

---

# A DENSITY RATIO FRAMEWORK FOR EVALUATING THE UTILITY OF SYNTHETIC DATA

---

A PREPRINT

**Thom Benjamin Volker** 

Methodology and Statistics | Methodology  
Utrecht University | Statistics Netherlands  
Utrecht, 3584CH  
[t.b.volker@uu.nl](mailto:t.b.volker@uu.nl)

**Peter-Paul de Wolf**

Methodology  
Statistics Netherlands  
The Hague, 2490HA  
[pp.dewolf@cbs.nl](mailto:pp.dewolf@cbs.nl)

**Erik-Jan van Kesteren**

Methodology and Statistics  
Utrecht University  
Utrecht, 3584CH  
[e.vankesteren1@uu.nl](mailto:e.vankesteren1@uu.nl)

November 28, 2023

## ABSTRACT

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

**Keywords** TODO • TODO

## 1 Introduction

Openly accessible research data accelerates scientific progress tremendously [ADD REFERENCES]. Open data allows third-party researchers to answer research questions with already collected data, substantially lowering the costs of research. Sharing data in combination with code allows others to validate research findings and build upon the work. Students can benefit from open data, as it fosters education with realistic data, as well as the general public, through stimulating citizen science projects. However, openly available research data is often at odds with privacy

constraints on the dissemination of the data. The information may cause harm to individuals or organizations if made publicly available. Additionally, sharing their information openly may withhold them from participating in future research. These constraints have been named among the biggest hurdles in the advancement of computational social science (Lazer et al. 2009), and among top reasons for companies to not share their data with researchers (Future of Privacy Forum 2017). To overcome these obstacles, data providers can employ a suite of different disclosure limitation techniques before sharing data, for example top-coding, record-swapping or adding noise (e.g., Hundepool et al. 2012; Willenborg and De Waal 2012).

Recently, synthetic data has gained substantial traction as a means to disclosure limitation. National statistical institutes and other government agencies have adopted the synthetic data framework to facilitate the use of their data (e.g., Abowd, Stinson, and Benedetto 2006; Hawala 2008; Drechsler 2012). Researchers started using synthetic data to train machine learning models (Nikolenko 2021), and began to share synthetic versions of their research data to comply with open science standards (e.g., van de Wiel et al. 2023; Obermeyer et al. 2019; Zettler et al. 2021). The idea of synthetic data is to replace some, or all, of the observed values in a data set by synthetic records that are generated from some model (e.g., Little 1993; Rubin 1993; Drechsler 2011). If only some values are replaced, disclosure risks can be reduced because the sensitive or identifying values do not correspond to their true values anymore. If all values are replaced, there is also no one-to-one mapping between the original and the synthetic data, further reducing the disclosure risk.

Many different methods have been proposed to generate synthetic data. Traditionally, these were closely connected to methods used for multiple imputation of missing data, such as joint modelling [misd2003], sequential regressions (Nowok, Raab, and Dibben 2016) or fully conditional specification (Drechsler and Reiter 2011; Volker and Vink 2021). The flexibility of sequential regressions and fully conditional specification is commonly combined with non-parametric imputation models, such as classification and regression trees (Reiter 2005), random forests (Caiola and Reiter 2010) or support vector machines (Drechsler 2010), to avoid distributional assumptions and easily model non-linear relationships. Recently, significant improvements in generative modelling sparked the scientific interest in synthetic data in the computer science community, leading to novel synthesis methods (e.g., Patki, Wedge, and Veeramachaneni 2016; Xu et al. 2019). Combined with work on formal privacy guarantees such as differential privacy, this resulted in new models that explicitly control the level of privacy risk in synthesis methods (Jordon, Yoon, and Schaar 2019; Torkzadehmahani, Kairouz, and Paten 2019).

## ***TODO***

All methods for statistical disclosure limitation alter the data before these are provided to the public. By doing so, the utility of the provided data is always lower than the utility of the original data, because some of the information in the data is sacrificed to protect the privacy of the respondents. The questions that naturally arise are how much information in the original data is actually sacrificed, and how useful the provided data are? Answering this question allows researchers to decide what the altered data can and cannot be used for, and to evaluate the worth of conclusions

drawn on the basis of these data. After all, inferences from the altered data are valid only up to the extent that the perturbation methods approximate the true data-generating mechanism. For data providers, a detailed assessment of the quality of the altered data can guide the procedure of altering the data. Statistical disclosure limitation is often an iterative process: some disclosure limitation technique is applied on the data, after which the result is investigated and modifications are made to applied process. Good measures of data quality are essential to determine the appropriate mechanisms used to protect the data, and can help to improve the utility of the data that will be disseminated.

