

Balancing the privacy-utility trade-off for synthetic time-to-event data

Generated with sequential regressions in Stata

Sigrid Leithe, Bjørn Møller, Bjarte Aagnes, Yngvar Nilssen, Tor Åge Myklebust

Synthetic data

Artificially generated data from a model that is trained to reproduce characteristics of the original data.

European Data Protection Supervisor (EDPS)



Public release of example data.



Publish data alongside journal articles to enable reproducibility.



IT development and testing without exposing sensitive information.



Education and training.



Methods and algorithm development.

Motivating dataset

British Journal of Cancer

ARTICLE

Improving communication of cancer survival statistics—feasibility of implementing model-based algorithms in routine publications

Tor Åge Myklebust^{1,2}, Bjarte Aagnes¹, Yngvar Nilssen¹, Mark Rutherford^{3,4}, Paul C. Lambert^{3,5}, Therese M. L. Andersson^{1,5}, Anna L. V. Johansson^{1,5}, Paul W. Dickman⁵ and Bjørn Møller¹

© The Author(s), under exclusive license to Springer Nature Limited

BACKGROUND: Routine inform about prognosis and targeting different a wider range of survival **METHODS:** We used data estimating flexible param expectancy across many **RESULTS:** For 21 of 23 estimates of all desired **DISCUSSION:** It may be modeling techniques. W estimates across a range *British Journal of Cancer*

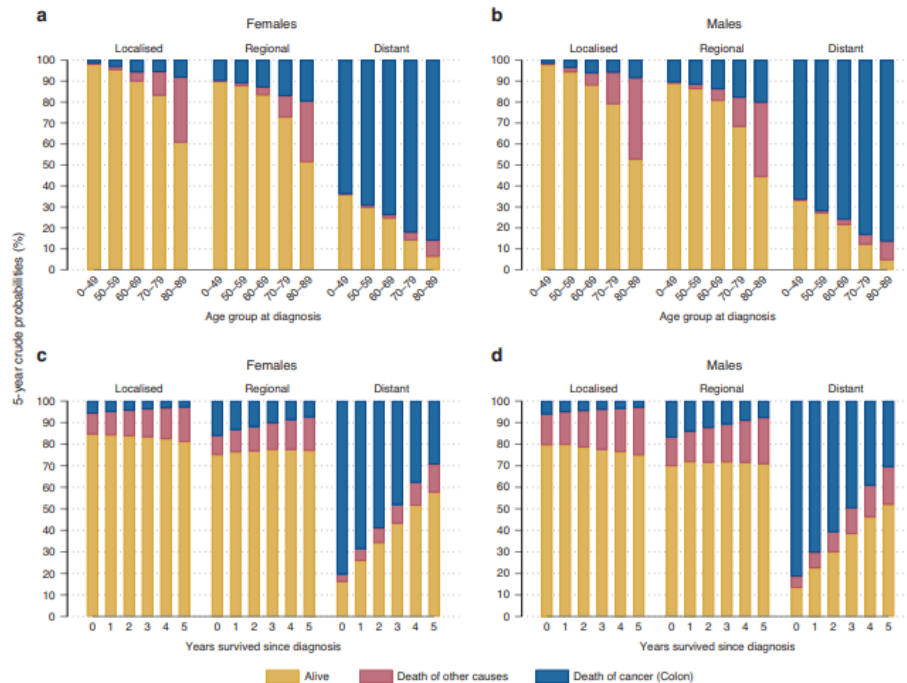


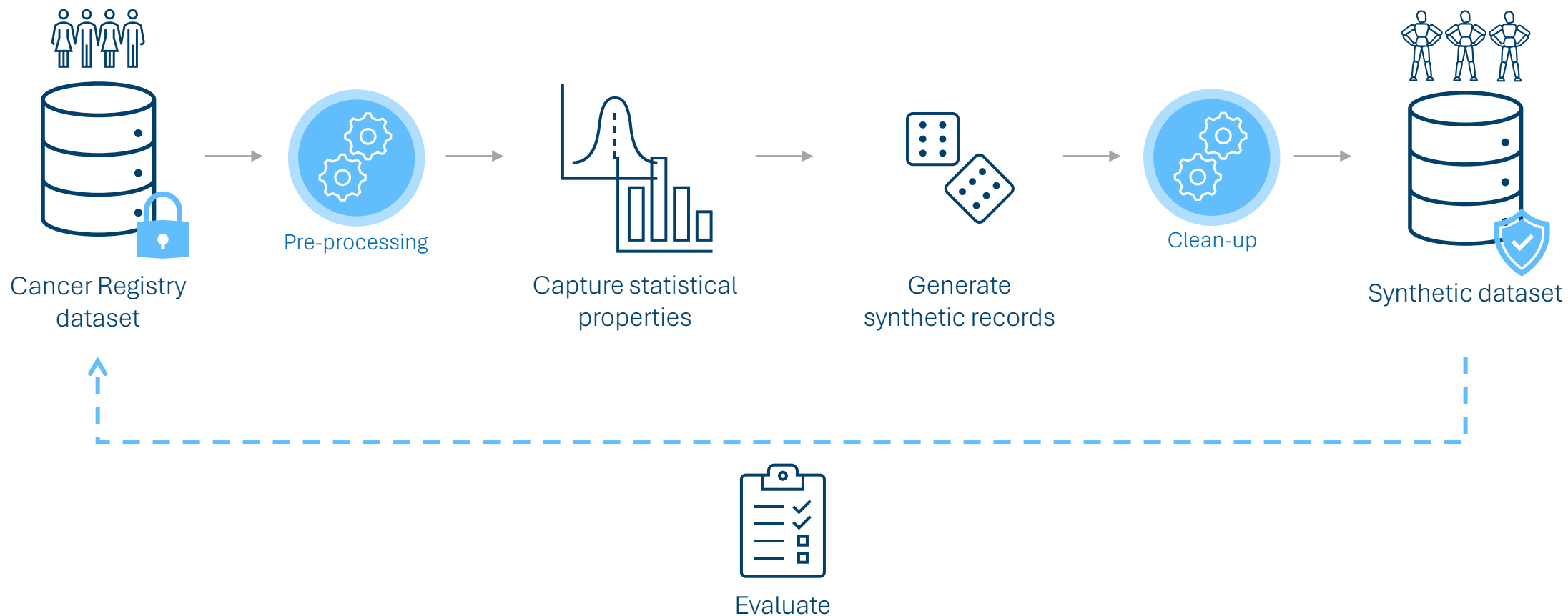
Fig. 1 The figure shows estimated 5-year crude probabilities of dying from colon cancer, dying from other causes and being alive. Results are stratified by sex, stage and age group at diagnosis (a, b) and by sex, stage and number of years survived since diagnosis (c, d). Crude probabilities are also referred to as cumulative incidence or "real-world probabilities".

