robin Burke, Bamshad Mobasher, and Runa Bhaumik. 2005. Limited knowledge shilling attacks in collaborative filtering systems. In Proceedings of 3rd international workshop on intelligent techniques for web personalization (ITWP 2005), 19th international joint conference on artificial intelligence (IJCAI 2005). 17–24.
- [4] Robin Burke, Bamshad Mobasher, Runa Bhaumik, and Chad Williams. 2005. Segment-based injection attacks against collaborative filtering recommender systems. In Fifth IEEE International Conference on Data Mining (ICDM'05). IEEE, 4-pp.
- [5] Ramiro Daniel Camino, Radu State, Leandro Montero, and Petko Valtchev. 2017. Finding suspicious activities in financial transactions and distributed ledgers. In 2017 IEEE International Conference on Data Mining Workshops (ICDMW). IEEE, 787–796.
- [6] Charitos Charitou, Simo Dragicevic, and Artur d'Avila Garcez. 2021. Synthetic Data Generation for Fraud Detection using GANs. arXiv preprint arXiv:2109.12546 (2021)
- [7] Zhiyuan Chen, Ee Na Teoh, Amril Nazir, Ettikan Kandasamy Karuppiah, Kim Sim Lam, et al. 2018. Machine learning techniques for anti-money laundering (AML) solutions in suspicious transaction detection: a review. Knowledge and Information Systems 57, 2 (2018), 245–285.
- [8] Nicola De Cao and Thomas Kipf. 2018. MolGAN: An implicit generative model for small molecular graphs. arXiv preprint arXiv:1805.11973 (2018).
- [9] Xinwei Deng, V Roshan Joseph, Agus Sudjianto, and CF Jeff Wu. 2009. Active learning through sequential design, with applications to detection of money laundering. J. Amer. Statist. Assoc. 104, 487 (2009), 969–981.
- [10] Zengan Gao. 2009. Application of cluster-based local outlier factor algorithm in anti-money laundering. In 2009 International Conference on Management and Service Science. IEEE, 1–4.
- [11] Siddhant Garg and Goutham Ramakrishnan. 2020. BAE: BERT-based Adversarial Examples for Text Classification. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). Association for Computational Linguistics, Online, 6174–6181. https://doi.org/10.18653/v1/2020.emnlp-main.498
- [12] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. Advances in neural information processing systems 27 (2014).
- [13] F Maxwell Harper and Joseph A Konstan. 2015. The movielens datasets: History and context. Acm transactions on interactive intelligent systems (tiis) 5, 4 (2015), 1–19
- [14] Martin Jullum, Anders Løland, Ragnar Bang Huseby, Geir Ånonsen, and Johannes Lorentzen. 2020. Detecting money laundering transactions with machine learning. Journal of Money Laundering Control (2020).
- [15] Shyong K Lam and John Riedl. 2004. Shilling recommender systems for fun and profit. In Proceedings of the 13th international conference on World Wide Web.

- 393-402.
- [16] Karel Lannoo and Richard Parlour. 2021. Anti-Money Laundering in the EU: Time to get serious. CEPS Task Force Report 28 Jan 2021. http://aei.pitt.edu/103318/
- [17] Xiangfeng Li, Shenghua Liu, Zifeng Li, Xiaotian Han, Chuan Shi, Bryan Hooi, He Huang, and Xueqi Cheng. 2020. Flowscope: Spotting money laundering based on graphs. In Proceedings of the AAAI Conference on Artificial Intelligence.
- [18] Chen Lin, Si Chen, Meifang Zeng, Sheng Zhang, Min Gao, and Hui Li. 2022. Shilling Black-Box Recommender Systems by Learning to Generate Fake User Profiles. IEEE Transactions on Neural Networks and Learning Systems (2022).
- [19] Joana Lorenz, Maria Inês Silva, David Aparício, João Tiago Ascensão, and Pedro Bizarro. 2020. Machine learning methods to detect money laundering in the Bitcoin blockchain in the presence of label scarcity. arXiv preprint arXiv:2005.14635 (2020)
- [20] Michael Luca. 2016. Reviews, Reputation, and Revenue: The Case of Yelp.com. American Economic Journal - Applied Economics (2016).
- [21] Lin-Tao Lv, Na Ji, and Jiu-Long Zhang. 2008. A RBF neural network model for anti-money laundering. In 2008 International Conference on Wavelet Analysis and Pattern Recognition, Vol. 1. IEEE, 209–215.