In the statistical disclosure control literature, two different branches of utility measures have been distinguished: specific utility measures and general utility measures. *Add one/two sentences on the merit of visualization when assessing utility of altered data.* Specific utility measures focus on similarity of results obtained from analyses performed on the altered data and the original data. For example, after fitting the same analysis model on both data sets, one can calculate the confidence interval overlap of the estimated parameters (Karr et al. 2006). Alternative measures are ellipsoidal overlap (Karr et al. 2006), which extends to confidence interval overlap to a measure that addresses the joint distribution of all model parameters simultaneously, the standardized absolute difference between estimates (Snok et al. 2018), and the ratio of estimates for tabular count data (taub2020impact?). As these measures quantify similarity between estimates from analyses performed on the observed and altered data, they are informed only to the extent that data users will recreate those analyses. This can be highly useful if the data is provided for reproducibility purposes (e.g., for third parties to evaluate analysis scripts). However, the goal of distributing the protected data is often to allow researchers to do novel research with the data. In many practical situations, data providers thus have only limited knowledge on the analyses that will be performed with the altered data. Covering the entire set of potentially relevant analyses is therefore not feasible. If it was, the data providers could simply report the (potentially privacy-protected) results of those analyses performed on the real data, so that access to the (perturbed) data no longer yields additional benefits (for a similar argument, see Drechsler 2022). Additionally, similarity between results on the analyses that have been performed gives no guarantee that the results will also be similar for other analyses. Hence, when determining how useful the altered data is for novel research, specific utility measures are only of limited use.

General utility measures attempt to capture how similar the multivariate distributions of the observed and altered data are. This can be done by, for instance, estimating the Kullback-Leibler divergence between the distributions of the observed and altered data (Karr et al. 2006). An alternative strategy is to try to discriminate between the observed and altered data, as is done with the  $pMSE$  (Snok et al. 2018; woo\_utility\_2009?). In essence, the  $pMSE$  quantifies how well one can predict whether observations are from the observed or the altered data. If better one can do this, the more pronounced the differences between the observed and altered data ought to be. However, various authors have criticized general utility measures for being too broad. That is, important discrepancies between the real and altered data might be missed, and an altered data set that is good in general (i.e., has high global utility) might still provide results that are far from the truth for some analyses (see, e.g., Drechsler 2022). Additionally, it is not straightforward to determine which prediction model to use for calculating the  $pMSE$ . Specifying a good prediction model in itself may be a challenging task, especially when the number of variables is large. When good models are available, different

models, or even different choices of hyperparameters, may yield different results, potentially rendering ambiguity with respect to which altered data set is best. Lastly, the output of global utility measures can be hard to interpret, and say little about the regions in which the synthetic data do not resemble the true data accurately enough. That is, they give little guidance on how the quality of the altered data can be improved.

---

## Section 6: our contribution

*Moet nog verder uitgewerkt worden*

1. We introduce density ratio estimation to the field of statistical disclosure control. Short remark on the idea that density ratio estimation is a complicated endeavor, especially if the goal is to compare distinct densities. Having to estimate just a single density (ratio) is generally much easier.
2. Note that density ratio estimation can capture specific and general utility measures into a common framework by being applicable on the level of the entire data, but also on the subset of variables that is relevant in an analysis. Additionally, note that confidence interval overlap, ellipsoidal overlap, but also  $pMSE$  and Kullback-Leibler divergence, are closely related to density ratio estimation, and can be considered from this perspective.
3. Create a new utility metric based on density ratio estimation (probability with respect to some reference distribution as in permutation testing).
4. Because density ratio estimation can be difficult when there are many variables, we use dimension reduction techniques to capture most of the variability in the data in fewer dimensions on which density ratio estimation can be applied. A by-product of this is that the lower-dimensional subspace allows to create visualizations of deviations from the observed data.
5. Perform a simulation study to give indications about which methods to use (think about how to do this).
6. Implement all this in an R-package

## Section 6: outline article

In the next section, we describe density ratio estimation and discuss how this method can be used as to measure utility. Subsequently, we provide multiple examples that show how density ratio estimation works in the context of evaluating the quality of synthetic data. Hereafter, we show in multiple simulations that the method is superior (HOPEFULLY) to current global utility measures as the  $pMSE$ . Lastly, we discuss the advantages and disadvantages of density ratio estimation as a utility measure.