www.nature.com/bj

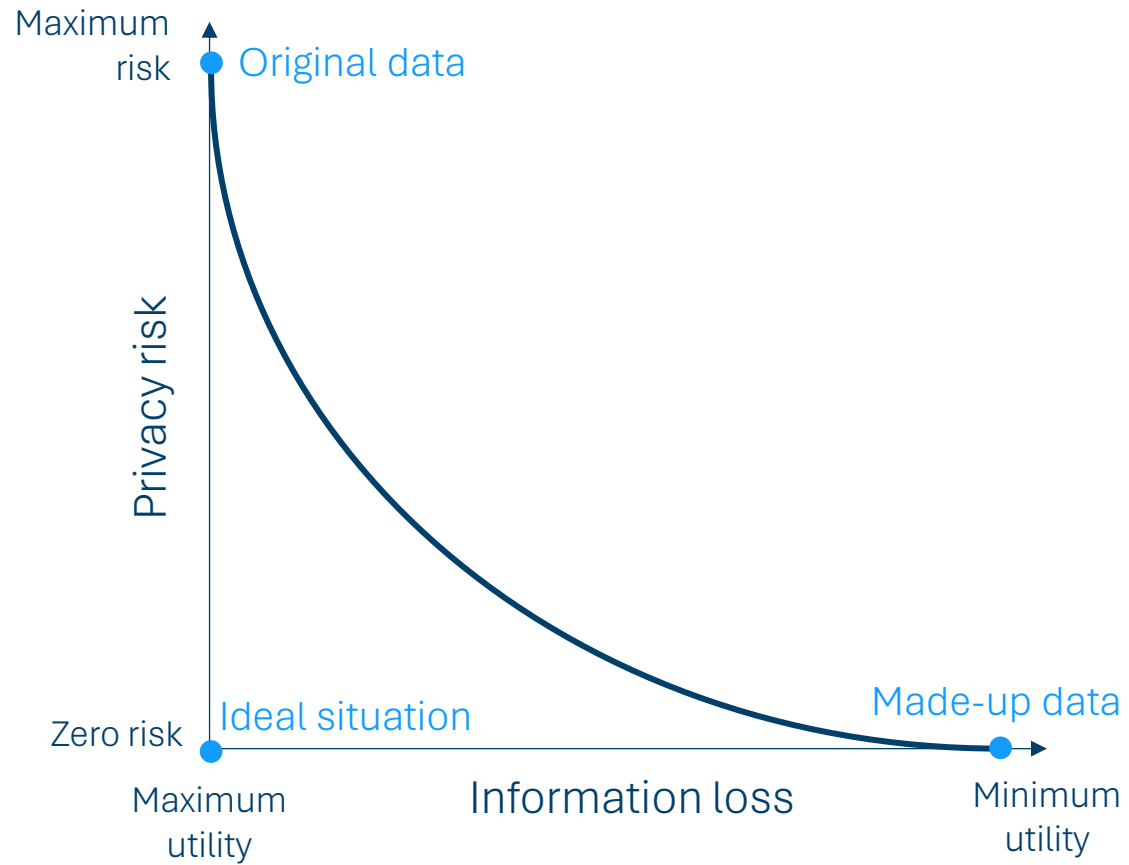
Check for updates

GitHub repository interface for **CancerRegistryOfNorway / cancer-survival-measures**. The repository is public and contains files for the BJC publication version 1.0, including folders for `ado`, `data`, `dofiles`, `results`, and `readme.md`. The README section is visible, showing the title **Supplementary materials** and the **Purpose** section, which states: "This repo is for code used in: Tor Åge Myklebust, Bjarte Aagnes, Yngvar Nilssen, Mark Rutherford, Paul C. Lambert, Therese M. L. Andersson, Anna L. V. Johansson, Paul W. Dickman & Bjørn Møller *Improving communication of cancer survival statistics—feasibility of implementing model-based algorithms in routine publications*. Br J Cancer (2023). <https://doi.org/10.1038/s41416-023-02360-5>".

Synthetic data generation



Privacy-utility trade-off



Synthetic data generators

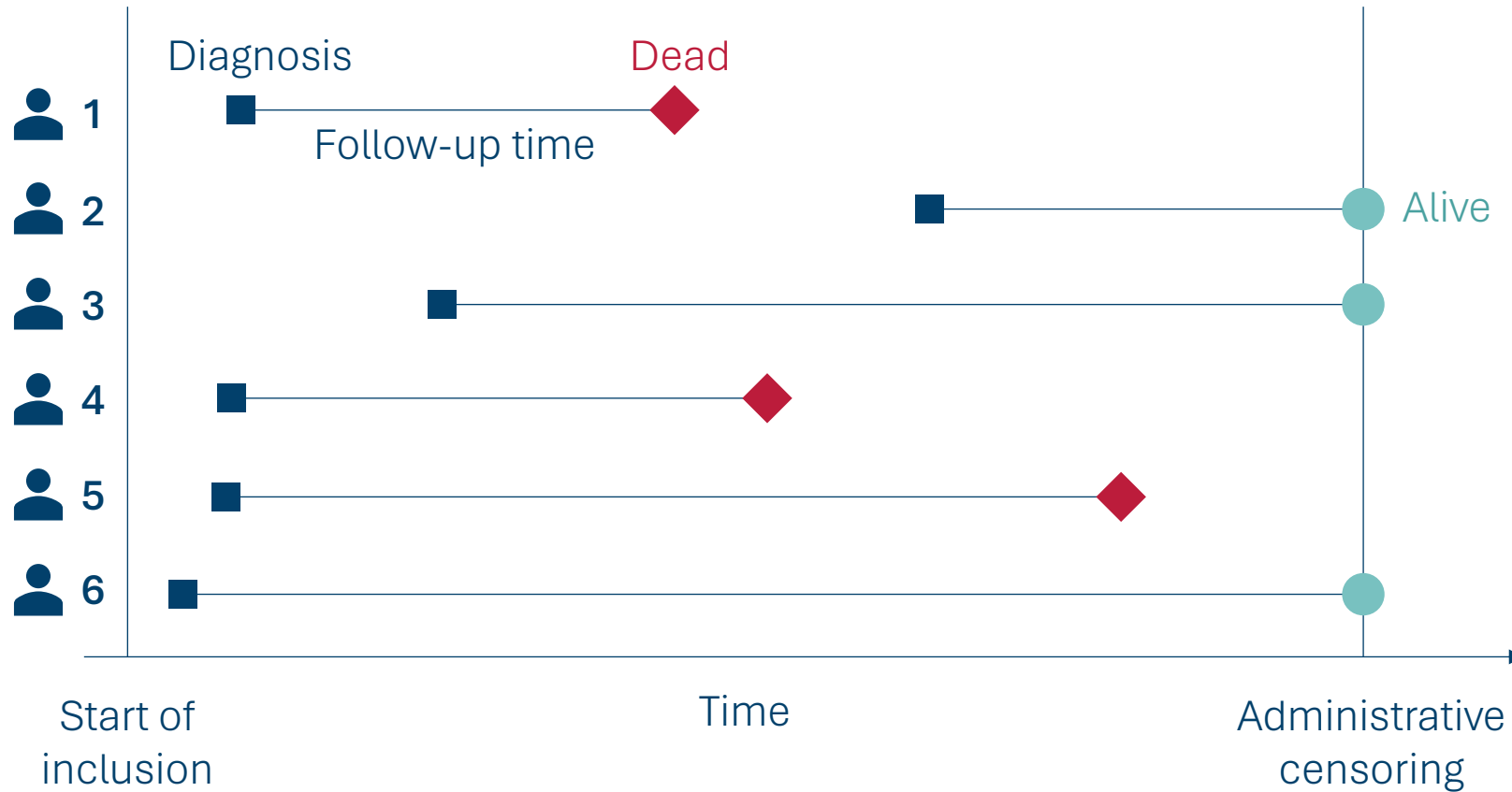
Statistical methods

- Imputation-based methods
- Bayesian networks
- Copula-based methods

Machine learning methods

- Generative Adversarial Networks (GAN)
- Variational Auto Encoders (VAE)
- Transformer-based models

Time-to-event data



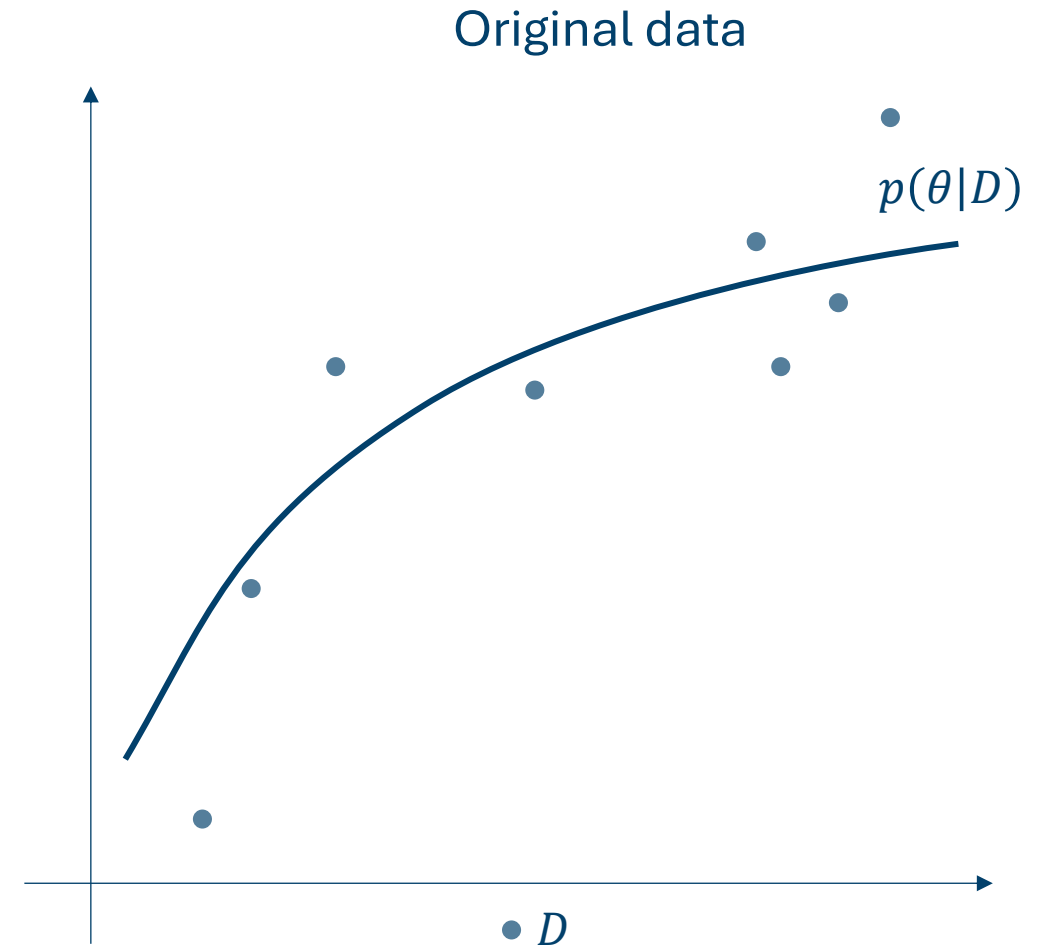
Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Original data: Colon cancer
52 162 patients
Diagnosed 2002-2021
Administrative censoring 31.12.2021

