

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 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 (NLP) [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 7.

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. Specifically, we focus on techniques used in studies published since 2019. 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 tasks like education, software development, or systems demonstrations [13].

We focus on data generation techniques in the tabular and feature space (*i.e.*, embedded inputs), given its breadth in scope. Related literature reviews are mostly focused on specific algorithmic or domain 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. There are some related reviews published since 2019. 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. Bayer et al. [17] distinguishes augmentation techniques for text classification into feature and data space, while providing an extensive overview of augmentation methods within this domain. However, the taxonomy proposed and feature space augmentation methods are not necessarily specific to the domain. Shorten et al. [18], Chen et al. [19], Feng et al. [7] and Liu et al. [20] also review data augmentation techniques for text data. Yi et al. [21] review Generative Adversarial Network architectures for medical imaging. Wang et al. [22] reviews face data augmentation techniques. Shorten et al. [23] and Khosla et al. [24] discuss techniques for image data augmentation. Iwana et al. [25] and Wen et al. [26] also review time series data augmentation techniques. Zhao et al. [27] review data augmentation techniques for graph data. The analysis of related literature reviews ¹ is shown in Table 1.

¹Results obtained using Google Scholar, limited to articles published since 2019, using the search

Table 1: Related literature reviews published since 2019.

Reference	Data type	ML problem	Domain	Observations
Assefa et al. [1]	—	Differential privacy	Finance	Analysis of applications, motivation and properties of synthetic data for anonymization.
Hernandez et al. [14]	Tabular	Differential privacy	Healthcare	Focus on GANs.
Raghunathan [15]	Tabular	Differential privacy	Statistics	Focus on general definitions such as differential privacy and statistical disclosure control.
Nalepa et al. [16]	Image	Segmentation	Medicine	Analysis of algorithmic applications on a 2018 brain-tumor segmentation challenge.
Bayer et al. [17]	Text	Classification	—	Distinguish 100 methods into 12 groups.
Shorten et al. [18]	Text	Deep Learning	—	General overview of text data augmentation.
Chen et al. [19]	Text	Few-shot Learning	—	Augmentation techniques for machine learning with limited data
Feng et al. [7]	Text	—	—	Overview of augmentation techniques and applications on NLP tasks.
Liu et al. [20]	Text	—	Various	Analysis of industry use cases of data augmentation in NLP. Emphasis on input level data augmentation.
Yi et al. [21]	Image	—	Medicine	Emphasis on GANs.
Wang et al. [22]	Image	Deep Learning	—	Regularization techniques using facial image data. Emphasis on Deep Learning generative models.
Shorten et al. [23]	Image	Deep Learning	—	Emphasis on data augmentation as a regularization technique.
Khosla et al. [24]	Image	—	—	Broad overview of image data augmentation. Emphasis on traditional approaches.
Iwana et al. [25]	Time series	Classification	—	Defined a taxonomy for time series data augmentation.
Wen et al. [26]	Time series	Various	—	Analysis of data augmentation methods for classification, anomaly detection and forecasting.
Zhao et al. [27]	Graph	Various	—	Graph data augmentation for supervised and self-supervised learning.
Khalifa et al. [28]	Image	—	Various	General overview of image data augmentation and relevant domains of application.

72 The different taxonomies established in the literature follow a similar philosophy, but vary in terminology
 73 and are often specific to the technique discussed. Regardless, it is possible to establish a broader taxonomy
 74 without giving up on specificity. This study provides a joint overview of the different data generation
 75 approaches, domains and ML techniques where data generation is being used, as well as a common
 76 taxonomy across domains. It extends the analyses found in these articles and uses the compiled knowledge
 77 to identify research gaps. We compare the strengths and weaknesses of the models developed within each
 78 of these fields. Finally, we identify possible future research directions to address some of the limitations
 79 found. The contributions of this paper are summarized below:

query ("synthetic data generation" OR "oversampling" OR "imbalanced learning" OR "data augmentation") AND ("literature review" OR "survey"). Retrieved on August 11th, 2022. More articles were added later whenever found relevant.

