

Generating synthetic data in finance: opportunities, challenges and pitfalls

Samuel A. Assefa*
samuel.a.assefa@jpmorgan.com
J.P. Morgan AI Research
New York City, NY, USA

Danial Dervovic*
danial.dervovic@jpmchase.com
J.P. Morgan AI Research
London, UK

Mahmoud Mahfouz
mahmoud.a.mahfouz@jpmchase.com
J.P. Morgan AI Research
London, UK

Robert E. Tillman
robert.e.tillman@jpmchase.com
J.P. Morgan AI Research
New York City, NY, USA

Prashant Reddy
prashant.reddy@jpmchase.com
J.P. Morgan AI Research
New York City, NY, USA

Manuela Veloso[†]
manuela.veloso@jpmchase.com
J.P. Morgan AI Research
New York City, NY, USA

ABSTRACT

Financial services generate a huge volume of data that is extremely complex and varied. These datasets are often stored in silos within organisations for various reasons, including but not limited to regulatory requirements and business needs. As a result, data sharing within different lines of business as well as outside of the organisation (e.g. to the research community) is severely limited. It is therefore critical to investigate methods for synthesising financial datasets that follow the same properties of the real data while respecting the need for privacy of the parties involved.

This introductory paper aims to highlight the growing need for effective synthetic data generation in the financial domain. We highlight three main areas of focus that are of particular importance while generating synthetic financial datasets: 1) Generating realistic synthetic datasets. 2) Measuring the similarities between real and generated datasets. 3) Ensuring the generative process satisfies any privacy constraints.

Although these challenges are also present in other domains, the additional regulatory and privacy requirements within financial services present unique questions that are not asked elsewhere. Due to the size and influence of the financial services industry, answering these questions has the potential for a great and lasting impact. Finally, we aim to develop a shared vocabulary and context for generating synthetic financial data using two types of financial datasets as examples.

CCS CONCEPTS

• **Security and privacy** → **Privacy-preserving protocols**; • **Computing methodologies** → *Agent / discrete models*; **Simulation evaluation**; • **Applied computing** → **Economics**.

*Equal contribution to this work.

[†]Also with Carnegie Mellon University, School of Computer Science.

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 ACM 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 '20, October 15–16, 2020, New York, NY, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7584-9/20/10...\$15.00

<https://doi.org/10.1145/3383455.3422554>

KEYWORDS

Synthetic data, Simulation, Privacy preserving data generation

ACM Reference Format:

Samuel A. Assefa, Danial Dervovic, Mahmoud Mahfouz, Robert E. Tillman, Prashant Reddy, and Manuela Veloso. 2020. Generating synthetic data in finance: opportunities, challenges and pitfalls. In *ACM International Conference on AI in Finance (ICAIF '20)*, October 15–16, 2020, New York, NY, USA. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3383455.3422554>

1 INTRODUCTION

Financial data contains some of the most sensitive and personally identifiable attributes of customers. Using and sharing such data, especially for research purposes outside of the organisations that generate it, is severely restricted. One approach to address this limitation is the generation of synthetic data. The primary directive in generating synthetic financial data is therefore protecting the privacy of customers and entities involved in generating a particular synthetic dataset. This directive has been enshrined into law in various jurisdictions, notably the GDPR laws in the European Union [44], and in the United States FERPA [75] and HIPAA [23], relating to educational and medical data privacy respectively. The issue of data privacy is of great importance in public opinion, as evidenced by reactions to the Facebook/Cambridge Analytica scandal [77] and the numerous data breaches [19] that have occurred over the last few decades, alongside the financial markets' response to these breaches [47].

We define synthetic data as data obtained from a generative process that learns the properties of the real data. Such processes are strictly different from the most commonly used data obfuscation techniques (e.g. anonymisation or removing certain sensitive attributes) as our intention is to synthesise new samples that are related to but can not be mapped back to the real data. Some of the requirements for such generative processes include:

- Capability for generating many different data types, such as numeric, binary and categorical, as well as complex data like images.
- The process should be able to generate an arbitrary number of data points, with as many features as desired to good fidelity.
- The privacy characteristics of the dataset can be precisely tuned against how realistic the data is.

The underlying data representations and structures also vary significantly based on the specific domain of interest ranging from structured tabular formats and time series data to networks, unstructured text and images.

We shall concentrate primarily on considerations of *privacy* over *security*, that is, we are in a setting where we wish to share a dataset without compromising information about any given entity within the dataset. An attacker is assumed to have full access to the synthetic dataset, as well as to the generating algorithm, but not the original dataset.

1.1 Motivation

Some of the basic use cases and motivations for synthetic data generation in finance are highlighted below.

Internal data use restrictions. Regulatory requirements may prevent data sharing between different lines of business within a company. Alternatively, teams may wish to begin working with data before the relevant approvals have been made.

Lack of historical data. There is a limited amount of historical data to study certain events (e.g. flash crashes in the market, recessions, new regimes of behaviour) that make studying the underlying mechanisms very challenging. It is useful in various such settings to have counterfactual data for testing strategies and inferences.

Tackling class imbalance. For use cases such as fraud detection, the datasets are usually highly imbalanced, and traditional machine learning and anomaly detection techniques often fail. Realistic synthetic data, along with appropriate data imputation techniques offer a promising approach to tackle this challenge.

Training advanced Machine Learning models. Large scale advanced machine learning (e.g. deep learning) is often carried out using cloud services, requiring compute resources and vast quantities of training data. Institutions may not be able to upload training data to these services for a number of reasons. Synthetic data can be used to train models, which can then be brought back on premise to be used on real data. Moreover, training on synthesised data offers some protection from “membership inference attacks”, wherein model parameters can be used to extract training data.

