

Synthetic data generation: A literature review

Joao Fonseca^{1*}, Fernando Bacao¹

¹NOVA Information Management School, Universidade Nova de Lisboa

*Corresponding Author

Postal Address: NOVA Information Management School, Campus de Campolide, 1070-312 Lisboa, Portugal

Telephone: +351 21 382 8610

The generation of synthetic data can be used for anonymization, regularization, oversampling, semi-supervised learning, self-supervised learning and various other tasks. The wide range of applications of these mechanisms motivated the development of new algorithms specialized in generating data for specific types of data and Machine Learning (ML) tasks. As a result, the analysis of the different types of generative models

1 Introduction

Synthetic data is defined as data obtained from a generative process based on properties of real data [1]. The generation of synthetic data is essential for various domains and tasks. For example, synthetic data is used as a form of regularizing neural networks (*i.e.*, data augmentation) [CITATION]. One form of anonymizing datasets is via the production of synthetic observations (*i.e.*, synthetic data generation) [CITATION]. In settings where only a small portion of training data is labeled, some techniques generate artificial data using both labeled and unlabeled data with a modified loss function to train neural networks (*i.e.*, semi-supervised learning) [2]. In imbalanced learning contexts, synthetic data can be used to balance the target classes' frequencies and reinforce the learning of minority classes (*i.e.*, oversampling) [3]. Some active learning frameworks use data generation to improve the quality of data selection and classifier training [4]. Other techniques employ data generation to produce deep neural networks without labeled data (*i.e.*, self-supervised learning) [5].

The breadth of these techniques span multiple domains, such as facial recognition [6], Land Use/Land Cover mapping [CITATION], medical image processing [CITATION], Natural Language Processing [7] or credit card default prediction [8]. According to the domain and data type, the data generation techniques used may vary significantly. Generally speaking, some data generation mechanisms are specific to some domains, data types or tasks. For example, ... Most, if not all, of these techniques are applied on the input or output space.

However, there are various data generation techniques that are invariant to the task or data types used. These techniques can be either applied in the feature space [9] or in problems using tabular data. On the one hand, data generation in the feature space uses a generative model to learn a manifold,

lower-dimensional abstraction over the input space [10], defined here as the feature space. At this level, any tabular data generation mechanism can be applied and reconstructed into the input space if necessary. On the other hand, synthetic data generation on tabular data can be applied to most problems. Although, the choice of the generation mechanism is still dependant on (1) the importance of the relationships found between the different features, (2) the ML task to be developed and (3) the motivation for the generation of synthetic data. For example, when generating data to address an imbalanced learning problem (*i.e.*, oversampling), the relationships between the different features are not necessarily kept since the goal is to reinforce the learning of the minority class by redefining an ML classifier’s decision boundaries. If the goal is to anonymize a dataset, perform some type of descriptive task, or ensure a consistent model interpretability, these relationships need to be kept.

Depending on the context, evaluating the quality of the generated data is a complex task. For example, for image and time series data, perceptually small changes in the original data can lead to large changes in the euclidean distance [1, 11]. The evaluation of generative models typically account primarily for the performance in a specific task, since good performance in one criterion does not imply good performance on another [11]. However, in computationally intensive tasks it is often impracticable to search for the optimal configurations of generative models. To address this limitation, other evaluation methods have been proposed to assist in this evaluation, which can be distinguished into statistical divergence metrics and precision/recall metrics [12]. The relevant performance metrics found in the literature are discussed in Section 6.

1.1 Motivation and Contributions

This literature review focuses on the generation mechanisms and generative models underlying the different techniques where synthetic data is generated. We focus on the ML perspective of synthetic data, as opposed to the practical perspective. From a practical sense, synthetic data is used as a proxy of real data. It is assumed to be inaccessible, essential and a secondary asset for purposes such as education, software development, or systems demonstrations [13].

