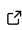# Metasyn: Transparent Generation of Synthetic Tabular Data with Privacy Guarantees

**Raoul Schram** [1*], **Samuel Spithorst** [1], and **Erik-Jan van Kesteren** [1,2*¶]

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

## 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 that removes multivariate information from the synthetic data. This choice enables transparency and auditability, keeps information leakage to a minimum, and enables privacy 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).

**Figure 1:** Logo of the metasyn project.

## Statement of need

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

These choices enable the software to generate synthetic data with **privacy and disclosure**

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

# Software features

At its core, metasyn has three main functions:

1. **Estimation**: Fit a generative model to a properly formatted tabular dataset, optionally with additional privacy guarantees.
2. **(De)serialization**: Create an intermediate representation of the fitted model for auditing, editing, and exporting.
3. **Generation**: Generate new synthetic datasets based on a fitted model.

## Estimation

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

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

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

Model estimation starts with an appropriately pre-processed data frame, meaning it is tidy (Wickham, 2014), each column has the correct data type, and missing data are represented by a missing value. Internally, our software uses the polars data frame library (Vink et al., 2024), as it is performant, has consistent data types, and natively supports missing data (i.e., null values). A simple example source table could look like this (note that categorical data has the 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      |
```

For each data type, a set of candidate distributions is fitted (see Table 1), and then metasyn selects the one with the lowest BIC (Neath & Cavanaugh, 2012). For distributions where BIC computation is impossible (e.g., for the string data type) a pseudo-BIC is created that trades off fit and complexity of the underlying models.

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

| Data type | Candidate distributions |
|---|---|
| Categorical | Categorical, Constant |
| Continuous | Uniform, Normal, LogNormal, TruncatedNormal, Exponential, Constant |
| Discrete | Poisson, Uniform, Normal, TruncatedNormal, Categorical, Constant |
| String | Regex, Categorical, Faker, FreeText, Constant |
| Date/time | Uniform, Constant |

From this table, the string distributions deserve special attention as they are not commonly encountered as probability distributions. The regex (regular expression) distribution uses the package regexmodel to automatically detect structure such as room numbers (A108, C122, B109), e-mail addresses, or websites. The FreeText distribution detects the language (using lingua) and randomly picks words from that language. The Faker distribution can generate specific data types such as localized names and addresses pre-specified by the user.

Generative model estimation with metasyn can be performed as follows:

```python
from metasyn import MetaFrame
mf = MetaFrame.fit_dataframe(df)
```

## Serialization and deserialization

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

```json
{
  "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" }
}
```

This .json can be manually audited, edited, and after exporting this file, an unlimited number of synthetic records can be created without incurring additional privacy risks. Serialization and deserialization with metasyn can be performed as follows:

```python
mf.export("fruits.json")
mf_new = MetaFrame.from_json("fruits.json")
```

## Data generation

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

follows:

```
df_syn = mf.synthesize(3)
```

This may result in the following `polars` data frame[1]. Note that missing values in the `optional` column are appropriately reproduced as well.

| ID | fruits | B | cars | optional |
| --- | --- | --- | --- | --- |
| i64 | cat | i64 | cat | i64 |
| 1 | banana | 4 | beetle | null |
| 2 | banana | 3 | audi | null |
| 3 | banana | 2 | beetle | 172 |

## Plug-ins and automatic privacy

In addition to its core features, the `metasyn` package allows for plug-ins: packages that alter the distribution fitting behaviour. Through this system, privacy guarantees can be built into metasyn (privacy plug-in template) and additional distributions can be supported (distribution plug-in template). The `metasyn-disclosure-control` plug-in implements output guidelines from Eurostat (Bond et al., 2015) by including micro-aggregation. In this way, information transfer from the sensitive real data to the synthetic public data can be further limited. Disclosure control is done as follows:

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

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

## Acknowledgements

## References

Bates, A., Spakulová, I., Dove, I., & Mealor, A. (2019). *ONS methodology working paper series number 16—synthetic data pilot*. https://www.ons.gov.uk/methodology/methodologicalpublications/generalmethodology/onsworkingpaperseries/onsmethodologyworkingpaperseriesnumber16syntheticdatapilot

Bond, S., Brandt, M., & Wolf, P. de. (2015). *Guidelines for the checking of output based on microdata research*. Eurostat. https://web.archive.org/web/20160408145718/http://dwbproject.org/export/sites/default/about/public_deliveraples/dwb_d11-8_synthetic-data_cta-ecta_output-checking-guidelines_final-reports.zip

Dwork, C. (2006). Differential privacy. *International Colloquium on Automata, Languages, and Programming*, 1–12. https://doi.org/10.1007/11787006_1

---

[1]This `polars` dataframe can be easily converted to a `pandas` dataframe using `df_syn.to_pandas()`

Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (Vol. 2). Springer. https://doi.org/10.1007/978-0-387-84858-7

Hundepool, A., Domingo-Ferrer, J., Franconi, L., Giessing, S., Nordholt, E. S., Spicer, K., & De Wolf, P.-P. (2012). *Statistical disclosure control*. Wiley & Sons, Chichester. https://doi.org/10.1002/9781118348239

Neath, A. A., & Cavanaugh, J. E. (2012). The bayesian information criterion: Background, derivation, and applications. *Wiley Interdisciplinary Reviews: Computational Statistics*, *4*(2), 199–203. https://doi.org/10.1002/wics.199

Nowok, B., Raab, G. M., & Dibben, C. (2016). Synthpop: Bespoke creation of synthetic data in r. *Journal of Statistical Software*, *74*, 1–26. https://doi.org/10.18637/jss.v074.i11

Ping, H., Stoyanovich, J., & Howe, B. (2017). Datasynthesizer: Privacy-preserving synthetic datasets. *Proceedings of the 29th International Conference on Scientific and Statistical Database Management*, 1–5. https://doi.org/10.1145/3085504.3091117

Sweeney, L. (2002). K-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, *10*(05), 557–570. https://doi.org/10.1142/S0218488502001648

Templ, M., Meindl, B., Kowarik, A., & Dupriez, O. (2017). Simulation of synthetic complex data: The r package simPop. *Journal of Statistical Software*, *79*(10), 1–38. https://doi.org/10.18637/jss.v079.i10

van Kesteren, E.-J. (2024). To democratize research with sensitive data, we should make synthetic data more accessible. *arXiv Preprint arXiv:2404.17271*. https://doi.org/10.48550/arXiv.2404.17271

Vink, R., Gooijer, S. de, Beedie, A., Gorelli, M. E., Guo, W., Zundert, J. van, Peters, O., Hulselmans, G., nameexhaustion, Grinstead, C., Marshall, Burghoorn, G., chielP, Turner-Trauring, I., Santamaria, M., Heres, D., Mitchell, L., Magarick, J., ibENPC, … Brannigan, L. (2024). *Pola-rs/polars: Python polars* (py-1.4.1). Zenodo. https://doi.org/10.5281/zenodo.7697217

Wickham, H. (2014). Tidy data. *Journal of Statistical Software*, *59*(10), 1–23. https://doi.org/10.18637/jss.v059.i10