

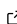


# Metasyn: Transparent Generation of Synthetic Tabular Data with Privacy Guarantees

Raoul Schram<sup>1\*</sup>, Samuel Spithorst<sup>1</sup>, and Erik-Jan van Kesteren<sup>1,2\*</sup>

<sup>1</sup> Utrecht University, The Netherlands <sup>2</sup> ODISSEI: Open Data Infrastructure for Social Science and Economic Innovations, The Netherlands ¶ Corresponding author \* These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

## Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

Synthetic data is a promising tool for improving the accessibility of datasets that are otherwise too sensitive to be shared publicly. To this end, we introduce metasyn, a Python package for generating synthetic data from tabular datasets. Unlike existing synthetic data generation software, metasyn is built on a simple generative model with a “naïve” marginal independence assumption — an explicit choice that removes multivariate information from the synthetic data. It makes this trade-off in order to maintain transparency and auditability, to keep information leakage to a minimum, and even to enable privacy or disclosure risk guarantees through a plug-in system. While the analytical validity of the generated data is thus intentionally limited, its potential uses are broad, including exploratory analyses, code development and testing, and external communication and teaching ([van Kesteren, 2024](#)). Metasyn is flexible, scalable, and easily extended to meet diverse privacy needs.



Figure 1: Logo of the metasyn project.

## Statement of need

Metasyn is a python package for generating synthetic data with a focus on privacy and disclosure control. It is aimed at owners of sensitive datasets such as public organisations, research groups, and individual researchers who want to improve the accessibility of their data for research and reproducibility by others. The goal of metasyn is to make it easy for data owners to share the structure and an approximation of the content of their data with others while keeping privacy concerns to a minimum.

With this goal in mind, metasyn distinguishes itself from existing software for generating synthetic data (e.g., [Nowok et al., 2016](#); [Ping et al., 2017](#); [Templ et al., 2017](#)) by strictly limiting the statistical information from the real data in the produced synthetic data. Metasyn explicitly avoids generating synthetic data with high analytical validity; instead, the synthetic data has realistic structure and plausible values, but multivariate relations are omitted (“augmented plausible synthetic data”; ([Bates et al., 2019](#))). Moreover, our system provides an **auditable**

31 and **editable intermediate representation** in the form of a .json metadata file from which new  
32 data can be synthesized.

33 These choices enable the software to generate synthetic data with **privacy and disclosure**  
34 **guarantees** through a plug-in system, recognizing that different data owners have different  
35 needs and definitions of privacy. A data owner can define under which conditions they would  
36 accept open distribution of their synthetic data — be it based on differential privacy (Dwork,  
37 2006), statistical disclosure control (Hundepool et al., 2012), k-anonymity (Sweeney, 2002), or  
38 another specific definition of privacy. As part of the initial release of metasyn, we publish a  
39 **plug-in** following the disclosure control guidelines from Eurostat (Bond et al., 2015).

## 40 Software features

41 At its core, metasyn has three main functions:

- 42 1. **Estimation**: Automatically select distributions and fit them to a properly formatted  
43 tabular dataset, optionally with additional privacy guarantees.
- 44 2. **(De)serialization**: Create an intermediate representation of the fitted model for auditing,  
45 editing, and exporting.
- 46 3. **Generation**: Generate new synthetic datasets based on a fitted model.

## 47 Estimation

48 The generative model for multivariate datasets in metasyn makes the simplifying assumption  
49 of marginal independence: each column is considered separately, just as is done in e.g., naïve  
50 Bayes classifiers (Hastie et al., 2009). Formally, this leads to the following generative model  
51 for the  $K$ -variate data  $\mathbf{x}$ :

$$p(\mathbf{x}) = \prod_{k=1}^K p(x_k) \quad (1)$$

52 There are many advantages to this naïve approach when compared to more advanced generative  
53 models: it is transparent and explainable, it is able to flexibly handle data of mixed types, and  
54 it is computationally scalable to high-dimensional datasets.

55 Model estimation starts with an appropriately pre-processed data frame, meaning it is tidy  
56 (Wickham, 2014), each column has the correct data type, and missing data are represented by  
57 a missing value. Internally, our software uses the polars data frame library (Vink et al., 2024),  
58 as it is performant, has consistent data types, and natively supports missing data (i.e., null  
59 values). A simple example source table could look like this (note that categorical data has the  
60 appropriate cat data type, not str):

ID	fruits	B	cars	optional
---	---	---	---	---
i64	cat	i64	cat	i64
1	banana	5	beetle	28
2	banana	4	audi	300
3	apple	3	beetle	null
4	apple	2	beetle	2
5	banana	1	beetle	-30

72 For each data type, a set of candidate distributions is fitted (see Table Table 1), and then  
73 metasyn selects the one with the lowest BIC (Neath & Cavanaugh, 2012). For distributions

74 where BIC computation is impossible (e.g., for the string data type) a pseudo-BIC is created  
75 that trades off fit and complexity of the underlying models.