We focus particularly on tabular and feature space (*i.e.*, embedded inputs) augmentation given its breadth in scope. Related literature reviews are mostly focused on specific algorithmic applications, with little to no emphasis on the core generative process. For this reason, these techniques often appear “sandboxed”, even though there is a significant overlap between them. Assefa et al. [1] provides a general overview of synthetic data generation for time series data anonymization in the finance sector. Hernandez et al. [14] reviews data generation techniques for tabular health records anonymization. Raghunathan [15] reviews synthetic data anonymization techniques that preserve the statistical properties of a dataset.

Nalepa et al. [16] reviews data augmentation techniques for brain-tumor segmentation.

The different taxonomies established are often specific to the technique discussed. However, it is possible to establish a broader taxonomy without giving up on specificity.

With this article, we aim to understand the current research gaps in the different data mining techniques that involve synthetic data generation. We compare the strengths and weaknesses of the models developed within each of these fields. Finally, we identify possible future research directions to address some of the limitations found.

Contributions of this paper are summarized below:

- Bridge different ML concepts using synthetic data generation in its core (Algorithmic applications + Review of the State-of-the-art).
- List the different synthetic data generation/data augmentation taxonomies and characterize all relevant methods accordingly (Data augmentation taxonomy).
- Discuss the ML techniques in which synthetic data generation/data augmentation is used, beyond regularization and consolidate the current data generation mechanisms across the different techniques (Algorithmic Applications).
- Bring to light the key challenges of synthetic data generation and put forward possible research directions in the future.

1.2 Paper Organization

TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO
 TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO
 TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO
 TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO
 TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO
 TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO
 TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO TODO
 TODO TODO TODO TODO TODO TODO TODO

2 Data Generation Taxonomy

Image data augmentation taxonomy [17]

There is a distinction between semantic and traditional image data augmentation [18], also discussed in [19]

Synthetic data generation for medical records taxonomy [14] which is incomplete

Data generation mechanisms can be characterized in 4 properties: Architecture, Application level, Scope and Data space. The overall definition of the proposed taxonomy is shown in Figure 1.

1. Level of application (External or Internal)
2. Scope (Local or Global augmentation)
3. Architectural approach (heuristic, network-based or others)
4. Data space (Input, feature or output). Within feature and output: Domain

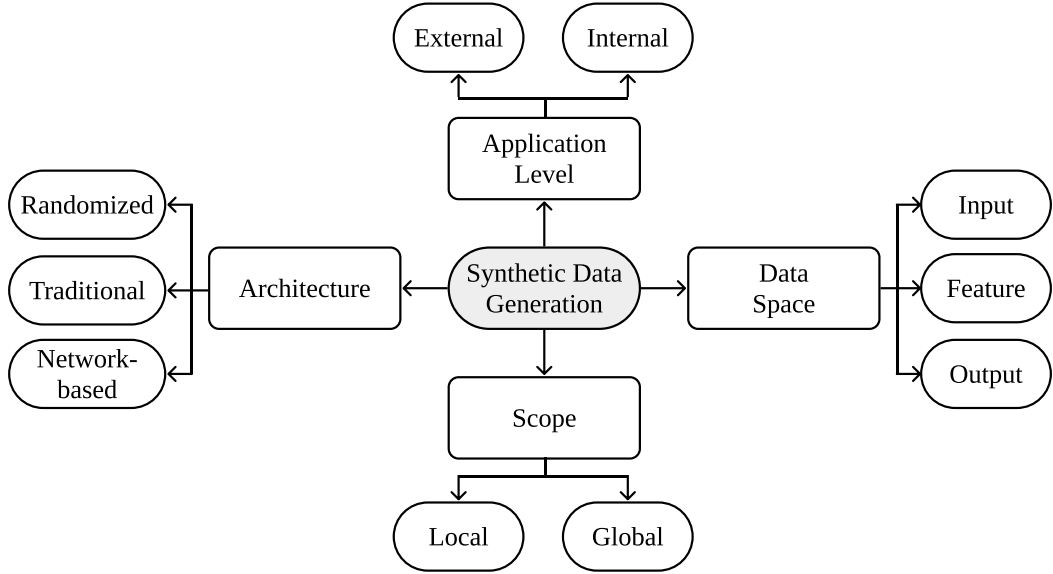


Figure 1: General taxonomy of data generation mechanisms proposed in this paper.