Data sharing. By sharing data between institutions and within the research community, better solutions can be found for technical problems faced by financial institutions. Sharing of realistic synthetic data allows financial institutions to do this in a way that satisfies their data sharing restrictions.

2 EXAMPLES OF FINANCIAL DATA

For the purposes of this paper, we consider two broad classes of financial data, *retail banking* and *market microstructure* data. Retail financial data arises from operations facing the general public. This includes, but is not limited, to transaction data, loan applications, customer service logs, etc. Market microstructure data refers to the data maintained and distributed by exchanges detailing historical limit order books for a given financial asset. Generating synthetic

data in each domain presents its own challenges, which we go through in turn in the following sections.

2.1 Retail Banking Data

2.1.1 Data Types. A famous example of a retail-type dataset is the UCI Adult dataset [28], where the binary feature of whether a particular person has an income of over \$50k/year is predicted by various census information, such as age, profession, marital status and gender. This data set is tabular, mimicking how data is typically stored in a relational database. Most examples of synthesised data sets are also tabular.

In retail banking there are also many examples of datasets that are graph-based, that is, they are most naturally described by a graph $G = (V, E)$. The nodes V and the edges E may be labelled and or weighted and E may consist of directed edges, undirected edges or a combination of both. Multiple edges between vertices are not uncommon. Notable applications include modelling the inter-dependencies between financial institutions [16, 42] and modeling payments networks for anti-money laundering [36] and fraud [66].

Other types of common data are JSON [24] structured documents and natural language data, as their use in finance is similar enough to that in other domains not to warrant inclusion in this review.

2.1.2 Privacy Model. There have been a number of technical approaches suggested for protecting privacy in retail-type data. The two most prominent approaches are *k-anonymity* [72] and *differential privacy* [30]. A dataset is said to have the *k-anonymity* property if the information for each entity contained in the dataset cannot be distinguished from $k - 1$ entities also appearing in the release. Differential privacy is a technical condition satisfied by a randomised algorithm ensuring that the distribution of the output of the algorithm is bounded when applied to two “adjacent” datasets, with the definition of adjacent being dependent on the application. In the most typical case, tabular datasets are adjacent when they differ by a single row. Typically, a differentially private algorithm is parameterized by the pair (ϵ, δ) , where smaller values of ϵ denote greater privacy and δ denotes a failure probability. Differential privacy enjoys the property of *composability*, whereby differentially private algorithms can be chained with their output remaining differentially private. Moreover, differentially private algorithms are provably safe against attackers with side information [30]. The *k-anonymity* paradigm suffers from the *curse of dimensionality* [6], in that one has to destroy a rapidly increasing fraction of a particular dataset as the number of columns in the table grows to ensure *k-anonymity*. As such, differential privacy is the most popular technical solution for releasing data privately. Data synthesis methods using differential privacy for tabular data are well covered in the review by Bowen and Liu [14].

One must also be cognizant of data *deanonymization*: whence an anonymized dataset is combined with existing data to personally identify individuals in the dataset. This has been demonstrated by Daries et al. [25] where data from massive open online courses (MOOCs) was deanonymized and by Archie et al. [8], who deanonymized the Netflix prize dataset using Amazon review data. Indeed, the US Census is recognised as having been vulnerable to such attacks in previous iterations [51].

2.2 Market Microstructure Data

2.2.1 Data Types. Market microstructure data typically come in the form of time series describing, for example, a stock price over time. The granularity of the data is dependent on the frequency of the trading activity by the market participants. In the past decade, the rise of algorithmic trading, and specifically, high frequency trading has resulted in a significant increase in the amount of data available for research. However, access to such data sets is very limited and therefore an effort to synthesize such data sets using real market data is needed.

Of particular interest to the research community is limit order book data. A limit order book is used by exchanges to match buyers and sellers of a particular security [40]. It is an electronic record of the outstanding orders in the market and represents a snapshot in time describing the supply and demand of the security. It is based on the continuous double auction mechanism whereby participants can submit both buy and sell orders and expect their trades to match instantaneously if a corresponding trade on the opposing side is present. Exchanges offer various order types. The two main types are market and limit orders. A market order is an instruction to buy/sell a specific amount of an asset without specifying the price. In contrast, a limit order specifies the price that should not be exceeded in the case of a buy order or gone below in case of a sell order.

A number of exchanges and market data vendors (e.g. [22, 59, 65]) provide limit order book data for commercial and/or academic research use. However, the cost associated with accessing highly granular data is typically a deterrent to many, thereby necessitating the need for high-quality free synthetic market data.

It is also important to note that there is a limit to the granularity of the data provided to market participants and the research community. Level 3 limit order book data providing a view of the full order book and the associated individual orders is now provided by market data vendors. However, Level 4 data, which associate each order to a specific trading entity is typically only available to regulators. This provides an interesting research opportunity for synthesising such data for use for various research purposes.

One technique to generate such data would be through the use of multi-agent simulators and the construction of agent-based models composed of various trading agents with different strategies, objectives and timescales [15]. This would allow for the generation of an agent-specific data describing the behaviour, for example, of a market maker. However, the limitation of such approach is in the complexity and realism of the strategy designed in the model. Another interesting approach would be to utilize techniques of imitation learning (e.g. behavioural cloning [11]) to synthesize this data utilising historical observations and actions taken by an expert trading agent.