## 2 Methodology

@Peter-Paul: Eventueel een korte beschrijving van data perturbation techniques/synthetic data generation hier. Denk je dat dit wat toevoegt hier?

### Section 1: density ratio estimation

In essence, the goal of utility measures is to quantify the similarity between the multivariate distribution of the observed data with the distribution of the altered data. If the used data perturbation techniques, or synthetic data generation models, approximate the distribution of the real data sufficiently, these distributions should be highly similar, and analyses on the two data sets should give similar results. However, estimating the probability distribution of a data set is known to be one of the most complicated challenges in statistics [E.G. Vapnik 1998]. Estimating the probability distribution for both observed and altered data can lead to errors in both, artificially magnifying discrepancies between the two. Hence, subsequent comparisons will be affected by these errors. The procedure can be simplified by using density ratio estimation, because this only requires to estimate a single density.

Introduce density ratio estimation as a utility measure. What does this measure mean/how to interpret it. How to make decisions based on this measure.

Say something on whether (and if so, how) categorical variables can be incorporated as well.

### Section 3: theoretical comparison with conventional approaches for general utility assessment

Relate density ratio estimation to specific and general utility measures. Pick one/two specific utility measures and relate these to density ratio estimation (ratio of estimates seems straightforward, as well as confidence interval overlap and ellipsoidal overlap).

Relate density ratio estimation to  $pMSE$  and KL divergence (to some extent, both are generalizations of density ratio estimation, or at least are conceptually similar). Give some more information on the  $pMSE$ , describe what its shortcomings are. The quality of the  $pMSE$  highly depends on the model used to calculate the propensity scores. Perhaps give an example of logistic regression, which basically estimates whether the conditional mean of the observed and altered data is the same. Explain how density ratio estimation can overcome the shortcomings of the previously mentioned methods.

### Section 4: Dimension reduction for utility

The difficulty of density ratio estimation increases with the dimensionality of the data. Therefore, we follow previous recommendations to incorporate dimensionality reduction techniques in density ratio estimation.

Shortly name examples of dimensionality reduction techniques (i.e., PCA; LFDA or UMAP).

A useful by-product of dimension reduction is that it allows to create visualizations, and these visualizations can be used to get more insight in discrepancies between observed and altered data. Show what such visualizations can look like, and how they can help.

### 3 Simulations

#### 3.1 Small illustration / example with multivariate Gaussian distributions.

1. Simple, multivariate normal simulation (e.g., two correlation structures, two sample sizes, so  $2 \times 2$  full factorial design); basically similar to what we did already.

#### 3.2 More complex simulation, more variables, non-linearities, perhaps using real data.

2. More advanced simulation (e.g., some non-linearities, different sample sizes)

Have to think about this in more detail still.

### 4 Real data example

Clinical records heart-failure data? Misschien ook niet, nog over nadenken.

Exemplify how utility measures could (should!) be used to improve the quality of the altered data (e.g., illustrate how models can be adjusted iteratively based on utility assessment).

*Some notes to self*

Current ways to assess the utility?

- pMSE - logistic, regression, CART models (Snoke, Raab, Nowok, Dibben & Slavkovic, 2018; General and specific utility measures for synthetic data AND Woo, Reiter, Oganian & Karr, 2009; Global measures of data utility for microdata masked for disclosure limitation)
- Kullback-Leiber divergence (Karr, Kohnen, Oganian, Reiter & Sanil, 2006; A framework for evaluating the utility of data altered to protect confidentiality).
- According to multiple authors, both specific and general utility measures have important drawbacks (see Drechsler Utility PSD; cites others). Narrow measures potentially focus on analyses that are not relevant for the end user, and do not generalize to the analyses that are relevant. Global utility measures are generally too broad, and important deviations in the synthetic data might be missed. Moreover, the measures are typically hard to interpret.
- See Drechsler for a paragraph on fit for purpose measures, that lie between general and specific utility measures (i.e., plausibility checks such as non-negativity; goodness of fit measures as  $\chi^2$  for cross-tabulations; Kolmogorov-Smirnov).
- Drechsler also illustrates that the standardized pMSE has substantial flaws, as the results are highly dependent on the model used to estimate the propensity scores, and unable to detect important differences in the utility for most of the model specifications. Hence, it is claimed that a thorough assessment of utility is required.