- Bridge different ML concepts using synthetic data generation in its core (Algorithmic applications + Review of the State-of-the-art).
- Propose a synthetic data generation/data augmentation taxonomy to resolve the ambiguity in the literature (Data augmentation taxonomy).
- Characterize all relevant data generation methods using the proposed 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

This paper is organized as follows: Section 2 defines and formalizes the different concepts, goals, trade-offs and motivations related to synthetic data generation. Section 3 establishes the taxonomy used to categorize all the methods described in the paper. Section 4 reviews synthetic data generation mechanisms in the feature space. Section 5 reviews synthetic data generation mechanisms in the input space. Section 6 describes the applications of synthetic data in ML methods. Section 7 reviews performance evaluation methods of synthetic data generation mechanisms. Section 8 summarizes the main findings and discusses limitations and possible research directions in the state-of-the-art. Section 9 presents the main conclusions drawn from this study.

2 Background

In this section we define basics concepts, common goals, trade-offs and motivations regarding the generation of synthetic data in ML. We define synthetic data generation as the production of observations using a generative model (regardless of its nature) that resemble naturally occurring observations within a certain domain. It requires access to either a training dataset, a generative process, or a data stream. However, additional requirements might be imposed depending on the ML task being developed. For example, to generate artificial data for regularization purposes in supervised learning (*i.e.*, data augmentation) the training dataset must be annotated [CITATION]. The generation of synthetic data for anonymization purposes assumes synthetic datasets to be different from the original data, while following the same statistical properties [CITATION]. Domain knowledge may also be necessary to encode specific relationships among features into the generative process.

The breach of sensitive information is an important barrier to the sharing of datasets, especially when it concerns personal information [29]. A common solution for this problem is the generation of synthetic data without identifiable information. Generally speaking, ML tasks that require data with sensitive information are not compromised when using synthetic data. The experiment conducted by Patki et al. [30] using relational datasets showed that in 11 out 15 comparisons ($\approx 73\%$), practitioners performing

116 predictive modelling tasks using fully synthetic datasets performed the same or better than those using
117 the original dataset.

118 A common problem in the training of deep neural networks are their capacity to generalize [31] (*i.e.*, reduce
119 the difference in classification performance between known and unseen observations). Data augmentation
120 is a common method to address this problem. The generation of synthetic observations increases the
121 range of the possible input space used in the training phase, which reduces the performance difference
122 between known and unseen observations. Although other regularization methods exist, data augmentation
123 is a useful method since it does not affect the choice in the architecture of the ML classifier and does not
124 exclude the usage of other regularization methods. In domains such as computer vision and NLP, data
125 augmentation is also used to improve the robustness of models against adversarial attacks [32, 33]. These
126 topics are discussed into higher detail in Section 6.2.

127 In supervised learning, synthetic data generation is often motivated by the need to balance target class
128 distributions (*i.e.*, oversampling). Since most ML classifiers are designed to perform best with balanced
129 datasets, defining an appropriate decision boundary to distinguish rare classes becomes difficult [34].
130 Although there are other approaches to address imbalanced learning, oversampling techniques are generally
131 easier to implement since they do not involve modifications to the classifier. This topic is discussed into
132 higher detail in Section 6.3.

133 In supervised learning projects where labeled data is not readily available, but can be labeled, an Active
134 Learning (AL) method may be used to improve the data labelling process. AL aims to find the most
135 informative observations to be labeled and fed into the classifier.