The main technical challenge in synthesising such order book data is that of representing aggregate decisions of many independent actors with differing risk tolerance, rationality and motives. In addition, generating realistic datasets requires defining an appropriate distance measure between datasets, then accurately measuring how close the two datasets are. The empirical properties of limit order book data have been studied extensively in the literature and are often referred to as *stylized facts* of the real limit order book

data. It is important to make sure that the empirical properties of the synthetic data follows, as close as possible, those of real order book data. For example, it is empirically shown that lower spreads (the difference between the best bid and ask prices) are observed during period of high trading volumes and that trading volumes are typically highest at the beginning and end of the trading day. These are two examples, of many, which would need to be taken into consideration when synthesizing synthetic data in the market microstructure domain [2].

2.2.2 Privacy model. With market microstructure data it is not immediately clear what the correct notion of privacy is. Exchanges as a rule do not publish the source of individual orders in the limit order book. Nonetheless, there are well-documented trading strategies [5, 12] that effectively ascertain the intentions of individual market participants and profit from this knowledge at the expense of specific actors. In real markets, institutional traders wish to be protected from such strategies. Work such as that by Asharov et al. [10] aim to protect investors using cryptographic techniques by operating a limit order book purely using encrypted orders sent in by clients. This work is invaluable for traders operating in real markets but the techniques are unsuitable for synthetic data generation, seeing as the orders are encrypted whereas in the synthetic data setting we require the data to resemble the source it derives from.

3 TECHNIQUES FOR SYNTHETIC DATA GENERATION

We shall focus mainly on synthetic data generation with privacy guarantees. The following sections highlight a selected list of techniques used for financial data generation.

3.1 Tabular Data

A number of techniques have been proposed for tabular data generation. For a comprehensive survey of these methods see the survey by Surendra and Mohan [76].

In the Data Mining literature Eno and Thompson [33] define an XML-based synthetic data definition language (SDDL) from which synthetic data may be generated. The algorithm generating the SDDL effectively “inverts” a decision-tree classifier. This method carries no guarantees on privacy. In a similar fashion, there are methods based on other classical machine learning classifiers, such as support vector machines [27] and Random Forests [17]. The drawback with these methods is that the better the classifier, the higher the risk of leaking data, with no tunable privacy parameter.

Abowd and Vilhuber [3] provide a Bayesian inspired differentially private synthetic data release of multidimensional tabular data. Zhang et al. [84] present a synthetic data generator based on Bayesian networks that is also differentially private. This method is effective, with tunable privacy parameters but suffers from the drawback of growing substantially with each new feature added. The method devised by Li et al. [54] generates data according to the histogram of each feature, linking the features via a copula. This method enjoys privacy guarantees but scales poorly with a growing number of features. The Gibbs sampling based method by Park and Ghosh [62] gives strong privacy guarantees and scales well. This technique is however limited to categorical variables.

These techniques for generating synthetic tabular data with privacy protection all suffer to a varying degree from the following limitations. Most differential privacy frameworks represent a row of a given table as a bit string with length equal to the *domain size*, which grows exponentially in the number of columns of the table. This representation quickly becomes impractical to use. The second limitation of this representation is that most high-dimensional datasets are very sparse, resulting in the noise being added to generate the privacy completely washing out the real data, rendering the released dataset unsuitable as an approximation to the true dataset. A more thorough discussion of these limitations can be found in Zhang [83].

Agent-based modelling (ABM) has been used in the context of synthesizing payments data, for instance in modelling a bank's payment processing system [38] and investigating the macroscopic impact of a disruptive event on the flow of interbank payments [9]. Synthetic data for a retail shoe store has been created using ABM by Lopez-Rojas and Axelsson [57]. This generated data intrinsically respects privacy constraints if calibration is manually carried out. If an automated quantitative method is used there is a risk of data leakage. Calibration methods in the ABM literature currently do not explicitly preserve privacy, a survey on ABM calibration methods can be found in [64].

Several deep generative models based on GANs, conditional GANs and variational autoencoders (VAEs) have been proposed for tabular data, such as MedGAN [21], TableGAN [61], PATE-GAN [45], CT-GAN and T-VAE [82]. CT-GAN, the current state of the art approach for mixed tabular data with discrete and continuous columns, overcomes the limitations of previous approaches by using a conditional GAN architecture and encoding scheme, for both discrete variables and continuous variables with multi-modal distributions, to improve the diversity of generated samples, particularly with there is a large class imbalance. CT-GAN, however, does not offer any privacy guarantees. PATE-GAN modifies the typical GAN training procedure to generate tabular data that is differentially private.

3.2 Synthetic financial time series

There has been a lot of work in generating synthetic financial time series data and in releasing differentially private data streams, but little in the way of synthetic financial time series data with privacy guarantees. The classical approach is to propose a simple statistical model such as autoregressive or GARCH (Generalized Autoregressive Conditional Heteroscedasticity) models for financial time series, and then fit real-world data to this model using maximum-likelihood. These methods carry the advantage of being easy to fit and interpret, but rely on strong assumptions and are unable to reproduce many of the statistical features of many financial time series [80]. Examples of more modern approaches using neural architectures are QuantGAN [81] and the work by Fu et al. [37] for modelling log returns of stocks and associated classical time series models. These methods provide no privacy guarantees and have yet to be conclusively shown not to be explicitly memorising data. Agent-based models are often used to replicate the dynamics of financial markets and then used to derive financial time series approximating those seen live. This topic lies outside the scope of

this paper, see [15], [52], [40, Section 5C] and references therein for more on agent based models in finance.