**Table 1:** Candidate distributions associated with data types in the core metasyn package.

Variable type	Example	Candidate distributions
categorical	yes/no, country	Categorical (Multinoulli), Constant
continuous	1.0, 2.1, ...	Uniform, Normal, LogNormal, TruncatedNormal, Exponential, Constant
discrete	1, 2, ...	Poisson, Uniform, Normal, TruncatedNormal, Categorical, Constant
string	A108, C122, some words	Regex, Categorical, Faker, FreeText, Constant
date/time	2021-01-13, 01:40:12	Uniform, Constant

76 From this table, the string distributions deserve special attention as they are not commonly  
77 encountered as probability distributions. Regex (regular expression) inference is performed  
78 on structured strings using the companion package [RegexModel](#). It is able to automatically  
79 detect structure such as room numbers (A108, C122, B109), e-mail addresses, websites, and  
80 more, which it summarizes using a probabilistic variant of regular expressions. The FreeText  
81 distribution detects the language (using [lingua](#)) and randomly picks words from that language.  
82 The [Faker](#) distribution can generate specific data types such as localized addresses, when  
83 pre-specified by the user.

84 Generative model estimation with metasyn can be performed as follows:

```

from metasyn import MetaFrame, VarSpec

# "ID" column is the primary key,
# thus should generate unique values.
# "B" column is not, despite unique
# values in the dataframe
specs = [
    VarSpec("ID", unique=True),
    VarSpec("B", unique=False),
]

# create metaframe
mf = MetaFrame.fit_dataframe(df, var_specs=specs)

```

## 85 Serialization and deserialization

86 After a fitted model object is created, metasyn allows it to be transparently stored in a  
87 human- and machine-readable .json file. This file can be considered as metadata: it contains  
88 dataset-level descriptive information as well as variable-level information. The header contains  
89 the following dataset-level information:

```

{"n_rows": 5,
 "n_columns": 5,
 "provenance": {
   "created by": {
     "name": "metasyn",
     "version": "1.0.2"
   }
 }
}

```

```
    },  
    "creation time": "2024-08-12T12:20:36.022017"  
  }  
}
```

90 Then, for each variable the file contains descriptive and statistical information, as in the  
91 following example:

```
{  
  "name": "fruits",  
  "type": "categorical",  
  "dtype": "Categorical(ordering='physical')",  
  "prop_missing": 0.0,  
  "distribution": {  
    "implements": "core.multinoulli",  
    "version": "1.0",  
    "provenance": "builtin",  
    "class_name": "MultinoulliDistribution",  
    "unique": false,  
    "parameters": {  
      "labels": ["apple", "banana"],  
      "probs": [0.4, 0.6]  
    }  
  },  
  "creation_method": { "created_by": "metasyn" }  
}
```

92 There are several advantages to creating such a serialized representation. It can be manually  
93 audited, edited, and after exporting this file, an unlimited number of synthetic records can be  
94 created without incurring additional privacy risks.

95 Serialization and deserialization with metasyn can be performed as follows:

```
# write a fitted MetaFrame to json  
mf.export("fruits.json")  
  
# then, audit and export json from secure environment  
  
# outside the secure environment, load json into MetaFrame  
mf_out = MetaFrame.from_json("fruits.json")
```

## 96 Data generation

97 For each variable in a fitted or deserialized model object, metasyn can randomly sample  
98 synthetic datapoints. Data generation (or synthetization) in metasyn can be performed as  
99 follows:

```
from metasyn import MetaFrame  
  
# load json into a metadataset object  
mf = MetaFrame.from_json("metasyn_example.json")  
  
# create a fake dataset  
df_syn = mf.synthesize(10)
```

100 This may result in the following polars data frame<sup>1</sup>. Note that missing values in the optional  
101 column are appropriately reproduced as well, courtesy of the "prop\_missing" entry in the  
102 metadata format.

<sup>1</sup>This polars dataframe can be easily converted to a pandas dataframe using `df_syn.to_pandas()`

103 shape: (10, 5)

ID	fruits	B	cars	optional
---	---	---	---	---
i64	cat	i64	cat	i64
1	banana	4	beetle	null
2	banana	3	audi	null
3	banana	1	beetle	223
4	banana	0	beetle	258
...	...	...	...	...
7	banana	3	beetle	298
8	banana	2	beetle	67
9	banana	4	beetle	-30
10	banana	2	beetle	172

## 119 Plug-ins and automatic privacy

120 In addition to the core features described above, the metasyn package allows for plug-ins: add-  
 121 on packages that alter the behaviour of the parameter estimation. Through this system, privacy  
 122 guarantees can be built into metasyn ([privacy plug-in template](#)) and additional distributions  
 123 can be supported ([distribution plug-in template](#)). The `metasyn-disclosure-control` plug-in  
 124 implements output guidelines from Eurostat ([Bond et al., 2015](#)) by including a micro-aggregation  
 125 step in the `fit()` method. In this way, information transfer from the sensitive real data to the  
 126 synthetic public data can be further limited. Disclosure control is done as follows:

```
from metasyn import MetaFrame
from metasyncontrib.disclosure import DisclosurePrivacy

mf = MetaFrame.fit_dataframe(df, privacy=DisclosurePrivacy())
```