- [22] Catarina Oliveira, João Torres, Maria Inês Silva, David Aparício, João Tiago Ascensão, and Pedro Bizarro. 2021. GuiltyWalker: Distance to illicit nodes in the Bitcoin network. arXiv preprint arXiv:2102.05373 (2021).
- [23] Arin Ray. 2021. IT and operational spending in AML-KYC: 2021 edition. https://www.celent.com/insights/428901357
- [24] Saleha Raza and Sajjad Haider. 2011. Suspicious activity reporting using dynamic bayesian networks. Procedia Computer Science 3 (2011), 987–991.
- [25] David Savage, Qingmai Wang, Pauline Chou, Xiuzhen Zhang, and Xinghuo Yu. 2016. Detection of money laundering groups using supervised learning in networks. arXiv preprint arXiv:1608.00708 (2016).
- [26] Reza Soltani, Uyen Trang Nguyen, Yang Yang, Mohammad Faghani, Alaa Yagoub, and Aijun An. 2016. A new algorithm for money laundering detection based on structural similarity. In 2016 IEEE 7th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON). IEEE, 1–7.
- [27] Xiaobing Sun, Jiabao Zhang, Qiming Zhao, Shenghua Liu, Jinglei Chen, Ruoyu Zhuang, Huawei Shen, and Xueqi Cheng. 2021. CubeFlow: Money Laundering Detection with Coupled Tensors.. In PAKDD (1). Springer, 78–90.
- [28] Jiaxi Tang, Hongyi Wen, and Ke Wang. 2020. Revisiting adversarially learned injection attacks against recommender systems. In Fourteenth ACM conference on recommender systems. 318–327.
- [29] Jun Tang and Jian Yin. 2005. Developing an intelligent data discriminating system of anti-money laundering based on SVM. In 2005 International conference on machine learning and cybernetics, Vol. 6. IEEE, 3453–3457.
- [30] Muhammad Usama, Muhammad Asim, Siddique Latif, Junaid Qadir, and Ala-Al-Fuqaha. 2019. Generative Adversarial Networks For Launching and Thwarting Adversarial Attacks on Network Intrusion Detection Systems. In 2019 15th International Wireless Communications and Mobile Computing Conference (IWCMC). 78–83. https://doi.org/10.1109/IWCMC.2019.8766353
- [31] Shiyu Wang, Yuanqi Du, Xiaojie Guo, Bo Pan, and Liang Zhao. 2022. Controllable Data Generation by Deep Learning: A Review. arXiv preprint arXiv:2207.09542 (2022).
- [32] Xingqi Wang and Guang Dong. 2009. Research on money laundering detection based on improved minimum spanning tree clustering and its application. In 2009 Second international symposium on knowledge acquisition and modeling, Vol. 2. IEEE, 62–64.
- [33] R Cory Watkins, K Michael Reynolds, Ron Demara, Michael Georgiopoulos, Avelino Gonzalez, and Ron Eaglin. 2003. Tracking dirty proceeds: exploring data mining technologies as tools to investigate money laundering. *Police Practice* and Research 4, 2 (2003), 163–178.
- [34] Mark Weber, Jie Chen, Toyotaro Suzumura, Aldo Pareja, Tengfei Ma, Hiroki Kanezashi, Tim Kaler, Charles E Leiserson, and Tao B Schardl. 2018. Scalable graph learning for anti-money laundering: A first look. arXiv preprint arXiv:1812.00076 (2018)
- [35] Fan Wu, Min Gao, Junliang Yu, Zongwei Wang, Kecheng Liu, and Xu Wang. 2021. Ready for emerging threats to recommender systems? A graph convolution-based generative shilling attack. *Information Sciences* 578 (2021), 683–701.
- [36] Han Xu, Yao Ma, Hao-Chen Liu, Debayan Deb, Hui Liu, Ji-Liang Tang, and Anil K. Jain. 2020. Adversarial Attacks and Defenses in Images, Graphs and Text: A Review. International Journal of Automation and Computing (2020), 151–178.
- [37] Xuxin Zhang, Jian Chen, Rui Zhang, Chen Wang, and Ling Liu. 2021. Attacking recommender systems with plausible profile. IEEE Transactions on Information Forensics and Security 16 (2021), 4788–4800.