

BUYER BEWARE: UNDERSTANDING THE TRADE-OFF BETWEEN UTILITY AND RISK IN CART BASED MODELS USING SIMULATION DATA

UNECE Expert Meeting on Statistical Data Collection and Sources,
Barcelona,
15-17 October 2025

Jonathan Latner, PhD
Dr. Marcel Neunhoeffer
Prof. Dr. Jörg Drechsler



SECTION 1: INTRODUCTION

BACKGROUND AND MOTIVATION

- Background:
 - Synthetic data are increasingly used to share data while preserving privacy.
 - Numerous synthetic data generators (SDGs) using variety of methods
 - CART-based SDGs: high statistical utility with high privacy protection (Little et al., 2025)
- Motivation:
 - There is a general perception that generating synthetic data are easy.
 - It is not always clear if the resulting synthetic data are in fact providing privacy protection.

OVERVIEW

- Research question:
 - Do common privacy measures capture disclosure risk in synthetic data generated by CART models?
- Evaluate 3 types of privacy measures:
 1. Identity disclosure risk
 2. Attribute disclosure risk
 3. Bayesian estimation of disclosure risk
- 2 types of data:
 1. Simulated dataset (Reiter et al., 2014 design: 1,000 obs., 4 binary vars., unique case).
 2. Public survey data: Social Diagnosis 2011 (SD2011).
- Contributions:
 - Commonly used disclosure risk measures may not capture disclosure risk.
 - We propose some solutions for measuring disclosure risk (Bayesian).
 - More generally, users interested in empirical measures of privacy risk should be aware of the challenges we describe here.

SECTION 2: GENERATE SIMULATED DATA (ORIGINAL AND SYNTHETIC)

ORIGINAL DATA SET: SIMULATED DATA

- Borrowing from Reiter et al. (2014), we create a data set with $n = 1000$ and 4 dichotomous, categorical variables (i.e. dummy variables).
- The first 999 observations are a random sample from all combinations of $var1(0, 1)$, $var2(0, 1)$, $var3(0, 1)$, $var4(0, 1)$ except the last one
- The last (1000^{th}) observation is ($var1 = 1$, $var2 = 1$, $var3 = 1$, $var4 = 1$).
- This is a vulnerable record in the original data that we would want to protect using synthetic data
- The value of the simulation is that we know there is a unique record because we created it.

SYNTHETIC DATA SET

- Generate 1 synthetic data set from a CART-based SDG using the Synthpop package in R
 - We use the default settings and hyperparameter values and set a seed=1237.
- As a sensitivity test, we create 10 synthetic data sets from the original simulated data.

COMPARE ORIGINAL AND 1 SYNTHETIC DATA COPY (SEED = 1237)

Figure 1: Frequency

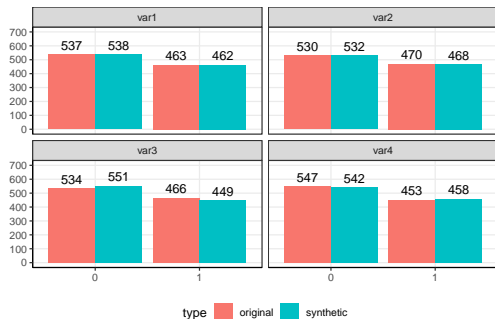
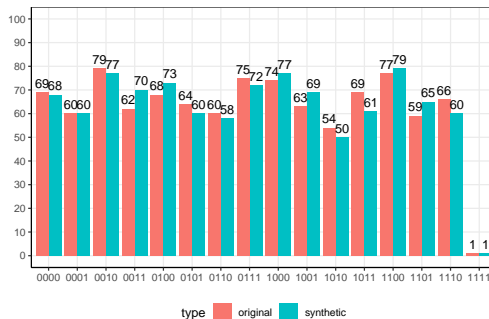
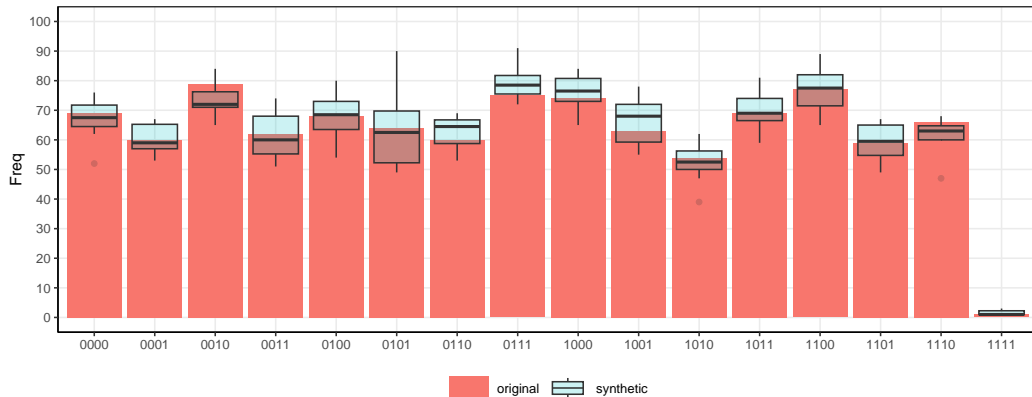