Capture statistical properties

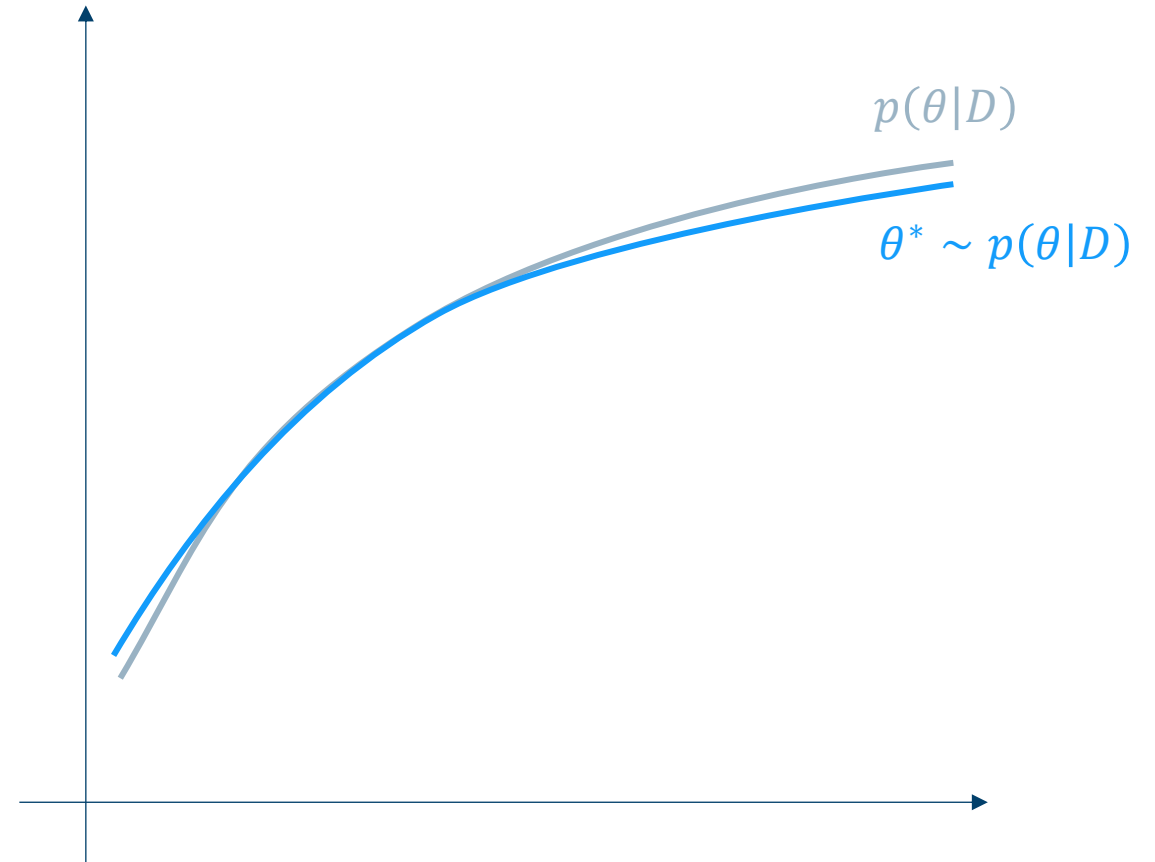
Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead



Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

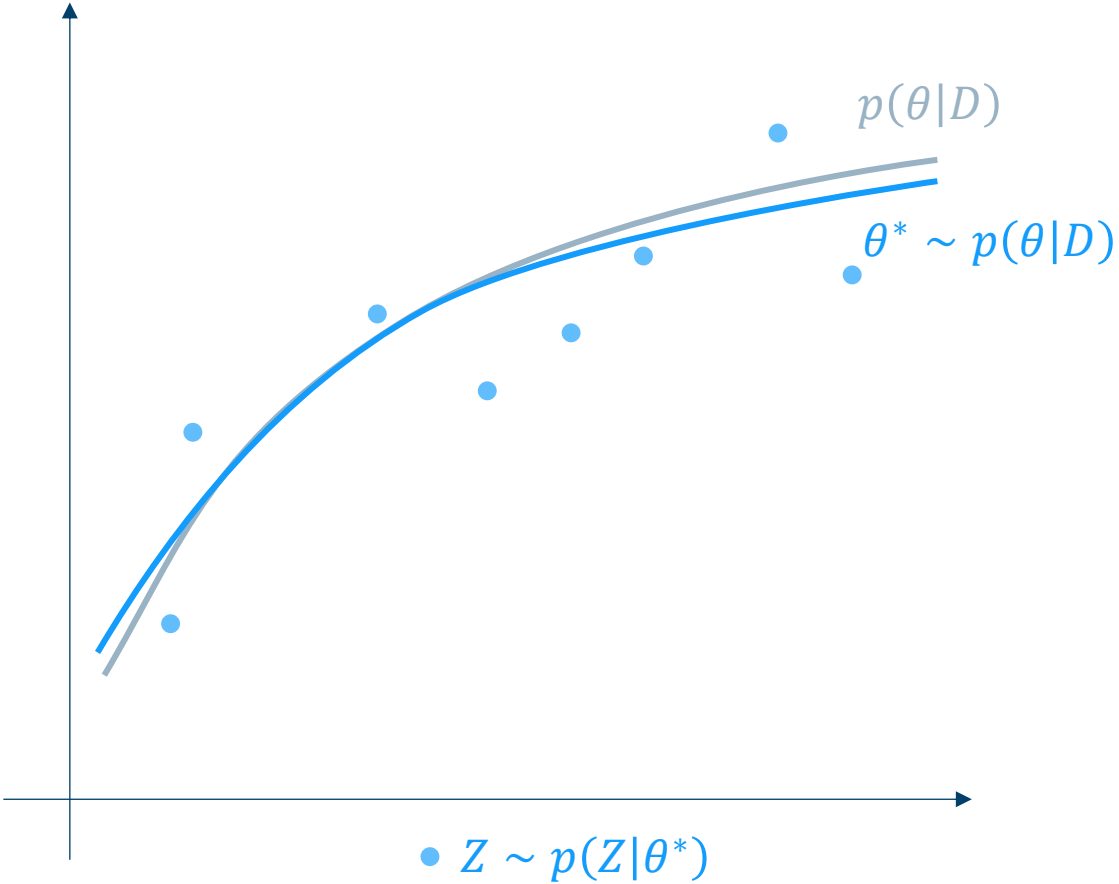
Statistical model



Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Synthetic data



Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

$$Age_i$$

Synthetic records:

54					
----	--	--	--	--	--

Smith, A., Lambert, P. C., & Rutherford, M. J. (2022). Generating high-fidelity synthetic time-to-event datasets to improve data transparency and accessibility. *BMC Med Res Methodol.*

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

Synthetic records:

54					
----	--	--	--	--	--

$P(\text{Female} \mid \text{Age} = 54) = 0.43$

$P(\text{Male} \mid \text{Age} = 54) = 0.57$

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

Synthetic records:

54	Male				
----	------	--	--	--	--

$P(\text{Female} \mid \text{Age} = 54) = 0.43$

$P(\text{Male} \mid \text{Age} = 54) = 0.57$

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

Synthetic records:

54	Male				
----	------	--	--	--	--

$P(\text{Local} \mid \text{Age} = 54, \text{Male}) = 0.45$

$P(\text{Regional} \mid \text{Age} = 54, \text{Male}) = 0.25$

$P(\text{Distant} \mid \text{Age} = 54, \text{Male}) = 0.19$

$P(\text{Unknown} \mid \text{Age} = 54, \text{Male}) = 0.11$

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

Synthetic records:

54	Male	Regional			
----	------	----------	--	--	--

$P(\text{Local} \mid \text{Age} = 54, \text{Male}) = 0.45$

$P(\text{Regional} \mid \text{Age} = 54, \text{Male}) = 0.25$

$P(\text{Distant} \mid \text{Age} = 54, \text{Male}) = 0.19$

$P(\text{Unknown} \mid \text{Age} = 54, \text{Male}) = 0.11$

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

$Year_i \sim Age_i, Sex_i, Stage_i$

Synthetic records:

54	Male	Regional			
----	------	----------	--	--	--

$P(2002 \mid Age = 54, Male, Regional) = 0.04$

$P(2003 \mid Age = 54, Male, Regional) = 0.06$

$P(2004 \mid Age = 54, Male, Regional) = 0.05$

$P(2005 \mid Age = 54, Male, Regional) = 0.07$

...

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

$Year_i \sim Age_i, Sex_i, Stage_i$

Synthetic records:

54	Male	Regional	06.03.2005		
----	------	----------	------------	--	--

$P(2002 \mid Age = 54, Male, Regional) = 0.04$

$P(2003 \mid Age = 54, Male, Regional) = 0.06$

$P(2004 \mid Age = 54, Male, Regional) = 0.05$

$P(2005 \mid Age = 54, Male, Regional) = 0.07$

...

