

Synthetic Data: Can We Trust Statistical Estimators?

Alexander Decruyenaere*, Heidelinde Dehaene*, Paloma Rabaey, Christiaan Polet, Johan Decruyenaere, Stijn Vansteelandt, Thomas Demeester

* Joint first authors

Background

Alongside great opportunities, great precaution should be taken regarding the possible sensitive nature of medical data and related privacy concerns.

Synthetic data are artificial data that mimic the original data in terms of statistical properties. As such, synthetic data might be able to replace the original data in statistical analysis, while **preserving the privacy** of the individual members of the original dataset.

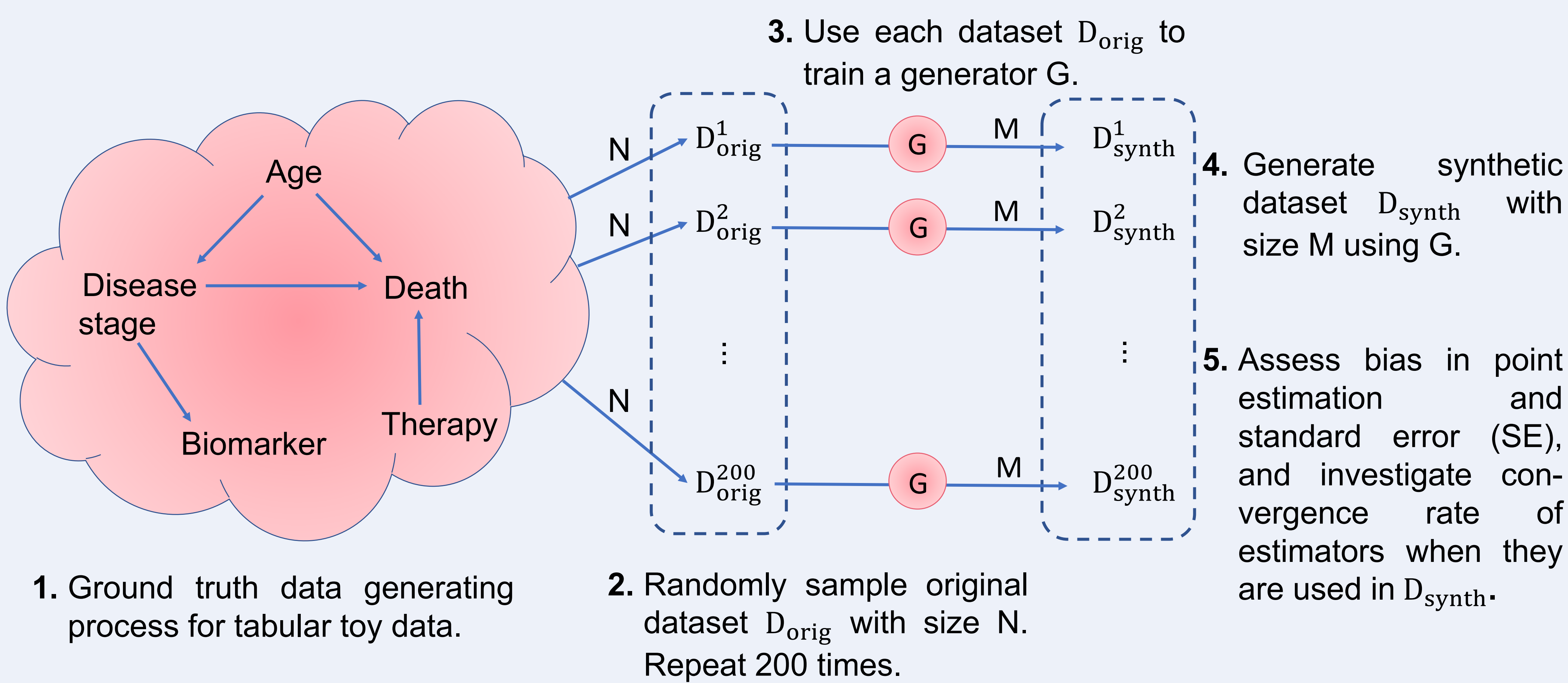


Problem statement

Can a synthetic sample be used to obtain valid estimates for a population parameter and to test hypotheses? We map two possible **pitfalls that may compromise this inferential utility of synthetic data**:

1. **Extra uncertainty** should be acknowledged since the distribution learned by the generative model is an approximation.
2. Statistical inference is typically based on **\sqrt{N} -consistency and asymptotic normality**. What is the effect of regularisation bias inherent to deep learning (DL) approaches on the default behaviour of estimators?