3.3 Stream data with privacy guarantees

Publishing stream data with privacy guarantees was first addressed by Dwork [29] and Dwork et al. [31]. In this approach, stream data is viewed as a bit string with a continuous counter of the number of ones observed being reported. This model is called *event-level* privacy as individual events in the stream are protected. One can also protect *user-level* privacy, where the presence of a particular user in a dataset is masked from an adversary. This is not necessarily the most natural model of privacy for time series data. A privacy model called *w-event* privacy, protecting events occurring within any window of w timestamps, was proposed by Kellaris et al. [48] with associated data release algorithms. The model used here is one in which an organisation wishes to publish continuous counts of a finite number of different events occurring. There is also the notion of *d*-privacy, where d is a metric defined on datasets that generalises adjacency of databases for differential privacy, defined by Chatzikokolakis et al. [18]. This model allows for preserving the privacy of continuous quantities. These two models were combined by Fioretto and Van Hentenryck [35] into the notion of (w, α) -indistinguishability, protecting blocks of continuously-valued data, which is more relevant to financial data. The authors present their OPTSTREAM algorithm for releasing time series data under this model of privacy. Limitations of these methods have not yet been settled upon by the research community, but representation of this data is likely to lead to same scaling issues as in tabular data as the number of streams being simultaneously released increases. There is a yet no established approach in synthesising financial time series data, in particular market microstructure data.

3.4 Unstructured Data

Neural network based methods such as those presented by Shokri and Shmatikov [74], Acs et al. [4] and [1] extend differentially private synthetic data generation to the domain of unstructured data such as images and audio. Although highly promising, these methods are currently in their infancy and suffer from the usual problems neural networks face [39, 67]. Moreover, these methods protect data on an individual level, that is, the presence of a singular data point (x, y) within the training set is protected. This notion of privacy is insufficient in the case of, say, a image generator for faces, with the training set containing the same face at different angles. In this example one would need to protect against an attacker ascertaining the presence of this entire group of images in the training set.

Controlling the the trade off between noise and privacy [30] in an optimal way for synthetic data remains an open problem.

3.5 Network/Graph Data

There is now a mature literature on generation of synthetic graphs. Well-established classical models for synthetic graph generation are the Erdős-Rényi model [34], the Barabási-Albert model [7] and Kronecker graphs [53]. These classical techniques offer no privacy guarantees on the synthetic data. Another generation technique for graph data garnering significant attention is maximum-entropy

methods [60]. This family of methods involve constraining various properties of the synthetic graph to be the same as that of a real graph, say, the spectrum of the graph's Laplacian and producing a distribution over graphs satisfying the constraints that maximises entropy. While not providing explicit privacy protection, the maximum entropy property of the generator distribution means that it is minimally informative, that is, it only reveals information encoded in the constraints. If the constraints are chosen in such a way that significant information is not revealed, then the generated graphs will not reveal this information either. The reader is referred to [56] for a survey focusing on synthetic graph generation.

Due to the often-sensitive nature of network data, the notion of privacy in this domain is well-developed. Particularly focusing on social networks, there has been work well over a decade now on *anonymization*, see for instance, Zhou et al. [85] or Tripathy et al. [79]. More powerful methods following differential privacy have been developed more recently, with the specialization to graphs well documented in [68]. A powerful recent method by Eliáš et al. [32] achieves the following: a polynomial-time differentially private method to generate synthetic graphs that have bounded *cut-distance* to the base graph. Cut-distance corresponds roughly to the maximum difference in weight between a chosen subset of weighted edges between two graphs.

4 EVALUATING GENERATIVE MODELS

As argued by Donoho [26], a large part of the success of machine learning as a field comes from the *common task framework*, namely having a set of benchmark datasets and a real-valued metric that was used to compare model quality. With this in mind, there has been a large body of work attempting to evaluate generative models using a single real-valued metric, as well as a standardisation of datasets being used.

In particular, we briefly focus on the success in the computer vision community with regard to generated image data. While of a very different nature to data that is of financial origin, image data shares many salient properties. Both types of data can often be qualitatively evaluated well by humans to some degree, and come from high-dimensional multi-modal distributions. Indeed, these properties are what lead classical methods such as density estimation to be poorly performing at generating convincing synthetic data. In the paper by Theis et al. [78], the authors argue against using classical methods such as log-likelihood and Parzen window estimates as an evaluation metric for generative models of high-dimensional image data. They provide explicit examples of where these metrics fail, for instance when using nearest neighbour evaluation in the dataset of generated images, whereby perceptually small changes can lead to large changes in Euclidean distance and vice versa. For datasets with a large number of images, small changes can make the nearest neighbour qualitatively different. There is theoretical justification for this being the rule rather than the exception, see for instance [73]. Metrics such as the Inception Score [71], Fréchet Inception Distance [43] and precision/recall [69] have been devised to combat these problems, but suffer from their own disadvantages. This type of analysis and creation of metrics is what is currently missing in the literature for generating financial data, time series in particular.

Evaluating generative models for tabular data is still a nascent discipline with, to our knowledge, no widely established benchmarks, datasets or metrics. At the time of writing, the most common method is to adapt classical statistical methods such as the Kolmogorov-Smirnov test [20, 58] or more recent nonparametric approaches like the MMD [41] to compare distances between the real and synthetic empirical distributions. Another approach taken in [82] is to train prediction models on the synthetic data and use the real data for test evaluation. However, none of these methods provide a single real-valued metric for evaluating such models.

A number of different techniques are used in evaluating generative models for time series-type data. For example, Wiese et al. [81] and Li et al. [55] estimate statistical distances between the test data and generated data. Indeed, in light of the discussion by Theis et al. [78] one can not rely on these measures as a metric for realism. Alternatively, the evaluation of certain quantities such as Value-at-Risk (VaR) [46] is “backtested”, whereby historical data for a time period is used to train a generative model, at which point the true VaR observed in the subsequent historical period is compared with that obtained by the generative model [37]. This method suffers from a dearth of data since there is only one history. There are alternatively, certain “stylized facts” that are widely observed in market microstructure data, that are often evaluated qualitatively [2, 13]. The stylized facts thus far have yet to be converted into a single real-valued metric in a convincing way.