Things to add in new version: - High dimensional example - Expansion of discussion points: empirical example with weighted analyses individual data utility - Privacy remark? - When does it not work?

## 5 Methodology

TO DO

## 6 Simulations

TO DO

## 7 Real data example

TO DO

## 8 Results

TO DO

## 9 Discussion and conclusion

TO DO

## References

- Abowd, John M., Martha Stinson, and Gary Benedetto. 2006. "Final Report to the Social Security Administration on the SIPP/SSA/IRS Public Use File Project." Longitudinal Employer-Household Dynamics Program, U.S. Bureau of the Census, Washington, DC. <https://ecommons.cornell.edu/bitstream/handle/1813/43929/SSAfinal.pdf?sequence=3&isAllowed=y>.
- Caiola, Gregory, and Jerome P. Reiter. 2010. "Random Forests for Generating Partially Synthetic, Categorical Data." *Transactions on Data Privacy* 3: 27–42. <https://doi.org/10.5555/1747335.1747337>.
- Drechsler, Jörg. 2010. "Using Support Vector Machines for Generating Synthetic Datasets." In *Privacy in Statistical Databases*, edited by Josep Domingo-Ferrer and Emmanouil Magkos, 148–61. Berlin, Heidelberg: Springer Berlin Heidelberg. [https://doi.org/10.1007/978-3-642-15838-4\\_14](https://doi.org/10.1007/978-3-642-15838-4_14).
- . 2011. *Synthetic Datasets for Statistical Disclosure Control: Theory and Implementation*. New York: Springer Science & Business Media. <https://doi.org/10.1007/978-1-4614-0326-5>.
- . 2012. "New Data Dissemination Approaches in Old Europe – Synthetic Datasets for a German Establishment Survey." *Journal of Applied Statistics* 39 (2): 243–65. <https://doi.org/10.1080/02664763.2011.584523>.
- . 2022. "Challenges in Measuring Utility for Fully Synthetic Data." In *Privacy in Statistical Databases*, edited by Josep Domingo-Ferrer and Maryline Laurent, 220–33. Cham: Springer International Publishing. [https://doi.org/10.1007/978-3-031-13945-1\\_16](https://doi.org/10.1007/978-3-031-13945-1_16).