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

$Year_i \sim Age_i, Sex_i, Stage_i$

$t_i^* \sim Age_i, Sex_i, Stage_i, Year_i$

$t_i = \min(t_i^*, C_i)$

$Status_i = \begin{cases} \text{Dead,} & t_i^* \leq C_i \\ \text{Alive,} & t_i^* > C_i \end{cases}$

Synthetic records:

54	Male	Regional	06.03.2005		
----	------	----------	------------	--	--

Capture statistical properties

Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
77	Female	Distant	05.03.2002	5.12	Dead
50	Male	Local	21.08.2011	9.51	Alive
63	Female	Unknown	30.12.2018	3.02	Alive
46	Male	Regional	02.07.2005	3.45	Dead
60	Male	Local	09.10.2016	13.83	Dead

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

$Year_i \sim Age_i, Sex_i, Stage_i$

$t_i^* \sim Age_i, Sex_i, Stage_i, Year_i$

$t_i = \min(t_i^*, C_i)$

$Status_i = \begin{cases} \text{Dead}, & t_i^* \leq C_i \\ \text{Alive}, & t_i^* > C_i \end{cases}$

Synthetic records:

54	Male	Regional	06.03.2005	4.31	Dead
----	------	----------	------------	------	------

Smith, A., Lambert, P. C., & Rutherford, M. J. (2022). Generating high-fidelity synthetic time-to-event datasets to improve data transparency and accessibility. *BMC Med Res Methodol*.

Experimental design

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

$Year_i \sim Age_i, Sex_i, Stage_i$

$t_i^* \sim Age_i, Sex_i, Stage_i, Year_i$

In all models: Main effects of Age, stage, sex and year(3-year periods).

Model #	Interactions	Time varying coefficients (TVCs)	Degrees of freedom		
			Age	Baseline excess hazard	TVCs
1	None	Age, stage	4	5	2
2	Stage \times Age	Age, stage	4	5	2
3	Stage \times Age	Age, stage, sex	4	5	3
4	Stage \times Age, Stage \times Sex, Age \times Sex	Age, stage, sex	4	5	3
5	Stage \times Age \times Sex	Age, stage, sex	4	5	3
6	Stage \times Age \times Sex	Age, stage, sex	6	8	6

Experimental design

Sequential regressions:

Age_i

$Sex_i \sim Age_i$

$Stage_i \sim Age_i, Sex_i$

$Year_i \sim Age_i, Sex_i, Stage_i$

$t_i^* \sim Age_i, Sex_i, Stage_i, Year_i$

Model #	Interactions	Time varying coefficients (TVCs)	Degrees of freedom		
			Age	Baseline excess hazard	TVCs
Independent marginals (lower reference model)					
1	None	Age, stage	4	5	2
2	Stage \times Age	Age, stage	4	5	2
3	Stage \times Age	Age, stage, sex	4	5	3
4	Stage \times Age, Stage \times Sex, Age \times Sex	Age, stage, sex	4	5	3
5	Stage \times Age \times Sex	Age, stage, sex	4	5	3
6	Stage \times Age \times Sex	Age, stage, sex	6	8	6
Resampling (upper reference model)					