For graphical synthetic data there are a number of common validation techniques. For model-based synthetic data generators, maximum likelihood is a common technique [53]. There is a large literature on *graph kernels*, a technique used predominantly in classification, where a graph kernel is a function that returns the distance between implicitly defined embeddings into \mathbb{R}^n of two graphs. The kernel function acts as a similarity measure between graphs. For a recent review on the various graph kernels, see Kriege et al. [50]. For maximum-entropy methods that return a distribution over graphs [60], the validation is implicit as the distribution returned is that which is minimally informative under the supplied constraints. In this case it is incumbent on the user to supply constraints that capture the salient features of the real network. For graph data based on social networks in particular, see Sala et al. [70] for a review on validation methods.

5 TOWARDS A COMPACT REPRESENTATION OF REAL DATA

Given the preceding discussion, it is worthwhile to consider a framework for how synthetic data can best be represented and transferred between different parties. Moreover, we can examine how existing techniques fit into this framework. The following properties of such a representation are desirable:

Privacy preserving. It is important to define the groups, entities and events whose privacy is to be protected, with tunable parameters for the trade-off between privacy and realism.

Human readable. It is desirable that the format in which synthetic data, or a generative model describing the synthetic data be readily interpretable without undue difficulty. In

finance this is especially relevant as regulators, internal colleagues and other parties must have confidence in the representation and its privacy properties. An overly technical representation is less likely to be relevant and foster trust.

Compact. The representation of synthetic data should be such that it is significantly smaller than the real data it is derived from, and should be reconstructible on the receiver's end, ideally with open source software. The synthetic data should also require little technical know-how to generate, so that staff with direct access can synthesise the data on-premise.

With regard to the first point in the framework above, (ϵ, δ) -differential privacy provides a good model for tabular and graphical data, as seen in Section 3.1. For stream data and more structured data there is still not a consensus in the literature as to the best way for privacy to be defined. Indeed, even for tabular data differential privacy is not the last word, with future directions suggested by Kifer and Machanavajjhala [49].

On the second point, Eno and Thompson [33] define a synthetic data language for tabular data, which is a promising direction although lacking any privacy guarantees. The *synthetic data vault* [63] is a more contemporary approach, with a synthetic database specified via JSON [24]. Some aspects of the correlations between variables are captured, although a solution more tailored to financial data driven from more exotic distributions is required.

The most common ways of sharing synthetic data are either providing the generated dataset itself, or sharing learned parameters for the generative process. These two approaches have contrasting readability and compactness properties that need to be balanced on a case by case basis.

6 CONCLUSIONS

In this review we have highlighted the unique challenges faced in generating financial data, and we have summarised the state of related literature at the time of writing.

We describe the common types of financial data that are encountered in the financial services sector, namely tabular and graphical data in retail banking and time series of market microstructure data. The state of the art in privacy-preserving generation and release of tabular and graph data is slightly more advanced than the equivalent task in stream data, perhaps due to the increased complexity of the latter.

Despite the various pitfalls described in our paper, we strongly believe in the importance and utility of researching methods to synthesize financial data as it will aid in building models to tackle issues of fairness and trustworthiness in financial market operations.

ACKNOWLEDGMENTS

This paper was prepared for informational purposes by the Artificial Intelligence Research group of JPMorgan Chase & Co and its affiliates ("J.P. Morgan"), and is not a product of the Research Department of J.P. Morgan. J.P. Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or

