


Supplementary Table 1: Existing categorisations of synthetic data currently used in the literature. A simplified categorisation is proposed in Table 1 of the main text

Data utility	Fully/partially synthetic	Common terms for synthetic data type	Description	Exemplar data case
Minimally disclosive, minimal analytic value, low fidelity	Fully synthetic data <sup>11</sup>	- Dummy data [31]	<ul style="list-style-type: none"> <li>-Preserves the type, structure and format of real data but not the statistical properties of the variables</li> <li>-Do not contain any original data</li> <li>- Almost impossible to identify any single entity</li> </ul>	<ul style="list-style-type: none"> <li>i) Basic code or advanced code testing including data management or cleaning</li> <li>ii) Technical development (API, tool or pipeline testing)</li> <li>iii) Education and training (for data analysis)</li> </ul>
		- Synthetic datasets <sup>12</sup> [32]	<ul style="list-style-type: none"> <li>- Preserves only the structure, format, and data types of the variables</li> <li>- Constructed based only on available metadata; values are generated from ad-hoc distributions and open sources</li> <li>- Contains only values present in the original (univariate) data.</li> <li>-No disclosure risk</li> </ul>	
		• Valid	<ul style="list-style-type: none"> <li>-Preserves the format and record-level plausibility (i.e., values use plausible distributions) and replicates marginal (univariate) distributions where possible</li> <li>-Produced dataset passes the sanity check (validation condition or edit rules) the real dataset would need to go through.</li> <li>- Missing value codes, errors and inconsistencies of the original data are present</li> <li>-Minimal disclosure risk</li> </ul>	
Fully & Partially synthetic data	- Population level synthetic data [87]	- Patient level synthetic data [46]	<ul style="list-style-type: none"> <li>- Key characteristics of variables in the original data (e.g., distributions) are preserved</li> <li>-Complex inter-relationships between variables are not considered</li> <li>-Preserved complex inter-relationships between variables (for each individual)</li> <li>-Close representation of values in real people in a specific population</li> </ul>	<ul style="list-style-type: none"> <li>i) Explore and compare populations</li> <li>ii) Education and training (for data analysis)</li> <li>iii) Testing experimental methods</li> <li>iv) Extended code testing</li> <li>v) Understanding and examining specific groups or populations for study or trial planning</li> <li>vi) Develop analysis plans</li> <li>vii) Produce preliminary results prior to accessing the real data</li> </ul>
	- Synthetically-augmented datasets [32]	• Synthetically-augmented plausible	<ul style="list-style-type: none"> <li>-Preserves the format and record-level plausibility and replicate marginal (univariate) distributions where possible</li> <li>-Constructed based on real dataset, values are generated based on original distributions (with added fuzziness and smoothing)</li> <li>-Does not preserve relationships</li> </ul>	
	• Synthetically-augmented multivariate plausible	<ul style="list-style-type: none"> <li>-Preserves the format and record-level plausibility and replicate multivariate distribution loosely for higher level geographies</li> <li>-Constructed based on real dataset, values are generated based on original distributions (with added fuzziness and smoothing)</li> <li>-Some key relationships are retained</li> </ul>		



Continued

Supplementary Table 1: Continued

Data utility	Fully/partially synthetic	Common terms for synthetic data type	Description	Exemplar data case
		<ul style="list-style-type: none"> <li data-bbox="451 338 687 421">• Synthetically-augmented multivariate detailed</li> <li data-bbox="451 510 667 566">• Synthetically-augmented replica</li> </ul>	<p data-bbox="724 297 1158 409">-Similar to previous but more effort to match the real relationships (joint distributions), e.g., in smaller geographies and household structure</p> <p data-bbox="724 465 1126 548">-Preserves format, structure, joint distributions, missingness patterns, low level geographies.</p> <p data-bbox="724 555 1126 667">-Constructed based on the real dataset, values are generated based on observed joint or conditional distributions, while de-identification methods are applied.</p> <p data-bbox="724 674 1171 763"><b>-In all types of synthetically-augmented datasets, missingness is to be preserved and disclosure control is necessary case by case</b></p>	
		<p data-bbox="411 797 616 853"><b>-Complex modality synthetic data</b></p>	<p data-bbox="724 797 1158 853">-This includes specific modalities such as radiology images, ECG time series data</p>	<p data-bbox="1198 797 1522 853">i) Machine learning for e.g., medicine and healthcare</p> <p data-bbox="1198 860 1522 947">ii) Facilitating the use of data in understanding social and human behaviour</p>