Figure 2: Histogram



COMPARE HISTOGRAM X 10 SYNTHETIC DATASETS

Figure 3: Multiple synthetic data sets does not reduce privacy risk



SUMMARY

- The problem: Synthetic data from CART models are disclosive in this simulation
- Next section: Can an attacker identify the disclosure?

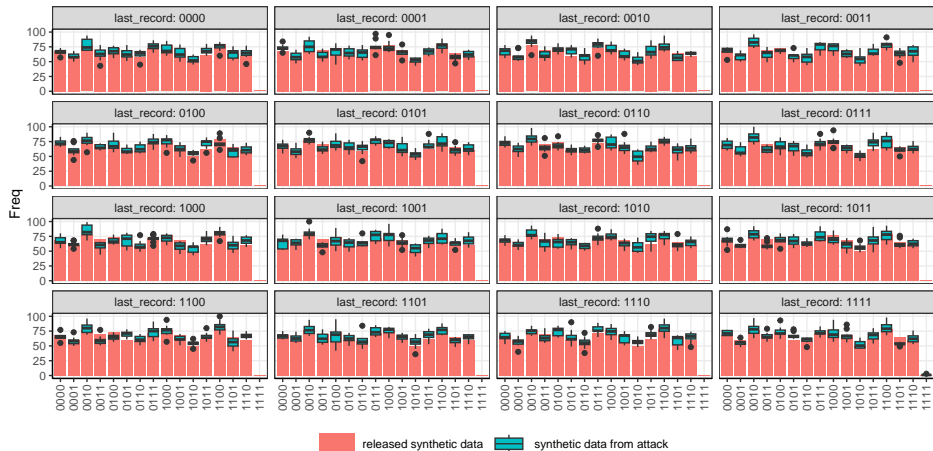
SECTION 3: THE ATTACK

DESCRIBING THE ATTACK

- We assume a 'strong' attacker similar to the attack model in differential privacy (DP).
- An attacker has the following knowledge
 - Knows the SDG model type (i.e. sequential CART).
 - Knowledge of the first 999 observations in the original data except the last one (1000th).
 - The 16 possible combinations that the last one could be.
- The attacker sees the synthetic data
- The attacker runs CART on the original data for all of the 16 different possibilities about the last record.
- Compares synthetic data from the attack to the released synthetic data
- Then they update their beliefs about what the last record could be

HISTOGRAM OF 16 WORLDS X 10 SYNTHETIC DATASETS

Figure 4



SUMMARY

- In our attack with our assumptions, the attacker can easily identify the last record
- Next section: Can we measure this disclosure risk?

SECTION 4: MEASURING DISCLOSURE RISK

THREE DISCLOSURE RISK MEASURES

- 2x Common disclosure risk measures reflect the current state of the art (Raab et al., 2025)
 - Identity risk (*repU*): the ability to identify individuals in the data from a set of known characteristics or ‘keys’ (q).
 - $q = \text{var1}(0, 1), \text{var2}(0, 1), \text{var3}(0, 1)$
 - Attribute risk (*DiSCO*): the ability to find out from the keys (q) something, not previously known or ‘target’ (t)
 - $t = \text{var4}(0, 1)$
- 1x Alternative disclosure risk measure
 - Bayesian approach (Reiter et al., 2014)
 - For example, the attacker has a prior (e.g., uniform distribution)
 - If posterior probability is close to the prior, little or no new information is revealed.
 - If posterior probability is substantially larger, the intruder has learned something about the last or unique record.