service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

REFERENCES

- [1] Nazmiye Ceren Abay, Yan Zhou, Murat Kantarcioglu, Bhavani Thuraisingham, and Latanya Sweeney. 2019. Privacy Preserving Synthetic Data Release Using Deep Learning. In *Machine Learning and Knowledge Discovery in Databases*, Michele Berlingerio, Francesco Bonchi, Thomas Gärtner, Neil Hurley, and Georgiana Ifrim (Eds.). Springer International Publishing, Cham, 510–526.
- [2] Frédéric Abergel, Marouane Anane, Anirban Chakraborti, Aymen Jedidi, and Ioane Muni Toke. 2016. *Limit Order Books*. Cambridge University Press.
- [3] John M. Abowd and Lars Vilhuber. 2008. How Protective Are Synthetic Data?. In *Privacy in Statistical Databases*, Josep Domingo-Ferrer and Yücel Saygın (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 239–246.
- [4] G. Acs, L. Melis, C. Castelluccia, and E. De Cristofaro. 2019. Differentially Private Mixture of Generative Neural Networks. *IEEE Transactions on Knowledge and Data Engineering* 31, 6 (2019), 1109–1121.
- [5] Jacob Adrian. 2016. Informational Inequality: How High Frequency Traders Use Premier Access to Information to Prey on Institutional Investors. *Duke Law Technology Review* 14 (2016), 256–279.
- [6] Charu C. Aggarwal. 2005. On K -anonymity and the Curse of Dimensionality. In *Proceedings of the 31st International Conference on Very Large Data Bases (Trondheim, Norway) (VLDB '05)*. VLDB Endowment, 901–909.
- [7] Réka Albert and Albert-László Barabási. 2002. Statistical mechanics of complex networks. *Rev. Mod. Phys.* 74 (Jan 2002), 47–97. Issue 1.
- [8] Maryam Archie, Sophie Gershon, Abigail Katcoff, and Aaron Zeng. 2018. *Who's Watching? De-anonymization of Netflix Reviews using Amazon Reviews*. Technical Report. MIT.
- [9] Luca Arciero, Claudia Biancotti, Leandro D'Aurizio, and Claudio Impenna. 2009. Exploring Agent-Based Methods for the Analysis of Payment Systems: A Crisis Model for StarLogo TNG. *Journal of Artificial Societies and Social Simulation* 12, 1 (2009), 2.
- [10] Gilad Asharov, Tucker Hybinette Balch, Antigoni Polychroniadou, and Manuela Veloso. 2020. Privacy-Preserving Dark Pools. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (Auckland, New Zealand) (AAMAS '20)*. International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1747–1749.
- [11] Michael Bain and Claude Sammut. 1995. A Framework for Behavioural Cloning. In *Machine Intelligence* 15. 103–129.
- [12] Andrea Barbon, Marco Di Maggio, Francesco Franzoni, and Augustin Landier. 2019. Brokers and Order Flow Leakage: Evidence from Fire Sales. *The Journal of Finance* 74, 6 (2019), 2707–2749.
- [13] Jean-Philippe Bouchaud, Julius Bonart, Jonathan Donier, and Martin Gould. 2018. *Trades, quotes and prices: financial markets under the microscope*. Cambridge University Press, Cambridge.
- [14] Claire McKay Bowen and Fang Liu. 2019. Comparative Study of Differentially Private Data Synthesis Methods. *arXiv:1602.01063 [stat.ME]* (2019).
- [15] David Byrd, Maria Hybinette, and Tucker Hybinette Balch. 2019. ABIDES: Towards High-Fidelity Market Simulation for AI Research. *arXiv:1904.12066 [cs.MA]* (2019).
- [16] Fabio Caccioli, Paolo Barucca, and Teruyoshi Kobayashi. 2017. Network models of financial systemic risk: A review. *Journal of Computational Social Science* (10 2017).
- [17] Gregory Caiola and Jerome P. Reiter. 2010. Random Forests for Generating Partially Synthetic, Categorical Data. *Trans. Data Privacy* 3, 1 (April 2010), 27–42.
- [18] Konstantinos Chatzikokolakis, Miguel E. Andrés, Nicolás Emilio Bordenabe, and Catuscia Palamidessi. 2013. Broadening the Scope of Differential Privacy Using Metrics. In *Privacy Enhancing Technologies*, Emiliano De Cristofaro and Matthew Wright (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 82–102.
- [19] Long Cheng, Fang Liu, and Danfeng (Daphne) Yao. 2017. Enterprise data breach: causes, challenges, prevention, and future directions. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 7, 5 (2017), e1211.
- [20] Edward Choi, Siddharth Biswal, Bradley Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun. 2017. Generating Multi-label Discrete Patient Records using Generative Adversarial Networks. In *Proceedings of the 2nd Machine Learning for Healthcare Conference (Proceedings of Machine Learning Research)*, Finale Doshi-Velez, Jim Fackler, David Kale, Rajesh Ranganath, Byron Wallace, and Jenna Wiens (Eds.), Vol. 68. PMLR, Boston, Massachusetts, 286–305.
- [21] Edward Choi, Siddharth Biswal, Bradley A. Malin, Jon Duke, Walter F. Stewart, and Jimeng Sun. 2017. Generating Multi-label Discrete Patient Records using Generative Adversarial Networks. In *Proceedings of the Machine Learning for Health Care Conference, MLHC 2017, Boston, Massachusetts, USA, 18-19 August 2017 (Proceedings of Machine Learning Research)*, Finale Doshi-Velez, Jim Fackler,