## 127 Acknowledgements

128 This research was conducted in whole or in part using ODISSEI, the Open Data Infrastructure  
 129 for Social Science and Economic Innovations (<https://ror.org/03m8v6t10>)

130 The metasyn project is supported by the FAIR Research IT Innovation Fund of Utrecht  
 131 University (March 2023)

## 132 References

- 133 Bates, A., Spakulová, I., Dove, I., & Meador, A. (2019). *ONS methodology work-*  
 134 *ing paper series number 16—synthetic data pilot*. [https://www.ons.gov.uk/](https://www.ons.gov.uk/methodology/methodologicalpublications/generalmethodology/onsworkingpaperseries/onsmethodologyworkingpaperseriesnumber16syntheticdatapilot)  
 135 [methodology/methodologicalpublications/generalmethodology/onsworkingpaperseries/](https://www.ons.gov.uk/methodology/methodologicalpublications/generalmethodology/onsworkingpaperseries/onsmethodologyworkingpaperseriesnumber16syntheticdatapilot)  
 136 [onsmethodologyworkingpaperseriesnumber16syntheticdatapilot](https://www.ons.gov.uk/methodology/methodologicalpublications/generalmethodology/onsworkingpaperseries/onsmethodologyworkingpaperseriesnumber16syntheticdatapilot)
- 137 Bond, S., Brandt, M., & Wolf, P. de. (2015). *Guidelines for the checking of output based*  
 138 *on microdata research*. Eurostat. [https://web.archive.org/web/20160408145718/http://dwbproject.org/export/sites/default/about/public\\_deliverables/dwb\\_d11-8\\_](https://web.archive.org/web/20160408145718/http://dwbproject.org/export/sites/default/about/public_deliverables/dwb_d11-8_synthetic-data_cta-ecta_output-checking-guidelines_final-reports.zip)  
 139 [synthetic-data\\_cta-ecta\\_output-checking-guidelines\\_final-reports.zip](https://web.archive.org/web/20160408145718/http://dwbproject.org/export/sites/default/about/public_deliverables/dwb_d11-8_synthetic-data_cta-ecta_output-checking-guidelines_final-reports.zip)
- 141 Dwork, C. (2006). Differential privacy. *International Colloquium on Automata, Languages,*  
 142 *and Programming*, 1–12. [https://doi.org/10.1007/11787006\\_1](https://doi.org/10.1007/11787006_1)

- 143 Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The elements of*  
144 *statistical learning: Data mining, inference, and prediction* (Vol. 2). Springer. <https://doi.org/10.1007/978-0-387-84858-7>  
145
- 146 Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K.,  
147 & De Wolf, P.-P. (2012). *Statistical disclosure control*. Wiley & Sons, Chichester.  
148 <https://doi.org/10.1002/9781118348239>
- 149 Neath, A. A., & Cavanaugh, J. E. (2012). The bayesian information criterion: Background,  
150 derivation, and applications. *Wiley Interdisciplinary Reviews: Computational Statistics*,  
151 4(2), 199–203. <https://doi.org/10.1002/wics.199>
- 152 Nowok, B., Raab, G. M., & Dibben, C. (2016). Synthpop: Bespoke creation of synthetic data  
153 in r. *Journal of Statistical Software*, 74, 1–26. <https://doi.org/10.18637/jss.v074.i11>
- 154 Ping, H., Stoyanovich, J., & Howe, B. (2017). Datasynthesizer: Privacy-preserving synthetic  
155 datasets. *Proceedings of the 29th International Conference on Scientific and Statistical*  
156 *Database Management*, 1–5. <https://doi.org/10.1145/3085504.3091117>
- 157 Sweeney, L. (2002). K-anonymity: A model for protecting privacy. *International Journal of*  
158 *Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05), 557–570. <https://doi.org/10.1142/S0218488502001648>  
159
- 160 Templ, M., Meindl, B., Kowarik, A., & Dupriez, O. (2017). Simulation of synthetic complex  
161 data: The r package simPop. *Journal of Statistical Software*, 79(10), 1–38. <https://doi.org/10.18637/jss.v079.i10>  
162
- 163 van Kesteren, E.-J. (2024). To democratize research with sensitive data, we should make  
164 synthetic data more accessible. *arXiv Preprint arXiv:2404.17271*. <https://doi.org/10.48550/arXiv.2404.17271>  
165
- 166 Vink, R., Gooijer, S. de, Beedie, A., Gorelli, M. E., Guo, W., Zundert, J. van, Peters,  
167 O., Hulselmans, G., nameexhaustion, Grinstead, C., Marshall, Burghoorn, G., chielP,  
168 Turner-Trauring, I., Santamaria, M., Heres, D., Mitchell, L., Magarick, J., ibENPC,  
169 ... Brannigan, L. (2024). *Pola-rs/polars: Python polars* (py-1.4.1). Zenodo. <https://doi.org/10.5281/zenodo.7697217>  
170
- 171 Wickham, H. (2014). Tidy data. *Journal of Statistical Software*, 59(10), 1–23. <https://doi.org/10.18637/jss.v059.i10>  
172