136 2.1 Problem Formulation

137 3 Data Generation Taxonomy

138

139 Image data augmentation taxonomy [28]

140 There is a distinction between semantic and traditional image data augmentation [35], also discussed
141 in [23]

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

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

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

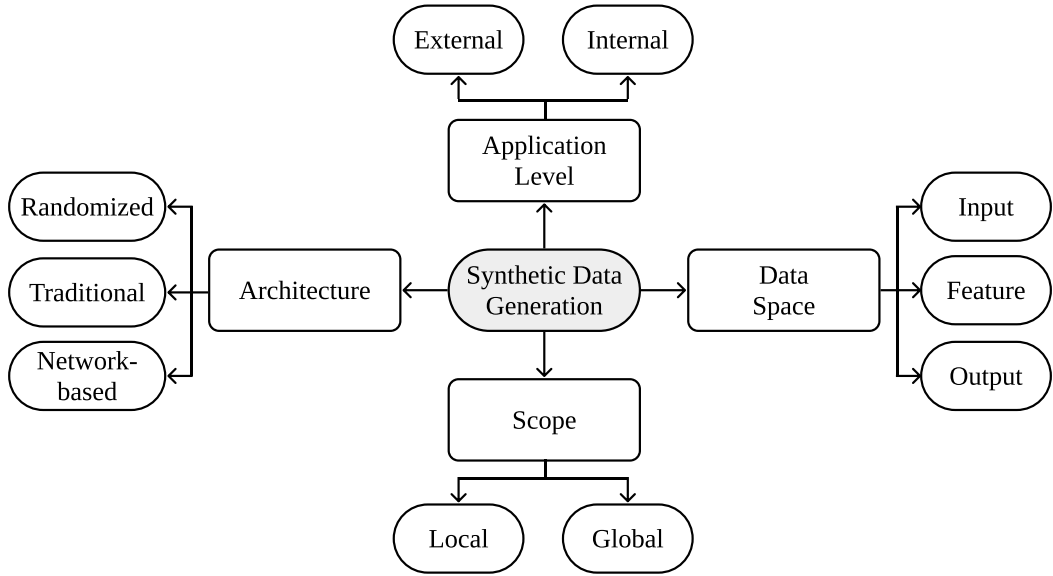


Figure 1: General taxonomy of data generation mechanisms proposed in this paper.

4 Data Generation in the Feature Space

The concept of data generation in the feature space was popularized with [9].

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

- Privacy preserving.
- Human readable.
- Compact.

Discuss Auto-augmentation (as mentioned in [22]) or meta learning (as mentioned in [23])

5 Data Generation in the Input Space

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.

161 5.1 Tabular

162 5.2 Time series

163 Generative adversarial networks in time series

164 5.3 Image

165 Image-specific data generation mechanisms can be further divided into traditional and semantic tech-
166 niques [35]. Traditional generation techniques comprise simple modifications such as translation, cropping
167 or random erasing [36]. Semantic generation methods involve more complex tasks, such changing colors of
168 specific attributes, backgrounds and visual angles [CITATION].

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

173 Generative Adversarial Networks in computer vision [37]

174 5.4 Text

175 NLP motivations [7]:

- 176 1. Low-resource languages (NLP)
- 177 2. Mitigate bias
- 178 3. Fixing class imbalance
- 179 4. Few-shot learning
- 180 5. Adversarial examples

181 NLP also benefit from data augmentation [7].

182 In NLP, there is the challenge of establishing universal rules for text transformations to provide new
183 linguistic patterns [38]

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

185 5.5 Graphs

186 Another relevant paper [39]

187 Various graph data augmentation methods can be applied to related data types such as text data [18].

188 An analysis on different graph data augmentation techniques and a new graph data augmentation
189 framework Zhao et al. [40]

190 List of papers about graph data augmentation: [https://github.com/zhao-tong/graph-data-augmentation-](https://github.com/zhao-tong/graph-data-augmentation-papers)
191 papers

192 5.6 Audio

193 6 Algorithmic applications

194

195 6.1 Data Privacy

196 SynSys [41], Sensegen [42], The Synthetic Data Vault [30]

197 Synthetic data generation is a technique used to produce synthetic, anonymized versions of datasets [29].
198 It is considered a good approach to share sensitive data without compromising significantly a given data
199 mining task [43, 44]. Traditional data anonymization techniques, as well as federated learning are two
200 other viable solutions for privacy-preserving data publishing tasks, but contain drawbacks [14]. On the
201 one hand, traditional data anonymization requires domain knowledge, is labor intensive and remains
202 susceptible to disclosure [45]. On the other hand, federated learning is a technically complex task that
203 consists on training ML classifiers on edge devices and aggregating temporarily updated parameters on a
204 centralized server, instead of aggregating the training data [46]. Although it prevents sharing sensitive
205 data, its applicability is dependent on the task. Dataset anonymization via synthetic data generation
206 attempts to balance disclosure risk and data utility in the final synthetic dataset. The goal is to ensure
207 observations are not identifiable and the relevant data mining tasks are not compromised [47, 48].

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

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

211 Guidelines for effective synthetic data generation [29]

212 6.2 Regularization in Supervised Learning

213

214 The performance of Machine Learning models is highly dependent on the quality of the training dataset
215 used [49, 50]. The presence of imbalanced and/or small datasets, target labels incorrectly assigned, outliers
216 and high dimensional input spaces reduce the prospects of a successful machine learning (ML) model
217 implementation [50, 51, 52]. In the case of deep learning, for example, these models are often limited by a

218 natural inclination to overfitting, label noise memorization and catastrophic forgetting [53]. Regularization
219 methods are the typical approach to address these problems, but producing robust ML solutions is still a
220 challenge [31].