- David C. Kale, Rajesh Ranganath, Byron C. Wallace, and Jenna Wiens (Eds.), Vol. 68. PMLR, 286–305.
- [22] IEX cloud. 2020. <https://iexcloud.io/>
- [23] I. Glenn Cohen and Michelle M. Mello. 2018. HIPAA and Protecting Health Information in the 21st Century. *JAMA* 320, 3 (2018), 231–232.
- [24] Douglas Crockford. 2001. JSON. <https://www.json.org>
- [25] Jon P. Daries, Justin Reich, Jim Waldo, Elise M. Young, Jonathan Whittinghill, Andrew Dean Ho, Daniel Thomas Seaton, and Isaac Chuang. 2014. Privacy, Anonymity, and Big Data in the Social Sciences. *Commun. ACM* 57, 9 (Sept. 2014), 56–63.
- [26] David Donoho. 2017. 50 Years of Data Science. *Journal of Computational and Graphical Statistics* 26, 4 (2017), 745–766.
- [27] Jörg Drechsler. 2010. Using Support Vector Machines for Generating Synthetic Datasets. In *Proceedings of the 2010 International Conference on Privacy in Statistical Databases* (Corfu, Greece) (PSD'10). Springer-Verlag, Berlin, Heidelberg, 148–161.
- [28] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- [29] Cynthia Dwork. 2010. Differential Privacy in New Settings. In *Proceedings of the Twenty-first Annual ACM-SIAM Symposium on Discrete Algorithms* (Austin, Texas) (SODA '10). Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 174–183.
- [30] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. 2006. Calibrating Noise to Sensitivity in Private Data Analysis. In *Theory of Cryptography*, Shai Halevi and Tal Rabin (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 265–284.
- [31] Cynthia Dwork, Moni Naor, Toniann Pitassi, and Guy N. Rothblum. 2010. Differential Privacy Under Continual Observation. In *Proceedings of the Forty-second ACM Symposium on Theory of Computing* (Cambridge, Massachusetts, USA) (STOC '10). ACM, New York, NY, USA, 715–724.
- [32] Marek Eliáš, Michael Kapralov, Janardhan Kulkarni, and Yin Tat Lee. 2020. Differentially Private Release of Synthetic Graphs. In *Proceedings of the 2020 ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 560–578.
- [33] J. Eno and C. W. Thompson. 2008. Generating Synthetic Data to Match Data Mining Patterns. *IEEE Internet Computing* 12, 3 (May 2008), 78–82.
- [34] P. Erdős and A. Rényi. 1959. On Random Graphs I. *Publicationes Mathematicae Debrecen* 6 (1959), 290.
- [35] Ferdinando Fioretto and Pascal Van Hentenryck. 2019. Optstream: Releasing Time Series Privately. *J. Artif. Int. Res.* 65, 1 (May 2019), 423–456.
- [36] Andrea Fronzetti Colladon and Elisa Remondi. 2017. Using social network analysis to prevent money laundering. *Expert Systems with Applications* 67 (2017), 49–58.
- [37] Rao Fu, Jie Chen, Shutian Zeng, Yiping Zhuang, and Agus Sudjianto. 2019. Time Series Simulation by Conditional Generative Adversarial Net. *arXiv:1904.11419 [stat.ML]* (2019).
- [38] M. Galbiati and K. Soramäki. 2011. An agent-based model of payment systems. *Journal of Economic Dynamics and Control* 35, 6 (2011), 859–875.
- [39] L. H. Gilpin, D. Bau, Yuan B. Z., A. Bajwa, M. Specter, and L. Kagal. 2018. Explaining Explanations: An Overview of Interpretability of Machine Learning. *The 5th IEEE International Conference on Data Science and Advanced Analytics (DSAA 2018)*. (2018).
- [40] M. D. Gould, M. A. Porter, S. Williams, M. McDonald, D. J. Fenn, and S. D. Howison. 2013. Limit order books. *Quantitative Finance* 13, 11 (2013), 1709–1742.
- [41] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. A Kernel Two-sample Test. *Journal of Machine Learning Research* 13, 1 (March 2012), 723–773.
- [42] Daniel Grigat and Fabio Caccioli. 2017. Reverse stress testing interbank networks. *Scientific Reports* 7, 1 (15 Nov 2017), 15616.
- [43] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. 2017. GANs Trained by a Two Time-scale Update Rule Converge to a Local Nash Equilibrium. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (Long Beach, California, USA) (NIPS'17). Curran Associates Inc., USA, 6629–6640.
- [44] Chris Jay Hoofnagle, Bart van der Sloot, and Frederik Zuiderveen Borgesius. 2019. The European Union general data protection regulation: what it is and what it means. *Information & Communications Technology Law* 28, 1 (2019), 65–98.
- [45] James Jordon, Jinsung Yoon, and Mihaela van der Schaar. 2019. PATE-GAN: Generating Synthetic Data with Differential Privacy Guarantees. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6–9, 2019*. OpenReview.net. <https://openreview.net/forum?id=S1zk9iRqF7>
- [46] Philippe Jorion. 2006. *Value at Risk: The New Benchmark for Managing Financial Risk* (3rd ed.). McGraw-Hill.
- [47] Niaz Kammoun, Ahmed Bounfour, Altay Özaygen, and Rokhaya Dieye. 2019. Financial market reaction to cyberattacks. *Cogent Economics & Finance* 7, 1 (2019), 1645584.
- [48] Georgios Kellaris, Stavros Papadopoulos, Xiaokui Xiao, and Dimitris Papadias. 2014. Differentially Private Event Sequences over Infinite Streams. *Proc. VLDB Endow.* 7, 12 (Aug. 2014), 1155–1166.
- [49] Daniel Kifer and Ashwin Machanavajjhala. 2012. A Rigorous and Customizable Framework for Privacy. In *Proceedings of the 31st ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems* (Scottsdale, Arizona, USA) (PODS '12). ACM, New York, NY, USA, 77–88.
- [50] Nils M. Kriege, Fredrik D. Johansson, and Christopher Morris. 2020. A survey on graph kernels. *Applied Network Science* 5, 1 (2020), 6.
- [51] Logan Kugler. 2019. Protecting the 2020 Census. *Commun. ACM* 62, 7 (June 2019), 17–19.
- [52] Blake LeBaron. 2006. Agent-based Computational Finance. In *Handbook of Computational Economics*, Leigh Tesfatsion and Kenneth L. Judd (Eds.). Handbook of Computational Economics, Vol. 2. Elsevier, Chapter 24, 1187–1233.
- [53] Jure Leskovec and Christos Faloutsos. 2007. Scalable Modeling of Real Graphs Using Kronecker Multiplication. In *Proceedings of the 24th International Conference on Machine Learning* (Corvallis, Oregon, USA) (ICML '07). Association for Computing Machinery, New York, NY, USA, 497–504.
- [54] Haoran Li, Li Xiong, Lifan Zhang, and Xiaoqian Jiang. 2014. DPSynthesizer: Differentially Private Data Synthesizer for Privacy Preserving Data Sharing. *Proc. VLDB Endow.* 7, 13 (Aug. 2014), 1677–1680.
- [55] Junyi Li, Xintong Wang, Yaoyang Lin, Arunesh Sinha, and Michael P. Wellman. 2019. Generating Realistic Stock Market Order Streams. <https://openreview.net/forum?id=rke41hC5Km>
- [56] Seung-Hwan Lim, Sangkeun Lee, Sarah S Powers, Mallikarjun Shankar, and Neena Imam. 2016. Survey of Approaches to Generate Realistic Synthetic Graphs. *Tech. Rep. Oak Ridge National Laboratory* (2016). Issue ORNL/TM-2016/3.
- [57] Edgar Alonso Lopez-Rojas and Stefan Axelsson. 2015. Using the RetSim Fraud Simulation Tool to Set Thresholds for Triage of Retail Fraud. In *Secure IT Systems*, Sonja Buchegger and Mads Dam (Eds.). Springer International Publishing, Cham, 156–171.
- [58] Alejandro Mottini, Alix Lheritier, and Rodrigo Acuna-Agost. 2019. Airline Passenger Name Record Generation using Generative Adversarial Networks. *Presented at the 2018 ICML Workshop on Theoretical Foundations and Applications of Deep Generative Models*. *arXiv:1807.06657 [cs.LG]* (2019).
- [59] NASDAQ TotalView. 2020. <https://www.nasdaq.com/solutions/nasdaq-totalview>
- [60] Juyong Park and M. E. J. Newman. 2004. Statistical mechanics of networks. *Phys. Rev. E* 70 (Dec 2004), 066117. Issue 6.
- [61] Noseong Park, Mahmoud Mohammadi, Kshitij Gorde, Sushil Jajodia, Hongkyu Park, and Youngmin Kim. 2018. Data Synthesis based on Generative Adversarial Networks. *Proc. VLDB Endow.* 11, 10 (2018), 1071–1083.
- [62] Yubin Park and Joydeep Ghosh. 2014. PeGS: Perturbed Gibbs Samplers That Generate Privacy-Compliant Synthetic Data. *Trans. Data Privacy* 7, 3 (Dec. 2014), 253–282.
- [63] N. Patki, R. Wedge, and K. Veeramachaneni. 2016. The Synthetic Data Vault. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, 399–410.
- [64] D. Platt. 2019. A Comparison of Economic Agent-Based Model Calibration Methods. *arXiv:1902.05938 [q-fin.CP]* (2019).
- [65] Polygon Financial Data Platform. 2020. <https://polygon.io/>
- [66] Tahereh Pourhabibi, Kok-Leong Ong, Boo H. Kam, and Yee Ling Boo. 2020. Fraud detection: A systematic literature review of graph-based anomaly detection approaches. *Decision Support Systems* 133 (2020), 113303.
- [67] A. Rahimi and B. Recht. 2017. Reflections on Random Kitchen Sinks. (2017). <http://www.argmin.net/2017/12/05/kitchen-sinks/>
- [68] Sofya Raskhodnikova and Adam Smith. 2014. *Private Analysis of Graph Data*. Springer Berlin Heidelberg, Berlin, Heidelberg, 1–6.
- [69] Mehdi S. M. Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. 2018. Assessing Generative Models via Precision and Recall. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems* (Montréal, Canada) (NIPS'18). Curran Associates Inc., USA, 5234–5243.
- [70] Alessandra Sala, Lili Cao, Christo Wilson, Robert Zablit, Haitao Zheng, and Ben Y. Zhao. 2010. Measurement-Calibrated Graph Models for Social Network Experiments. In *Proceedings of the 19th International Conference on World Wide Web* (Raleigh, North Carolina, USA) (WWW '10). Association for Computing Machinery, New York, NY, USA, 861–870.
- [71] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. 2016. Improved Techniques for Training GANs. In *Proceedings of the 30th International Conference on Neural Information Processing Systems* (Barcelona, Spain) (NIPS'16). Curran Associates Inc., USA, 2234–2242.
- [72] Pierangela Samarati and Latanya Sweeney. 1998. *Protecting Privacy when Disclosing Information: k-Anonymity and Its Enforcement through Generalization and Suppression*. Technical Report.
- [73] Ali Shafahi, W. Ronny Huang, Christoph Studer, Soheil Feizi, and Tom Goldstein. 2019. Are adversarial examples inevitable?. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=r1IWUoA9FQ>
- [74] R. Shokri and V. Shmatikov. 2015. Privacy-preserving deep learning. In *2015 53rd Annual Allerton Conference on Communication, Control, and Computing* (Allerton), 909–910.