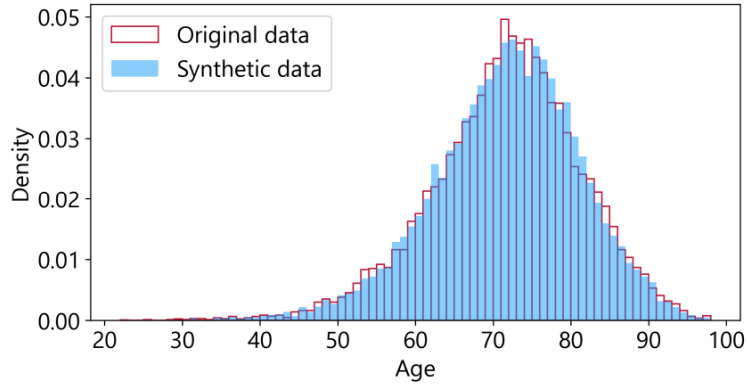
50 datasets from each model

Synthetic data evaluation

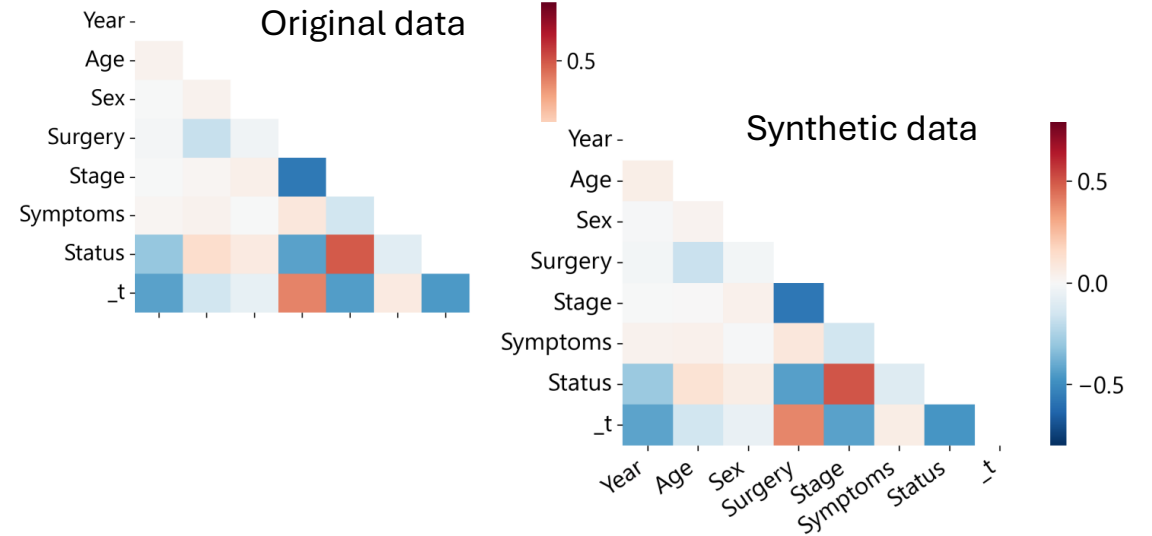


Synthetic data utility

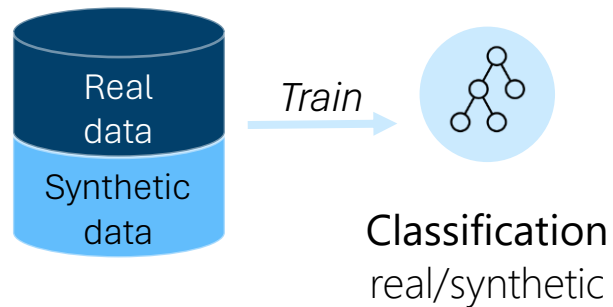
Univariate



Bivariate

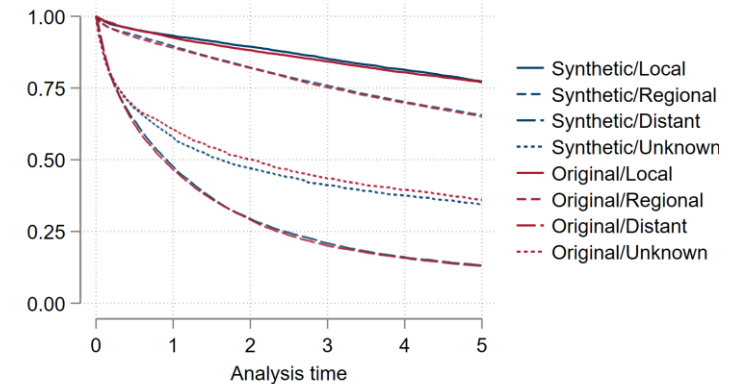


Multivariate

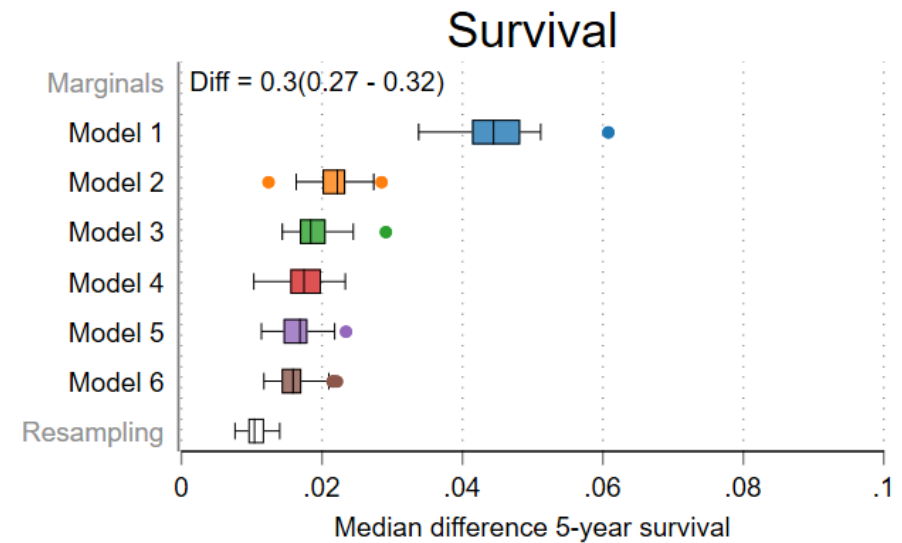
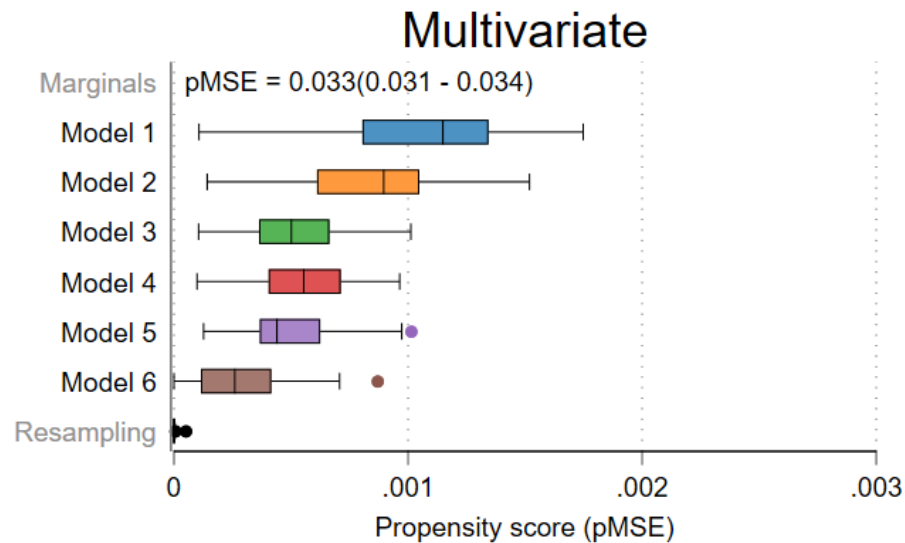
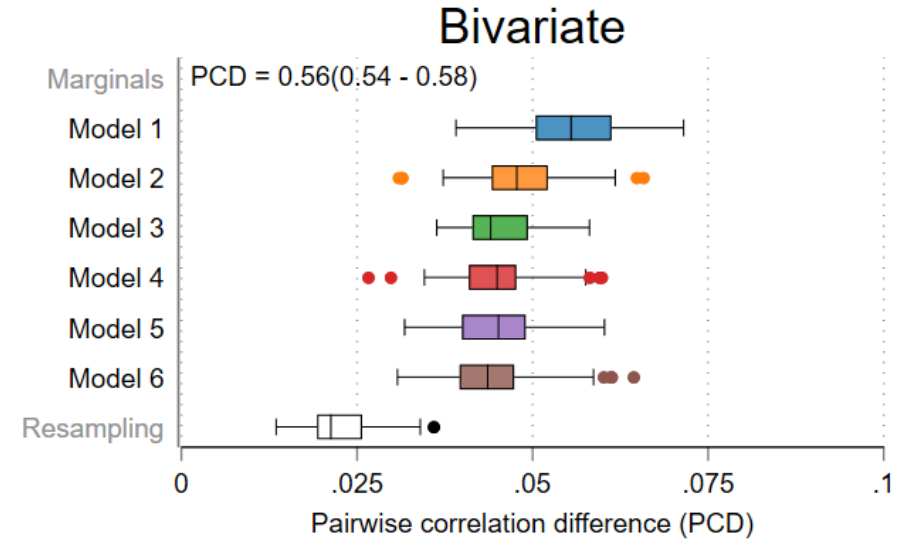
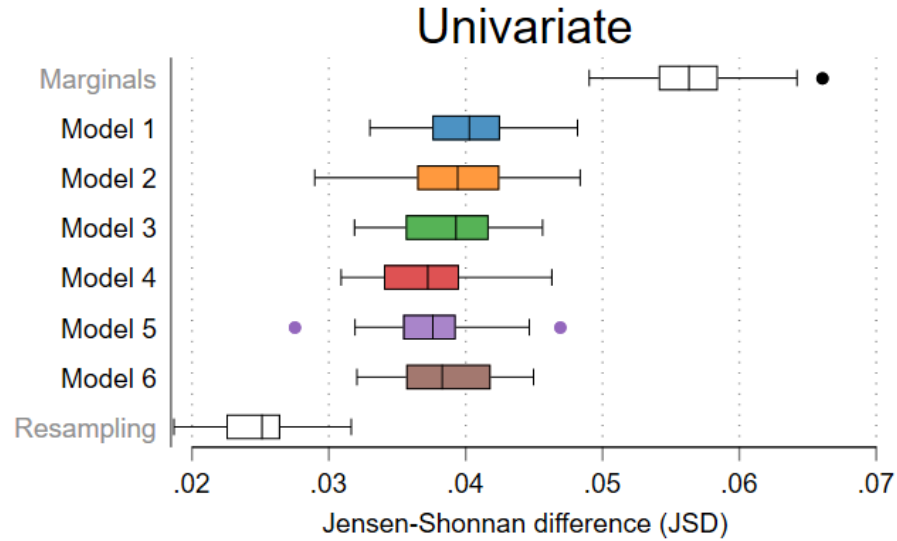