221 It is frequently assumed that the training data is sampled from a fixed data source, it is balanced and does
222 not contain label noise. Under these conditions, the resulting ML classifier is expected to achieve good
223 generalization performance [54]. Although, in practical applications, this is rarely the case. When the
224 training data is not representative of the true population, or the model is over-parametrized, it becomes
225 particularly prone to overfitting [55]. Regularization methods attempt to address these limitations. They
226 can be divided into three categories [56]:

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

231 The last approach, input level modifications, is known as data augmentation. Data augmentation is used to
232 increase the size and data variability of data in a training dataset, by producing synthetic observations [57,
233 58]. Since it is applied at the data level, it can be used for various types of problems and classifiers [59].

234 Problems such as fraud detection and healthcare are frequently tackled via synthetic data generation [60].

235 “Su et al. [78] show that 70.97% of images can be misclassified by changing just one pixel” Shorten et al.
236 [23]

237 “Moreover, the current research about so called adversarial attacks on CNNs showed that deep neural
238 networks can be easily fooled into misclassification of images just by partial rotations and image translation
239 [1], adding the noise to images [5] and even changing one, skillfully selected pixel in the image [6].”
240 Mikołajczyk et al. [61]

241 Data augmentation can also be used to improve a model’s robustness against adversarial attacks.

242 6.3 Oversampling

243

244 KernelADASYN [62]

245 The original author of SMOTE recently published the paper “Efficient Augmentation for Imbalanced Deep
246 Learning” [63]

247 6.4 Active Learning

248 6.5 Few-shot Learning

249

250 6.6 Semi-supervised Learning

251

252 Synthetic data generation for semi-supervised learning given limited labeled data regarding the COVID-19
253 pandemic [64].

254 6.7 Self-supervised Learning

255

256 7 Evaluating the Quality of Synthetic Data

257

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

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

- 264 1. Fidelity
- 265 2. Diversity
- 266 3. Generalization

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

269 7.1 Statistical Divergence Metrics

270 7.2 Precision/Recall Metrics

271 8 Discussion

272

273 8.1 Main Findings

274 The combination of data generation strategies is an approach commonly found in different problems, such
275 as self-supervised learning [5]. It can be more frequently found in text data applications [17] and image
276 data [CITATION].

277 8.1.1 RQ1: bla bla bla

278 8.1.2 RQ2: bla bla bla

279 8.1.3 RQ3: bla bla bla

280 8.2 Limitations

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

283 Given the breadth and complexity of input-level and feature-level data generation mechanisms, it is
284 increasingly important to find a method to efficiently determine the most appropriate data generation
285 policies. However, the complexity of this task is determined by various factors: different data types, ML
286 problems, model architectures, computational resources, performance metrics and contextual constraints.
287 Auto-augmentation and meta learning aim to address this challenge and are still subject to active
288 research.

289 The evaluation of anonymization techniques lack standardized, objective and reliable performance metrics
290 and benchmark datasets to allow an easier comparison across classifiers to evaluate key aspects of data
291 anonymization (resemblance, utility, privacy and performance). These datasets should contain mixed data
292 types (*i.e.*, a combination of categorical, ordinal, continuous and discrete features) and the metrics should
293 evaluate the performance of different data mining tasks along with the anonymization reliability. This
294 problem appears to be universal across domains. For example, Hernandez et al. [14] observed the lack of
295 a universal method or metric to report the performance synthetic data generation algorithms for tabular
296 health records.

297 Computational cost and inconsistent quality of synthetic data generated with GANs (*e.g.*, mode collapse).

298 Unlike with data privacy solutions, data augmentation techniques generally do not consider the similar-
299 ity/dissimilarity of synthetic data.

300 There is not a clear understanding of what types of data augmentation methods are more appropriate
301 according to different model architectures, ML tasks or domains and the reason why they work better or
302 worse depending on the task. In addition, it is still unclear *why* data augmentation works. Research on
303 this topic lacks depth and fails to address the theoretical underpinnings [7].

304 “Dao et al. (2019) note that “data augmentation is typically performed in an ad-hoc manner with little
305 understanding of the underlying theoretical principles”, and claim the typical explanation of DA as
306 regularization to be insufficient.” [7]

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

309 There is no clear understanding of the most appropriate data augmentation techniques used to train
310 self-supervised models and how their behavior and performance varies according to the data generation
311 method used.