Experimental set-up



Desired properties of an estimator:

Standard error (SE) goes to zero when sample size increases, at rate $1/\sqrt{N}$ and bias at faster rate. True and estimated SE are according.

These properties affect inferential utility, here measured by **type 1 error rate**.

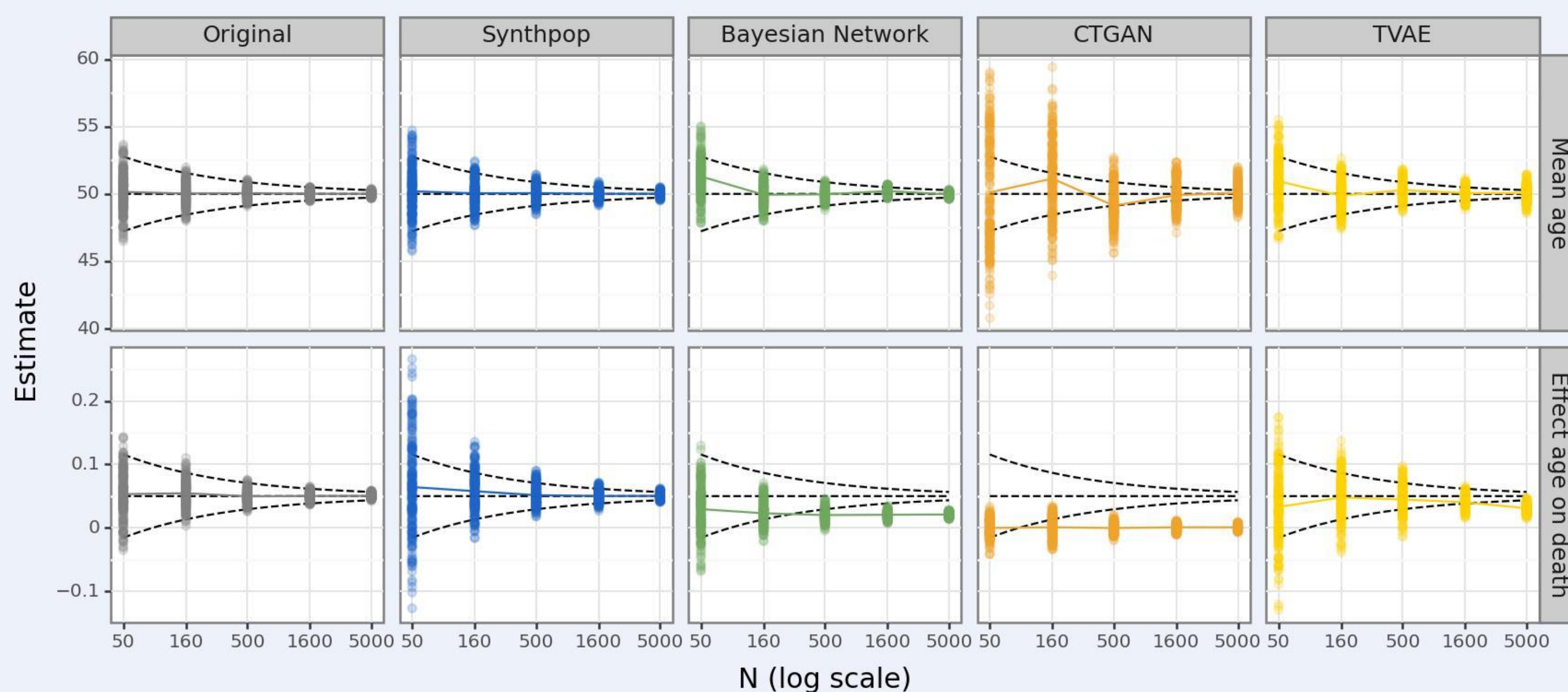
Minimal adaptation estimator SE:

$$\sigma_{\hat{\theta}, corrected} = \sigma_{\hat{\theta}, naive} \sqrt{1 + \frac{M}{N}}$$