Survival

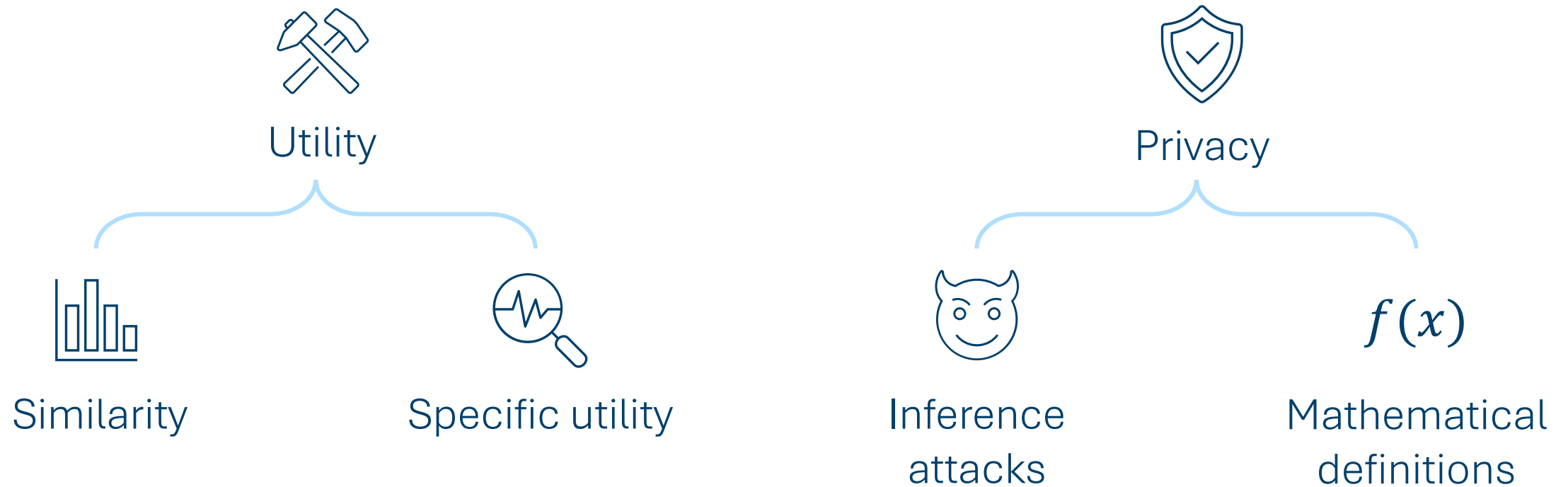
Kaplan–Meier survival estimates



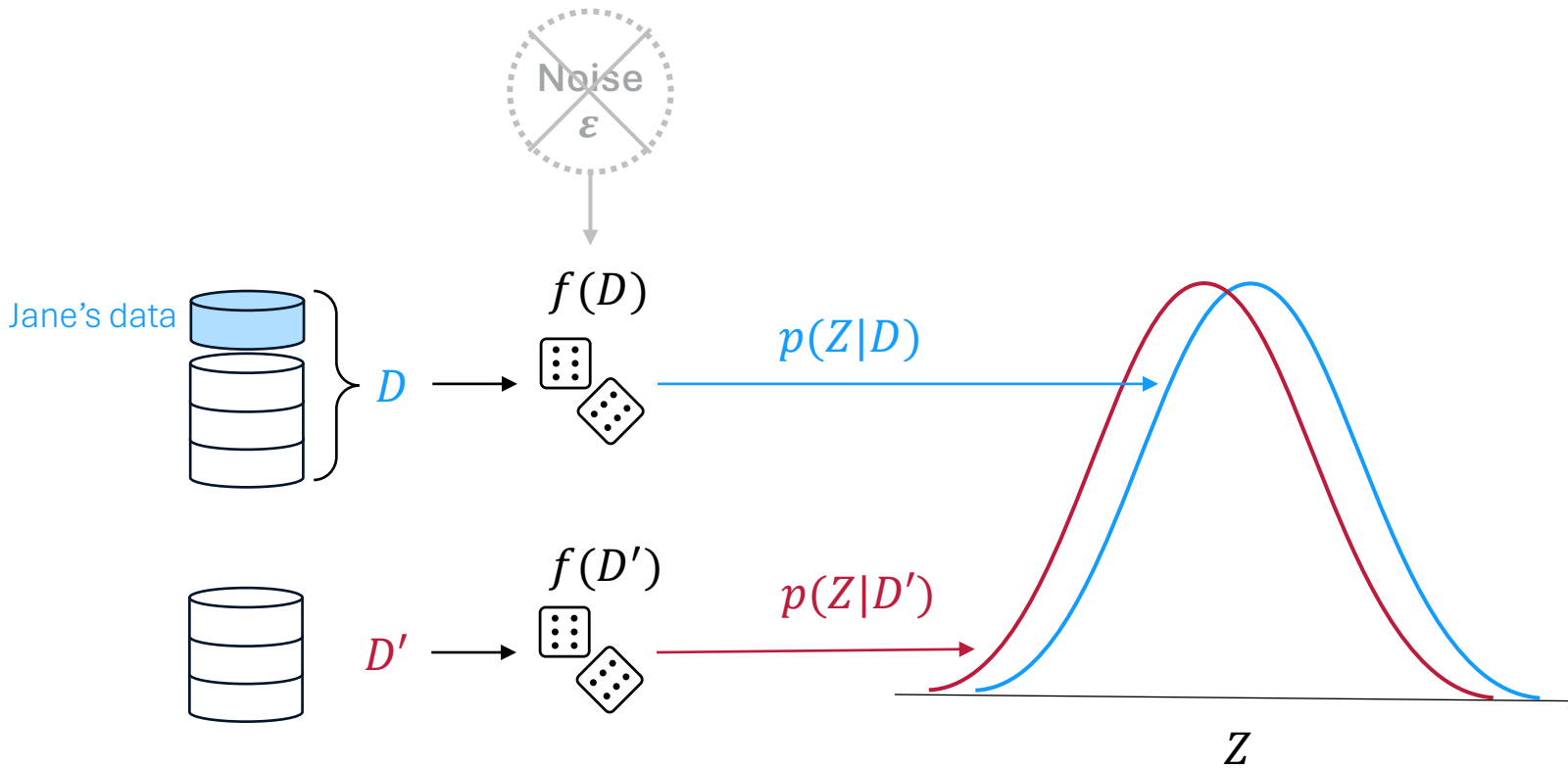
Synthetic data utility



Synthetic data evaluation



Differential privacy



For all $Z \in \text{Range}(f)$:

$$\frac{p(Z|D)}{p(Z|D')} \leq e^\epsilon$$

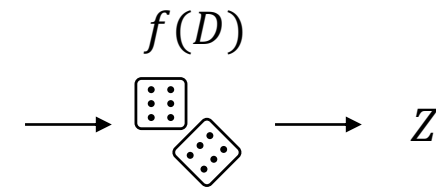
$$|\log(p(Z|D)) - \log(p(Z|D'))| \leq \epsilon$$

$$\epsilon_{emp.} = \max |\log(p(Z|D)) - \log(p(Z|D'))|$$

$$\epsilon_{emp.} \ll \epsilon$$

Differential privacy auditing

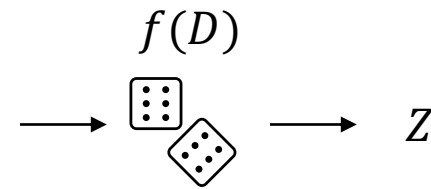
i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
Jane 1	77	Female	Distant	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead



Reiter, J. P., Wang, Q., & Zhang, B. (2014). Bayesian Estimation of Disclosure Risks for Multiply Imputed, Synthetic Data. *Journal of Privacy and Confidentiality*, 6(1).

Differential privacy auditing

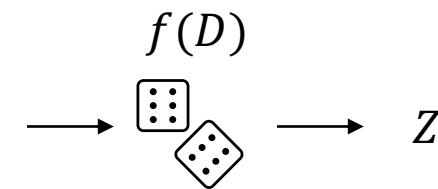
i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
Jane 1	77	Female	?	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead



Reiter, J. P., Wang, Q., & Zhang, B. (2014). Bayesian Estimation of Disclosure Risks for Multiply Imputed, Synthetic Data. *Journal of Privacy and Confidentiality*, 6(1).

Differential privacy auditing

i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
Jane 1	77	Female	?	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead



Estimate:

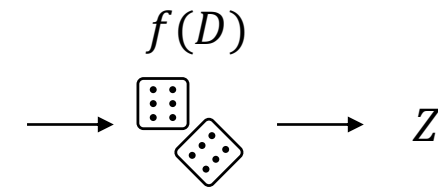
- $p(Z|D, \text{stage}_1 = \text{Localised})$
- $p(Z|D, \text{stage}_1 = \text{Regional})$
- $p(Z|D)$ ($\text{stage}_1 = \text{Distant}$)
- $p(Z|D, \text{stage}_1 = \text{Unknown})$

Information leakage:

$$\max_y | \log(p(Z|D, \text{stage}_1 = y)) - \log(p(Z|D)) |$$

Differential privacy auditing

i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
1	77	Female	Distant	05.03.2002	5.12	Dead
2	50	Male	?	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead

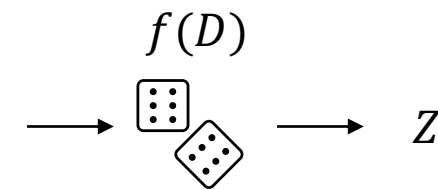


Information leakage:

$$\max_y \left| \log(p(Z|D, \text{stage}_2 = y)) - \log(p(Z|D)) \right|$$

Differential privacy auditing

i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
1	77	Female	Distant	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	?	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead

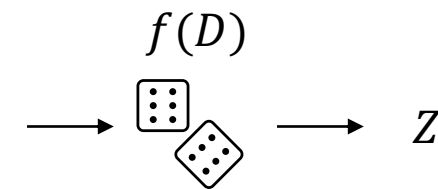


Information leakage:

$$\max_y \left| \log(p(Z|D, \text{stage}_3 = y)) - \log(p(Z|D)) \right|$$

Differential privacy auditing

i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
1	77	Female	Distant	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	?	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead

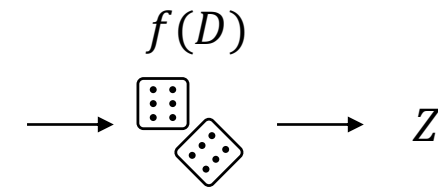


Information leakage:

$$\max_y | \log(p(Z|D, \text{stage}_4 = y)) - \log(p(Z|D)) |$$

Differential privacy auditing

i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
1	77	Female	Distant	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	?	09.10.2016	13.83	Dead

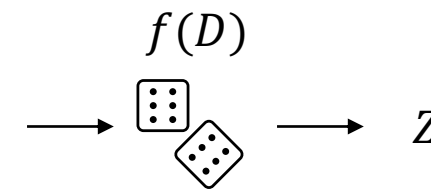


Information leakage:

$$\max_y \left| \log(p(Z|D, \text{stage}_5 = y)) - \log(p(Z|D)) \right|$$

Differential privacy auditing

i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
1	77	Female	Distant	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead



Differential privacy audit:

$$\epsilon_{emp.} = \max_{i,y} | \log(p(Z|D, \text{stage}_i = y)) - \log(p(Z|D)) |$$

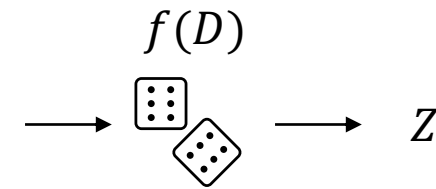
Differential privacy auditing

i	Age	Sex	Stage	Date of diagnosis	Follow-up time	Status
1	77	Female	Distant	05.03.2002	5.12	Dead
2	50	Male	Local	21.08.2011	9.51	Alive
3	63	Female	Unknown	30.12.2018	3.02	Alive
4	46	Male	Regional	02.07.2005	3.45	Dead
5	60	Male	Local	09.10.2016	13.83	Dead
\vdots						

$52\,162 \times 3 = 156\,486$

Differential privacy audit:

$$\epsilon_{emp.} = \max_{i,y} | \log(p(Z|D, \text{stage}_i = y)) - \log(p(Z|D)) |$$