3 Synthetic Data Generation

According to [1]. The generation of synthetic data should aim to fulfil the conditions below:

- Privacy preserving.
- Human readable.
- Compact.

4 Data Generation Mechanisms

In this section, we describe some popular domain and data type-specific data generation techniques. For each data type we include a table with related literature reviews specific to different domains.

4.1 Tabular

4.2 Time series

Generative adversarial networks in time series

4.3 Image

Image-specific data generation mechanisms can be further divided into traditional and semantic techniques [18]. Traditional generation techniques comprise simple modifications such as translation, cropping or random erasing [20]. Semantic generation methods involve more complex tasks, such as changing colors of specific attributes, backgrounds and visual angles [CITATION].

Data generation by modifying specific attributes in data points with known perturbations [6]. For example, overlaying facial elements into a picture containing a human face (*e.g.*, adding sunglasses and different hairstyles), introducing perturbations in facial landmarks, different illumination and artificial misalignment are different approaches to generate artificial observations for facial recognition.

Generative Adversarial Networks in computer vision [21]

4.4 Text

NLP also benefit from data augmentation [7].

In NLP, there is the challenge of establishing universal rules for text transformations to provide new linguistic patterns [22]

<https://github.com/styfeng/DataAug4NLP>

5 Algorithmic applications

5.1 Data Privacy

Synthetic data generation is a technique used to produce synthetic, anonymized versions of datasets [23]. It is considered a good approach to share sensitive data without compromising significantly a given data mining task [24, 25]. Traditional data anonymization techniques, as well as federated learning are two other viable solutions for privacy-preserving data publishing tasks, but contain drawbacks [14]. On the one hand, traditional data anonymization requires domain knowledge, is labor intensive and remains susceptible to disclosure [26]. On the other hand, federated learning is a technically complex task that consists on training ML classifiers on edge devices and aggregating temporarily updated parameters on a centralized server, instead of aggregating the training data [27]. Although it prevents sharing sensitive data, its applicability is dependent on the task. Dataset anonymization via synthetic data generation attempts to balance disclosure risk and data utility in the final synthetic dataset. The goal is to ensure observations are not identifiable and the relevant data mining tasks are not compromised [28, 29].

The generation of synthetic datasets allow a more flexible approach to the successful implementation of ML tasks. However,

Anonymizing data using synthetic data generation in the financial sector [1].

Guidelines for effective synthetic data generation [23]

5.2 Regularization in Supervised Learning

The performance of Machine Learning models is highly dependent on the quality of the training dataset used [30, 31]. The presence of imbalanced and/or small datasets, target labels incorrectly assigned, outliers and high dimensional input spaces reduce the prospects of a successful machine learning (ML) model implementation [31, 32, 33]. In the case of deep learning, for example, these models are often limited by a natural inclination to overfitting, label noise memorization and catastrophic forgetting [34]. Regularization methods are the typical approach to address these problems, but producing robust ML solutions is still a challenge [35].

It is frequently assumed that the training data is sampled from a fixed data source, it is balanced and does not contain label noise. Under these conditions, the resulting ML classifier is expected to achieve good generalization performance [36]. Although, in practical applications, this is rarely the case. When the training data is not representative of the true population, or the model is over-parametrized, it becomes particularly prone to overfitting [37]. Regularization methods attempt to address these limitations. They can be divided into three categories [38]:

1. Output level modifications. Transforms the labels in the training data.
2. Algorithmic level modifications. Modifies the classifier’s architecture, loss function or other components in the training procedure.
3. Input level modifications. Modifies the training dataset by expanding it with synthetic data.

The last approach, input level modifications, is known as data augmentation. Data augmentation is used to increase the size and data variability of data in a training dataset, by producing synthetic observations [39, 40]. Since it is applied at the data level, it can be used for various types of problems and classifiers [41].