RESULTS DISCLOSURE RISK MEASURES

Table 1: x 1 synthetic data set (seed = 1237)

Data	Unique	Identity Risk (<i>repU</i>)	Attribute Risk (<i>DiSCO</i>)	Bayesian Estimate of Risk
Original	1	0.00	0.00	1.00
Synthetic	1	0.00	0.00	1.00

Table 2: x 10 synthetic data sets

Data	Unique	Identity Risk (<i>repU</i>)	Attribute Risk (<i>DiSCO</i>)	Bayesian Estimate of Risk
Original	1	0.00	0.00	1.00
Synthetic 1	1	0.00	0.00	1.00
Synthetic 2	0	0.00	6.60	0.02
Synthetic 3	1	0.00	0.00	1.00
Synthetic 4	3	0.00	0.00	1.00
Synthetic 5	2	0.00	0.00	1.00
Synthetic 6	1	0.00	0.00	1.00
Synthetic 7	3	0.00	0.00	1.00
Synthetic 8	0	0.00	6.60	0.03
Synthetic 9	1	0.00	0.00	1.00
Synthetic 10	1	0.00	0.00	1.00
Average	-	0.00	1.32	-

SUMMARY

- Common privacy measures do not capture the disclosure risk in our data
- However (and this is the point): We know there is a problem (because we created it)
- Only Bayesian approach captures disclosure risk

SECTION 5: IS THIS SCENARIO REALISTIC?

REAL WORLD DATA (SD2011)

- Replicate the approach the authors of Synthpop (Raab, 2024; Raab et al., 2024)
- Data are from Social Diagnosis 2011 (SD2011).
- Measure disclosure risk
 - 4 keys (q): `sex age region placesize`.
 - 1 target (t): `depress`
- Generate 5 synthetic copies with a 'bad' synthesizer
 - Generate 5 synthetic copies using CART with default parameters in Synthpop
 - Modify synthetic copies by setting $t = 0$, or constant for all observations in all 5 synthetic data sets.
 - Therefore, we know risk declined (because we reduced it).
- Do common disclosure risk measures ($repU$, $DiSCO$) capture this decline?
 - Why not Bayesian approach? High-dimensional, real data is too computationally complex. Only good for low-dimensional data

RESULTS

Table 3: Risk measures for depress from keys: sex, age, region, placesize (SD2011)

Data	Identity risk (<i>repU</i>)		Attribute risk (<i>DiSCO</i>)	
	Raab et al., 2024	Modified	Raab et al., 2024	Modified
Original data	48.38	48.38	53.30	53.30
Synthetic 1	14.82	14.82	8.96	14.74
Synthetic 2	14.20	14.20	9.90	14.82
Synthetic 3	15.16	15.16	10.46	14.94
Synthetic 4	14.12	14.12	9.68	14.50
Synthetic 5	14.30	14.30	8.88	14.66
Average	14.52	14.52	9.58	14.73

Note: Modified indicates that values of $\text{depress}=0$ for all records in the synthetic data

SUMMARY

- When we modify synthetic data to reduce attribute disclosure risk, *DiSCO* measure increases
- The package authors are aware of the problem that the *DiSCO* measure of attribute disclosure risk can indicate a high level of risk for a target variable where a high proportion of records have one level (Raab et al., 2025).
- This is good, but our example illustrates a more general concern: *DiSCO* may mismeasure risk by indicating it is rising, when it declined

SECTION 6: CONCLUSION

SUMMARY

- Key contribution: Common disclosure risk metrics may fail to detect or even misstate risk.
 - Suggests no risk, when we know there is a risk (simulation data)
 - Suggests risk is rising, when we know it declined (real data)
 - Bayesian approach can be a good solution, but only in low-dimensional data
- Key point: users must understand how disclosure risk measures operate.
 - Empirical disclosure risk measures always have problems
 - There is no one-size-fits-all solution.

THANK YOU

Jonathan Latner: jonathan.latner@iab.de

Reproducible code: https://github.com/jonlatner/KEM_GAN/tree/main/latner/projects/simulation