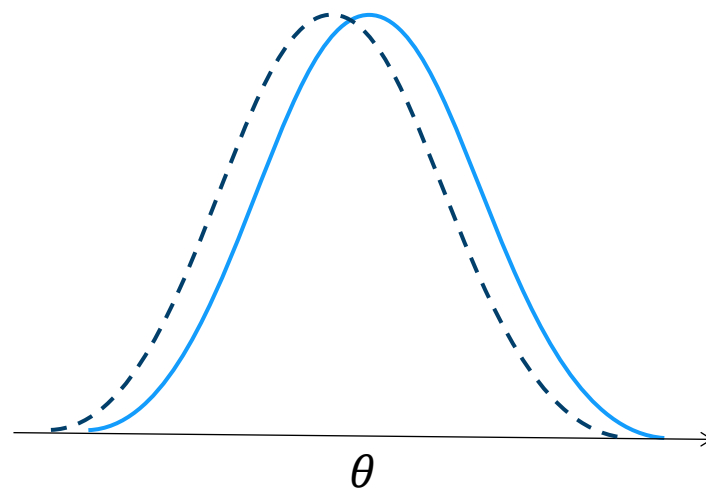


Differential privacy auditing

$$\epsilon_{emp.} = \max_{i,y} | \log(p(Z|D, \text{stage}_i = y)) - \log(p(Z|D)) |$$

$D_{\text{stage}_i = y} \approx D$

$p(\theta|D_{\text{stage}_i = y}) \approx p(\theta|D)$



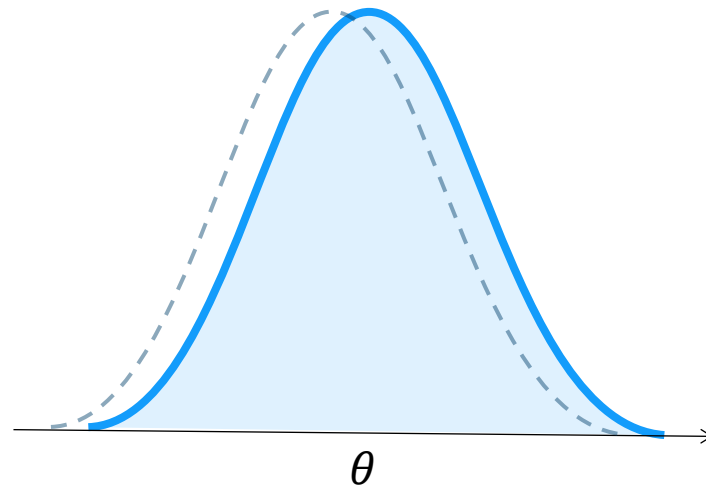
Hu, J., Reiter, J. P., & Wang, Q. (2014). Disclosure risk evaluation for fully synthetic categorical data. *International conference on privacy in statistical databases*.

Differential privacy auditing

$$\epsilon_{emp.} = \max_{i,y} | \log(p(Z|D, \text{stage}_i = y)) - \log(p(Z|D)) |$$

$D_{\text{stage}_i = y} \approx D$

$p(\theta|D_{\text{stage}_i = y}) \approx p(\theta|D)$



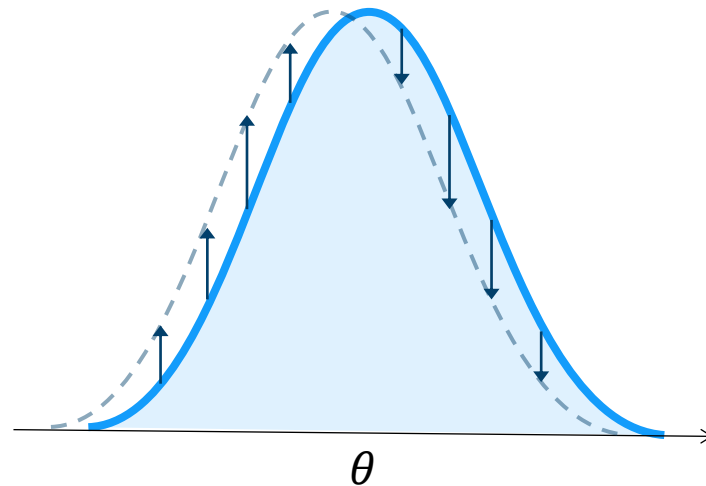
Hu, J., Reiter, J. P., & Wang, Q. (2014). Disclosure risk evaluation for fully synthetic categorical data. *International conference on privacy in statistical databases*.

Differential privacy auditing

$$\epsilon_{emp.} = \max_{i,y} | \log(p(Z|D, \text{stage}_i = y)) - \log(p(Z|D)) |$$

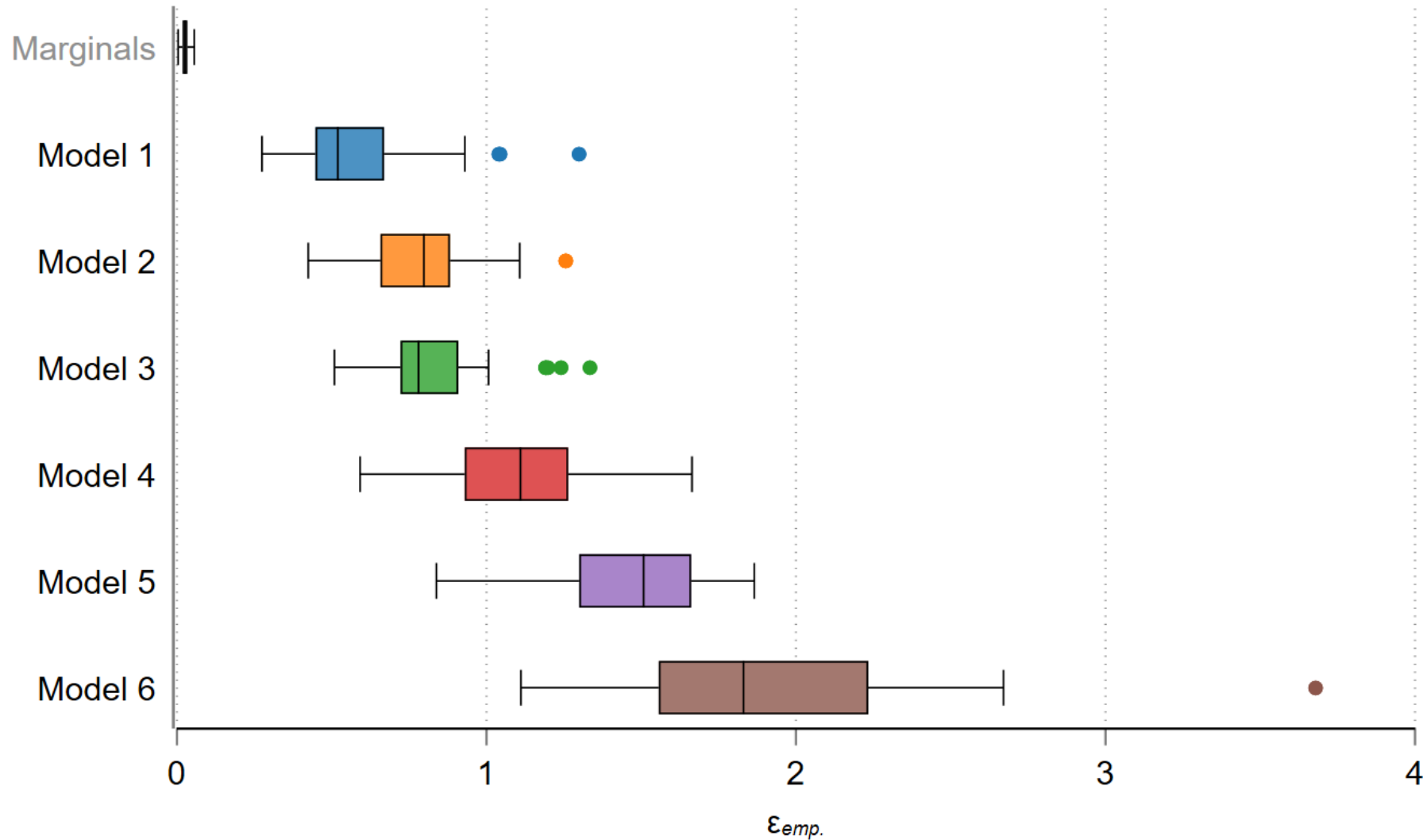
$$D_{\text{stage}_i = y} \approx D$$

$$p(\theta|D_{\text{stage}_i = y}) \approx p(\theta|D)$$

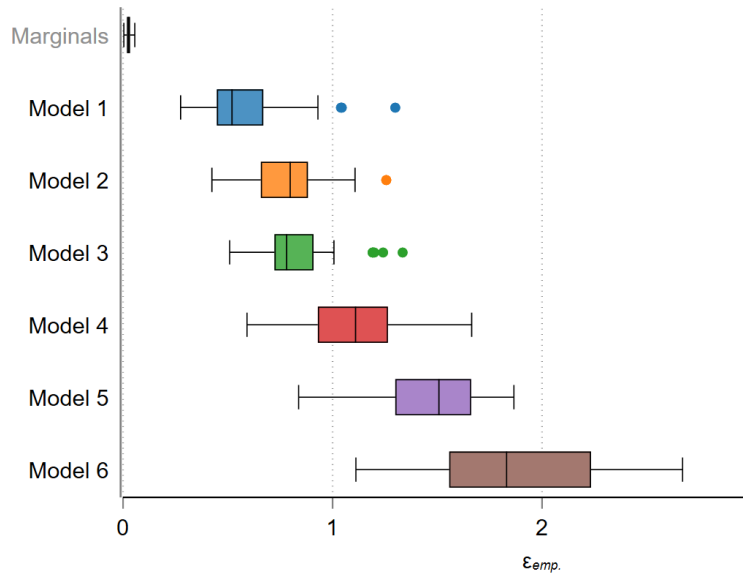


Hu, J., Reiter, J. P., & Wang, Q. (2014). Disclosure risk evaluation for fully synthetic categorical data. *International conference on privacy in statistical databases*.