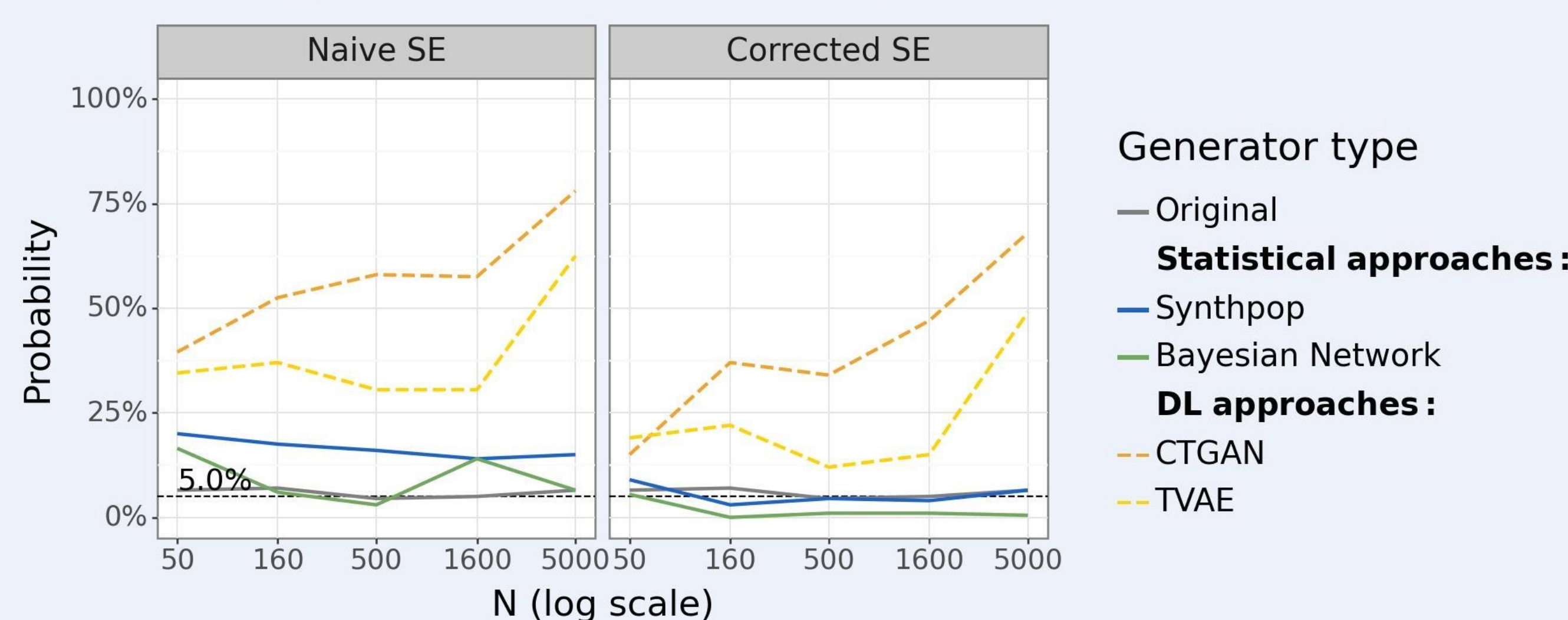
Is this valid for all estimators and generators? What about \sqrt{N} -consistency?

Results inferential utility

Finite sample behaviour of estimator



Type 1 error rate for mean age



Main results:

1. Generative model **misspecification** introduces **bias**.
2. True **SEs are larger** for D_{synth} than for D_{orig} and extra variability varies over generative models.
3. Therefore, **naive estimation of SE leads to its underestimation**.
4. **Convergence rate** of the SE of various estimators **differs** between statistical and DL approaches. For the statistical approaches, estimators remain roughly \sqrt{N} -consistent. In the **DL approaches**, estimators converge **slower**.
5. **Naive analyses** lead to **inflation of type 1 error rate** (compromising inferential utility).
6. **Adaptation** for the SE **controls type 1 error rate only in statistical and not in DL** approaches (due to slower-than- \sqrt{N} -convergence).

Conclusion:

Before publishing synthetic data, it is essential to develop statistical inference tools for such data.