5.3 Oversampling

The original author of SMOTE recently published the paper “Efficient Augmentation for Imbalanced Deep Learning” [42]

5.4 Active Learning

5.5 Semi-supervised Learning

5.6 Self-supervised Learning

6 Evaluating the Quality of Synthetic Data

The log-likelihood (and equivalently the Kullback-Leibler Divergence) is a de-facto standard to train and evaluate generative models [11]. Other common metrics include Parzen window estimates, which Theis et al. [11] show that these metrics behave independently and should generally be avoided. Therefore, it is necessary to evaluate generative models with respect to the application these models are being developed for.

The evaluation of generative models should quantify three key aspects of synthetic data [12]:

1. Fidelity
2. Diversity
3. Generalization

The 3-dimensional metric proposed by Alaa et al. [12] quantifies these aspects via the combination of three metrics (α -Precision, β -Recall and Authenticity) for various application domains.

6.1 Statistical Divergence Metrics

6.2 Precision/Recall Metrics

7 Discussion

7.1 Main Findings

7.1.1 RQ1: bla bla bla

7.1.2 RQ2: bla bla bla

7.1.3 RQ3: bla bla bla

7.2 Limitations

Research across the different applications appears to be sandboxed even though all techniques integrate synthetic data in its core.

The evaluation of anonymization techniques lack standardized, objective and reliable performance metrics and benchmark datasets to allow an easier comparison across classifiers to evaluate key aspects of data anonymization (resemblance, utility, privacy and performance). These datasets should contain mixed data types (*i.e.*, a combination of categorical, ordinal, continuous and discrete features) and the metrics should evaluate the performance of different data mining tasks along with the anonymization reliability.

Unlike with data privacy solutions, data augmentation techniques generally do not consider the similarity/dissimilarity of synthetic data.

There is not a clear understanding of what types of data augmentation methods are more appropriate according to different model architectures, ML tasks or domains and the reason why they work better or worse depending on the task.

There is a lack of research on oversampling solutions to generate synthetic data with mixed data types and datasets with exclusively non metric features.

oversampling does not seem to be a relevant source of bias in behavioral research and does not appear to have an appreciably different effect on results for directly versus indirectly oversampled variables [43]

7.3 Research directions

Quantifying the quality of the generated data:

1. Realistic
2. Similarity
3. Usefulness (determine purpose and relevant performance metric)
4. Understand the relationship between the 3 factors

8 Conclusions

References

- [1] Samuel A Assefa, Danial Dervovic, Mahmoud Mahfouz, Robert E Tillman, Prashant Reddy, and Manuela Veloso. “Generating synthetic data in finance: opportunities, challenges and pitfalls”. In: *Proceedings of the First ACM International Conference on AI in Finance*. 2020, pp. 1–8.
- [2] Samuli Laine and Timo Aila. “Temporal ensembling for semi-supervised learning”. In: *International Conference on Learning Representations (ICLR)*. Vol. 4. 5. 2017, p. 6.
- [3] Joao Fonseca, Georgios Douzas, and Fernando Bacao. “Improving imbalanced land cover classification with K-Means SMOTE: Detecting and oversampling distinctive minority spectral signatures”. In: *Information* 12.7 (2021), p. 266.
- [4] Yoon-Yeong Kim, Kyungwoo Song, JoonHo Jang, and Il-Chul Moon. “LADA: Look-Ahead Data Acquisition via Augmentation for Deep Active Learning”. In: *Advances in Neural Information Processing Systems* 34 (2021), pp. 22919–22930.
- [5] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. “Bootstrap your own latent-a new approach to self-supervised learning”. In: *Advances in neural information processing systems* 33 (2020), pp. 21271–21284.
- [6] Jiang-Jing Lv, Xiao-Hu Shao, Jia-Shui Huang, Xiang-Dong Zhou, and Xi Zhou. “Data augmentation for face recognition”. In: *Neurocomputing* 230 (2017), pp. 184–196.