- [75] Joan Sieber. 2007. Family Educational Rights and Privacy Act (FERPA). *Journal of Empirical Research on Human Research Ethics* 2, 1 (2007), 101–101.
- [76] H. Surendra and H. S. Mohan. 2015. A Review Of Synthetic Data Generation Methods For Privacy Preserving Data Publishing. *International Journal of Scientific & Technology Research* 4, 8 (2015), 95–101.
- [77] Brian Tarran. 2018. What can we learn from the Facebook–Cambridge Analytica scandal? *Significance* 15, 3 (2018), 4–5.
- [78] L. Theis, A. van den Oord, and M. Bethge. 2016. A note on the evaluation of generative models. *arXiv:1511.01844 [stat.ML]* (2016).
- [79] B. K. Tripathy, M. S. Sishodia, Sumeet Jain, and Anirban Mitra. 2014. *Privacy and Anonymization in Social Networks*. Springer International Publishing, Cham, 243–270.
- [80] Ruey Tsay. 2010. *Analysis of Financial Time Series* (3rd ed.). Wiley & Sons.
- [81] Magnus Wiese, Robert Knobloch, Ralf Korn, and Peter Kretschmer. 2019. Quant GANs: Deep Generation of Financial Time Series. *arXiv:1907.06673 [q-fin.MF]* (2019).
- [82] Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. 2019. Modeling Tabular data using Conditional GAN. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8–14 December 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 7333–7343.
- [83] Jun Zhang. 2016. *Algorithms for Synthetic Data Release under Differential Privacy*. Ph.D. Dissertation. Nanyang Technological University.
- [84] Jun Zhang, Graham Cormode, Cecilia M. Procopiuc, Divesh Srivastava, and Xiaokui Xiao. 2014. PrivBayes: Private Data Release via Bayesian Networks. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data (Snowbird, Utah, USA) (SIGMOD ’14)*. ACM, New York, NY, USA, 1423–1434.
- [85] Bin Zhou, Jian Pei, and WoShun Luk. 2008. A Brief Survey on Anonymization Techniques for Privacy Preserving Publishing of Social Network Data. *SIGKDD Explor. Newsl.* 10, 2 (Dec. 2008), 12–22.