312 oversampling does not seem to be a relevant source of bias in behavioral research and does not appear to
313 have an appreciably different effect on results for directly versus indirectly oversampled variables [65]

314 8.3 Research directions

315 Quantifying the quality of the generated data:

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

320 9 Conclusions

321

322 References

- 323 [1] Samuel A Assefa, Danial Dervovic, Mahmoud Mahfouz, Robert E Tillman, Prashant Reddy, and
324 Manuela Veloso. “Generating synthetic data in finance: opportunities, challenges and pitfalls”. In:
325 *Proceedings of the First ACM International Conference on AI in Finance*. 2020, pp. 1–8.
- 326 [2] Samuli Laine and Timo Aila. “Temporal ensembling for semi-supervised learning”. In: *International*
327 *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, Suhui 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] Markus Bayer, Marc-André Kaufhold, and Christian Reuter. “A survey on data augmentation for text classification”. In: *ACM Computing Surveys* (2021).
- [18] Connor Shorten, Taghi M Khoshgoftaar, and Borko Furht. “Text data augmentation for deep learning”. In: *Journal of big Data* 8.1 (2021), pp. 1–34.
- [19] Jiaao Chen, Derek Tam, Colin Raffel, Mohit Bansal, and Diyi Yang. “An empirical survey of data augmentation for limited data learning in NLP”. In: *arXiv preprint arXiv:2106.07499* (2021).

- [20] Pei Liu, Xuemin Wang, Chao Xiang, and Weiye Meng. “A survey of text data augmentation”. In: *2020 International Conference on Computer Communication and Network Security (CCNS)*. IEEE. 2020, pp. 191–195.
- [21] Xin Yi, Ekta Walia, and Paul Babyn. “Generative adversarial network in medical imaging: A review”. In: *Medical image analysis* 58 (2019), p. 101552.
- [22] Xiang Wang, Kai Wang, and Shiguo Lian. “A survey on face data augmentation for the training of deep neural networks”. In: *Neural computing and applications* 32.19 (2020), pp. 15503–15531.
- [23] 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.
- [24] Cherry Khosla and Baljit Singh Saini. “Enhancing performance of deep learning models with different data augmentation techniques: A survey”. In: *2020 International Conference on Intelligent Engineering and Management (ICIEM)*. IEEE. 2020, pp. 79–85.
- [25] Brian Kenji Iwana and Seiichi Uchida. “An empirical survey of data augmentation for time series classification with neural networks”. In: *Plos one* 16.7 (2021), e0254841.
- [26] Qingsong Wen, Liang Sun, Fan Yang, Xiaomin Song, Jingkun Gao, Xue Wang, and Huan Xu. “Time series data augmentation for deep learning: a survey”. In: *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*. Ed. by Zhi-Hua Zhou. International Joint Conferences on Artificial Intelligence Organization, Aug. 2021, pp. 4653–4660.
- [27] Tong Zhao, Gang Liu, Stephan Günnemann, and Meng Jiang. “Graph Data Augmentation for Graph Machine Learning: A Survey”. In: *arXiv preprint arXiv:2202.08871* (2022).
- [28] 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.
- [29] 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.
- [30] Neha Patki, Roy Wedge, and Kalyan Veeramachaneni. “The synthetic data vault”. In: *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*. IEEE. 2016, pp. 399–410.
- [31] 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.
- [32] Yi Zeng, Han Qiu, Gerard Memmi, and Meikang Qiu. “A data augmentation-based defense method against adversarial attacks in neural networks”. In: *International Conference on Algorithms and Architectures for Parallel Processing*. Springer. 2020, pp. 274–289.
- [33] John X Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. “Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp”. In: *arXiv preprint arXiv:2005.05909* (2020).
- [34] José A Sáez, Bartosz Krawczyk, and Michał Woźniak. “Analyzing the oversampling of different classes and types of examples in multi-class imbalanced datasets”. In: *Pattern Recognition* 57 (2016), pp. 164–178.
- [35] 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).
- [36] 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.