- Drechsler, Jörg, and Jerome P Reiter. 2011. "An Empirical Evaluation of Easily Implemented, Nonparametric Methods for Generating Synthetic Datasets." *Computational Statistics & Data Analysis* 55 (12): 3232–43.
- Future of Privacy Forum. 2017. "Understanding Corporate Data Sharing Decisions: Practices, Challenges, and Opportunities for Sharing Corporate Data with Researchers."
- Hawala, Sam. 2008. *Producing Partially Synthetic Data to Avoid Disclosure*. <http://www.asasrms.org/Proceedings/y2008/Files/301018.pdf>.
- Hundepool, Anco, Josep Domingo-Ferrer, Luisa Franconi, Sarah Giessing, Eric Schulte Nordholt, Keith Spicer, and Peter-Paul De Wolf. 2012. *Statistical Disclosure Control*. John Wiley & Sons. <https://doi.org/10.1002/9781118348239>.
- Jordon, James, Jinsung Yoon, and Mihaela van der Schaar. 2019. "PATE-GAN: Generating Synthetic Data with Differential Privacy Guarantees." In *International Conference on Learning Representations*. <https://openreview.net/forum?id=S1zk9iRqF7>.
- Karr, Alan F., Christine N. Kohnen, Anna Oganian, Jerome P. Reiter, and Ashish P. Sanil. 2006. "A Framework for Evaluating the Utility of Data Altered to Protect Confidentiality." *The American Statistician* 60 (3): 224–32. <https://doi.org/10.1198/000313006X124640>.
- Lazer, David, Alex Pentland, Lada Adamic, Sinan Aral, Albert-László Barabási, Devon Brewer, Nicholas Christakis, et al. 2009. "Computational Social Science." *Science* 323 (5915): 721–23. <https://doi.org/10.1126/science.1167742>.
- Little, Roderick J. A. 1993. "Statistical Analysis of Masked Data." *Journal of Official Statistics* 9 (2): 407–7. <https://www.scb.se/contentassets/ca21efb41fee47d293bbee5bf7be7fb3/statistical-analysis-of-masked-data.pdf>.
- Nikolenko, Sergey I. 2021. *Synthetic Data for Deep Learning*. Springer International Publishing. <https://doi.org/10.1007/978-3-030-75178-4>.
- Nowok, Beata, Gillian M. Raab, and Chris Dibben. 2016. "Synthpop: Bespoke Creation of Synthetic Data in R." *Journal of Statistical Software* 74 (11). <https://doi.org/10.18637/jss.v074.i11>.
- Obermeyer, Ziad, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. 2019. "Dissecting Racial Bias in an Algorithm Used to Manage the Health of Populations." *Science* 366 (6464): 447–53. <https://doi.org/10.1126/science.aax2342>.
- Patki, Neha, Roy Wedge, and Kalyan Veeramachaneni. 2016. "The Synthetic Data Vault." *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, October. <https://doi.org/10.1109/dsaa.2016.49>.
- Reiter, Jerome P. 2005. "Using CART to Generate Partially Synthetic Public Use Microdata." *Journal of Official Statistics* 21 (3): 441–62. <https://www.scb.se/contentassets/ca21efb41fee47d293bbee5bf7be7fb3/using-cart-to-generate-partially-synthetic-public-use-microdata.pdf>.
- Rubin, Donald B. 1993. "Statistical Disclosure Limitation." *Journal of Official Statistics* 9 (2): 461–68. <https://www.scb.se/contentassets/ca21efb41fee47d293bbee5bf7be7fb3/discussion-statistical-disclosure-limitation2.pdf>.
- Snoke, Joshua, Gillian M. Raab, Beata Nowok, Chris Dibben, and Aleksandra Slavkovic. 2018. "General and Specific Utility Measures for Synthetic Data." *Journal of the Royal Statistical Society. Series A (Statistics in Society)* 181 (3): pp. 663–688. <https://doi.org/10.1111/rssa.12358>.



- Torkzadehmahani, Reihaneh, Peter Kairouz, and Benedict Paten. 2019. “DP-CGAN: Differentially Private Synthetic Data and Label Generation.” In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*. <https://doi.org/10.1109/cvprw.2019.00018>.
- van de Wiel, Mark A., Gwenaël G. R. Leday, Jeroen Hoogland, Martijn W. Heymans, Erik W. van Zwet, and Ailko H. Zwinderman. 2023. “Think Before You Shrink: Alternatives to Default Shrinkage Methods Can Improve Prediction Accuracy, Calibration and Coverage.” <https://doi.org/10.48550/ARXIV.2301.09890>.
- Volker, Thom Benjamin, and Gerko Vink. 2021. “Anonymiced Shareable Data: Using Mice to Create and Analyze Multiply Imputed Synthetic Datasets.” *Psych* 3 (4): 703–16. <https://doi.org/10.3390/psych3040045>.
- Willenborg, Leon, and Ton De Waal. 2012. *Elements of Statistical Disclosure Control*. Springer Science & Business Media. <https://doi.org/10.1007/978-1-4613-0121-9>.
- Xu, Lei, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. 2019. “Modeling Tabular Data Using Conditional GAN.” In *Advances in Neural Information Processing Systems*, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. dAlché-Buc, E. Fox, and R. Garnett. Vol. 32. Curran Associates, Inc. [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/254ed7d2de3b23ab10936522dd547b78-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/254ed7d2de3b23ab10936522dd547b78-Paper.pdf).
- Zettler, Ingo, Christoph Schild, Lau Lilleholt, Lara Kroencke, Till Utesch, Morten Moshagen, Robert Böhm, Mitja D. Back, and Katharina Geukes. 2021. “The Role of Personality in COVID-19-Related Perceptions, Evaluations, and Behaviors: Findings Across Five Samples, Nine Traits, and 17 Criteria.” *Social Psychological and Personality Science* 13 (1): 299–310. <https://doi.org/10.1177/19485506211001680>.