- [7] Steven Y Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. “A Survey of Data Augmentation Approaches for NLP”. In: *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. Online: Association for Computational Linguistics, Aug. 2021, pp. 968–988. DOI: [10.18653/v1/2021.findings-acl.84](https://doi.org/10.18653/v1/2021.findings-acl.84). URL: <https://aclanthology.org/2021.findings-acl.84>.
- [8] Talha Mahboob Alam, Kamran Shaukat, Ibrahim A Hameed, Suhuai Luo, Muhammad Umer Sarwar, Shakir Shabbir, Jiaming Li, and Matloob Khushi. “An investigation of credit card default prediction in the imbalanced datasets”. In: *IEEE Access* 8 (2020), pp. 201173–201198.
- [9] Terrance DeVries and Graham W Taylor. “Dataset augmentation in feature space”. In: *arXiv preprint arXiv:1702.05538* (2017).
- [10] Diederik P Kingma, Max Welling, et al. “An introduction to variational autoencoders”. In: *Foundations and Trends® in Machine Learning* 12.4 (2019), pp. 307–392.
- [11] L Theis, A van den Oord, and M Bethge. “A note on the evaluation of generative models”. In: *International Conference on Learning Representations (ICLR 2016)*. 2016, pp. 1–10.
- [12] Ahmed Alaa, Boris Van Breugel, Evgeny S Saveliev, and Mihaela van der Schaar. “How faithful is your synthetic data? sample-level metrics for evaluating and auditing generative models”. In: *International Conference on Machine Learning*. PMLR. 2022, pp. 290–306.
- [13] Miro Mannino and Azza Abouzied. “Is this real? Generating synthetic data that looks real”. In: *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology*. 2019, pp. 549–561.
- [14] Mikel Hernandez, Gorka Epelde, Ane Alberdi, Rodrigo Cilla, and Debbie Rankin. “Synthetic Data Generation for Tabular Health Records: A Systematic Review”. In: *Neurocomputing* (2022).
- [15] Trivellore E Raghunathan. “Synthetic data”. In: *Annual Review of Statistics and Its Application* 8 (2021), pp. 129–140.
- [16] Jakub Nalepa, Michal Marcinkiewicz, and Michal Kawulok. “Data augmentation for brain-tumor segmentation: a review”. In: *Frontiers in computational neuroscience* 13 (2019), p. 83.
- [17] Nour Eldeen Khalifa, Mohamed Loey, and Seyedali Mirjalili. “A comprehensive survey of recent trends in deep learning for digital images augmentation”. In: *Artificial Intelligence Review* (2021), pp. 1–27.
- [18] Yulin Wang, Gao Huang, Shiji Song, Xuran Pan, Yitong Xia, and Cheng Wu. “Regularizing deep networks with semantic data augmentation”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [19] Connor Shorten and Taghi M Khoshgoftaar. “A survey on image data augmentation for deep learning”. In: *Journal of big data* 6.1 (2019), pp. 1–48.
- [20] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. “Random erasing data augmentation”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 34. 07. 2020, pp. 13001–13008.
- [21] Zhengwei Wang, Qi She, and Tomas E Ward. “Generative adversarial networks in computer vision: A survey and taxonomy”. In: *ACM Computing Surveys (CSUR)* 54.2 (2021), pp. 1–38.
- [22] Markus Bayer, Marc-André Kaufhold, Björn Buchhold, Marcel Keller, Jörg Dallmeyer, and Christian Reuter. “Data augmentation in natural language processing: a novel text generation approach for long and short text classifiers”. In: *International Journal of Machine Learning and Cybernetics* (2022), pp. 1–16.
- [23] Fida K Dankar and Mahmoud Ibrahim. “Fake it till you make it: Guidelines for effective synthetic data generation”. In: *Applied Sciences* 11.5 (2021), p. 2158.