- [37] 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.
- [38] 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.
- [39] Jiajun Zhou, Jie Shen, and Qi Xuan. “Data augmentation for graph classification”. In: *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 2020, pp. 2341–2344.
- [40] Tong Zhao, Yozen Liu, Leonardo Neves, Oliver Woodford, Meng Jiang, and Neil Shah. “Data augmentation for graph neural networks”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 35. 12. 2021, pp. 11015–11023.
- [41] Jessamyn Dahmen and Diane Cook. “SynSys: A synthetic data generation system for healthcare applications”. In: *Sensors* 19.5 (2019), p. 1181.
- [42] Moustafa Alzantot, Supriyo Chakraborty, and Mani Srivastava. “Sensegen: A deep learning architecture for synthetic sensor data generation”. In: *2017 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops)*. IEEE. 2017, pp. 188–193.
- [43] 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.
- [44] 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).
- [45] Jerome P Reiter. “New approaches to data dissemination: A glimpse into the future (?)” In: *Chance* 17.3 (2004), pp. 11–15.
- [46] 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.
- [47] Kalpana Singh and Lynn Batten. “Aggregating privatized medical data for secure querying applications”. In: *Future Generation Computer Systems* 72 (2017), pp. 250–263.
- [48] 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.
- [49] 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.
- [50] Alon Halevy, Peter Norvig, and Fernando Pereira. “The unreasonable effectiveness of data”. In: *IEEE Intelligent Systems* 24.2 (2009), pp. 8–12.
- [51] Pedro Domingos. “A few useful things to know about machine learning”. In: *Communications of the ACM* 55.10 (2012), pp. 78–87.
- [52] Shaeke Salman and Xiuwen Liu. “Overfitting mechanism and avoidance in deep neural networks”. In: *arXiv preprint arXiv:1901.06566* (2019).
- [53] 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.

- 463 [54] Martin Benning and Martin Burger. “Modern regularization methods for inverse problems”. In: *Acta*
464 *Numerica* 27 (2018), pp. 1–111.
- 465 [55] Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. “Deep learning: a statistical viewpoint”.
466 In: *Acta numerica* 30 (2021), pp. 87–201.
- 467 [56] Claudio Filipi Gonçalves dos Santos and João Paulo Papa. “Avoiding Overfitting: A Survey on
468 Regularization Methods for Convolutional Neural Networks”. In: *ACM Computing Surveys (CSUR)*
469 (2022).
- 470 [57] David A Van Dyk and Xiao-Li Meng. “The art of data augmentation”. In: *Journal of Computational*
471 *and Graphical Statistics* 10.1 (2001), pp. 1–50.
- 472 [58] Sebastien C Wong, Adam Gatt, Victor Stamatescu, and Mark D McDonnell. “Understanding data
473 augmentation for classification: when to warp?” In: *2016 international conference on digital image*
474 *computing: techniques and applications (DICTA)*. IEEE. 2016, pp. 1–6.
- 475 [59] Sima Behpour, Kris M Kitani, and Brian D Ziebart. “Ada: Adversarial data augmentation for object
476 detection”. In: *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE.
477 2019, pp. 1243–1252.
- 478 [60] Hadi Keivan Ekbatani, Oriol Pujol, and Santi Segui. “Synthetic Data Generation for Deep Learning
479 in Counting Pedestrians.” In: *ICPRAM*. 2017, pp. 318–323.
- 480 [61] Agnieszka Mikołajczyk and Michał Grochowski. “Data augmentation for improving deep learning
481 in image classification problem”. In: *2018 international interdisciplinary PhD workshop (IIPhDW)*.
482 IEEE. 2018, pp. 117–122.
- 483 [62] Bo Tang and Haibo He. “KernelADASYN: Kernel based adaptive synthetic data generation for
484 imbalanced learning”. In: *2015 IEEE congress on evolutionary computation (CEC)*. IEEE. 2015,
485 pp. 664–671.
- 486 [63] Damien Dablain, Colin Bellinger, Bartosz Krawczyk, and Nitesh Chawla. “Efficient Augmentation
487 for Imbalanced Deep Learning”. In: *arXiv e-prints* (2022), arXiv–2207.
- 488 [64] Hari Prasanna Das, Ryan Tran, Japjot Singh, Xiangyu Yue, Geoffrey Tison, Alberto Sangiovanni-
489 Vincentelli, and Costas J Spanos. “Conditional synthetic data generation for robust machine learning
490 applications with limited pandemic data”. In: *Proceedings of the AAAI Conference on Artificial*
491 *Intelligence*. Vol. 36. 11. 2022, pp. 11792–11800.
- 492 [65] Katherina K Hauner, Richard E Zinbarg, and William Revelle. “A latent variable model approach
493 to estimating systematic bias in the oversampling method”. In: *Behavior Research Methods* 46.3
494 (2014), pp. 786–797.