<sup>11</sup>The term “Fully synthetic data” refers to datasets in which all variables are generated, and original values are included. The term “Partially synthetic data” refer to datasets in which only some variables, typically those with sensitive information, are generated while some of the original variables are still present. Also, for the terms assigned in the “Fully and Partially synthetic data” category, datasets can either contain fully synthesized data or partially synthesized data based on the purpose of the research and the statistical modelling applied.

<sup>12</sup>The terms “Synthetic datasets” and “Synthetically-augmented datasets” refer to a high-level scale to evaluate the synthetic data based on how closely they resemble the original data, their purpose and disclosure risk and is proposed by the Office of National Statistics (ONS). More details can be found here: <https://www.ons.gov.uk/methodology/methodologicalpublications/generalmethodology/onsworkingpaperseries/onsmethodologyworkingpaperseriesnumber16syntheticdatapilot>



Supplementary Table 2: Glossary of synthetic data generation techniques

Terms	Definition
<b>Deep Neural Networks</b>	Deep learning is part of the machine learning methods based on artificial neural networks (a technology built to simulate the activity of the human brain) with representation learning. As such, deep neural networks are networks with a certain level of complexity and consist of an input layer, an output layer and at least one hidden layer in between. Each layer performs specific types of sorting and ordering ('feature hierarchy process'). These networks use sophisticated mathematical modelling to process data in complex ways [88].
<b>General Adversarial Networks (GANs)</b>	GANs belong to a class of generative models within the artificial intelligence (AI) and machine learning field and represent a powerful way of learning any kind of data distribution based on unsupervised learning. GANs aim to learn the true data distribution of the training (input) dataset and attempt to generate new data points from this distribution with some variations and not just reproducing the old data the model has been trained on. GANs try to use the power of neural networks (as described above) to learn a function to approximate the approach to model a distribution as close as possible to the real data [80].
<b>Autoencoders</b>	Autoencoders were originally introduced as a method for learning meaningful representations from data in an unsupervised manner and the concept of autoencoders in the context of artificial neural networks was first presented by Ballard....[89]. An autoencoder is a feed-forward deep neural network that first compresses the input data into a more compact representation and then attempt to reconstruct the original input by using an in-between layer which restricts the amount of information that travels within the network. Autoencoders have been frequently used for data compression and dimensionality reduction and can learn nonlinear relationships [84].
<b>Minority class</b>	This term refers to classification predictive modelling in machine learning (ML) and involves predicting a class label for a given observation. Most of the ML algorithms used for classification in predictive modelling were designed with the assumption of an equal number of examples for each class. However, imbalanced classification problems might occur (i.e., where the distribution of examples across the known classes is bias or skewed) and one of the target classes can contain a much smaller number of instances than the other classes (minority class) [90].
<b>Bayesian Network</b>	Bayesian networks represent systems as a network of interactions between variables from primary cause to final outcome, with all cause-effect assumptions made explicit....[91]. They are a type of probabilistic graphical model that uses Bayesian inference for probability computations. Bayesian networks aim to model conditional dependence, and therefore causation, by representing conditional dependence by edges in a directed graph. Through these relationships, one can efficiently conduct inference on the random variables in the graph through the use of factors [92].
<b>Function approximation</b>	Function approximation is a technique for estimating an unknown underlying function using historical or available observations from the domain. Artificial neural networks learn to approximate a function [93].