- [24] Jennifer Taub, Mark Elliot, Maria Pampaka, and Duncan Smith. “Differential correct attribution probability for synthetic data: an exploration”. In: *International Conference on Privacy in Statistical Databases*. Springer. 2018, pp. 122–137.
- [25] Noseong Park, Mahmoud Mohammadi, Kshitij Gorde, Sushil Jajodia, Hongkyu Park, and Youngmin Kim. “Data Synthesis based on Generative Adversarial Networks”. In: *Proceedings of the VLDB Endowment* 11.10 (2018).
- [26] Jerome P Reiter. “New approaches to data dissemination: A glimpse into the future (?)” In: *Chance* 17.3 (2004), pp. 11–15.
- [27] Bin Yu, Wenjie Mao, Yihan Lv, Chen Zhang, and Yu Xie. “A survey on federated learning in data mining”. In: *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 12.1 (2022), e1443.
- [28] Kalpana Singh and Lynn Batten. “Aggregating privatized medical data for secure querying applications”. In: *Future Generation Computer Systems* 72 (2017), pp. 250–263.
- [29] Ping Li, Tong Li, Heng Ye, Jin Li, Xiaofeng Chen, and Yang Xiang. “Privacy-preserving machine learning with multiple data providers”. In: *Future Generation Computer Systems* 87 (2018), pp. 341–350.
- [30] Giuseppe Fenza, Mariacristina Gallo, Vincenzo Loia, Francesco Orciuoli, and Enrique Herrera-Viedma. “Data set quality in Machine Learning: Consistency measure based on Group Decision Making”. In: *Applied Soft Computing* 106 (2021), p. 107366.
- [31] Alon Halevy, Peter Norvig, and Fernando Pereira. “The unreasonable effectiveness of data”. In: *IEEE Intelligent Systems* 24.2 (2009), pp. 8–12.
- [32] Pedro Domingos. “A few useful things to know about machine learning”. In: *Communications of the ACM* 55.10 (2012), pp. 78–87.
- [33] Shaeke Salman and Xiuwen Liu. “Overfitting mechanism and avoidance in deep neural networks”. In: *arXiv preprint arXiv:1901.06566* (2019).
- [34] Zeke Xie, Fengxiang He, Shaopeng Fu, Issei Sato, Dacheng Tao, and Masashi Sugiyama. “Artificial neural variability for deep learning: On overfitting, noise memorization, and catastrophic forgetting”. In: *Neural computation* 33.8 (2021), pp. 2163–2192.
- [35] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. “Understanding deep learning (still) requires rethinking generalization”. In: *Communications of the ACM* 64.3 (2021), pp. 107–115.
- [36] Martin Benning and Martin Burger. “Modern regularization methods for inverse problems”. In: *Acta Numerica* 27 (2018), pp. 1–111.
- [37] Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. “Deep learning: a statistical viewpoint”. In: *Acta numerica* 30 (2021), pp. 87–201.
- [38] Claudio Filipi Gonçalves dos Santos and João Paulo Papa. “Avoiding Overfitting: A Survey on Regularization Methods for Convolutional Neural Networks”. In: *ACM Computing Surveys (CSUR)* (2022).
- [39] David A Van Dyk and Xiao-Li Meng. “The art of data augmentation”. In: *Journal of Computational and Graphical Statistics* 10.1 (2001), pp. 1–50.
- [40] Sebastien C Wong, Adam Gatt, Victor Stamatescu, and Mark D McDonnell. “Understanding data augmentation for classification: when to warp?” In: *2016 international conference on digital image computing: techniques and applications (DICTA)*. IEEE. 2016, pp. 1–6.
- [41] Sima Behpour, Kris M Kitani, and Brian D Ziebart. “Ada: Adversarial data augmentation for object detection”. In: *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE. 2019, pp. 1243–1252.

- [42] Damien Dablain, Colin Bellinger, Bartosz Krawczyk, and Nitesh Chawla. “Efficient Augmentation for Imbalanced Deep Learning”. In: *arXiv e-prints* (2022), arXiv-2207.
- [43] Katherina K Hauner, Richard E Zinbarg, and William Revelle. “A latent variable model approach to estimating systematic bias in the oversampling method”. In: *Behavior Research Methods* 46.3 (2014), pp. 786–797.