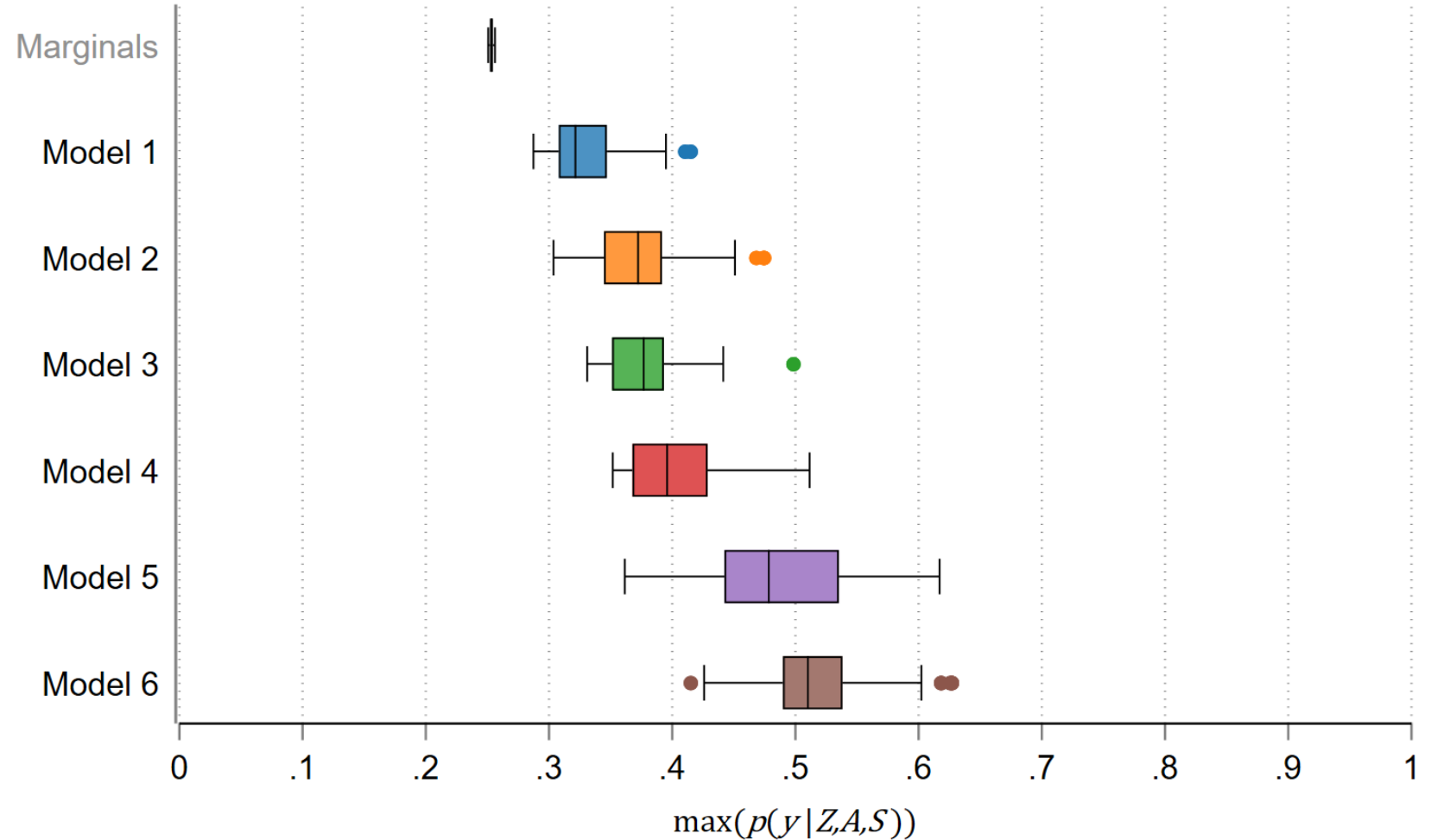
Differential privacy auditing

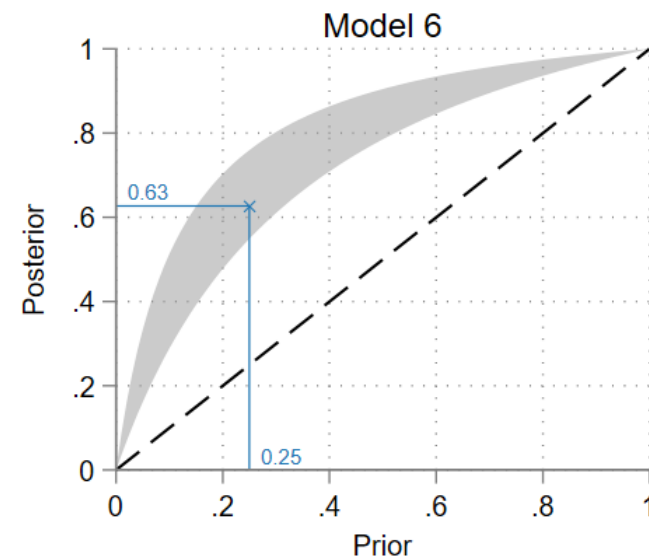
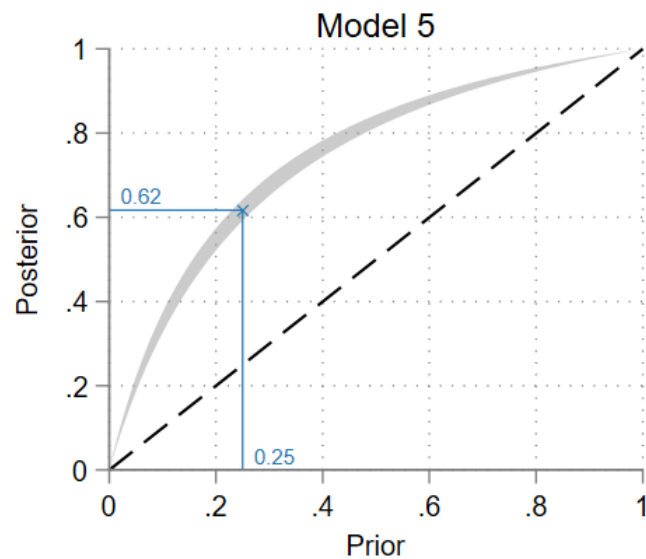
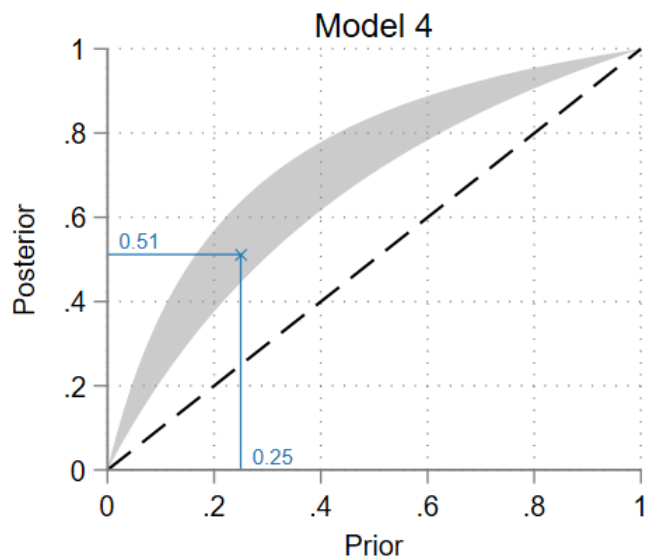
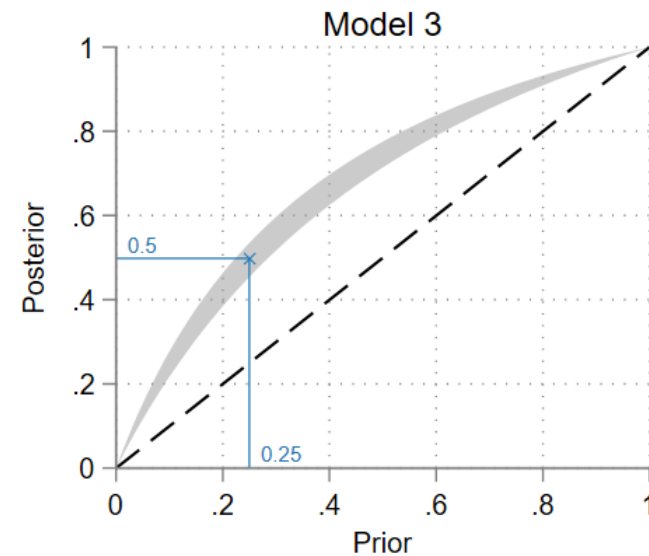
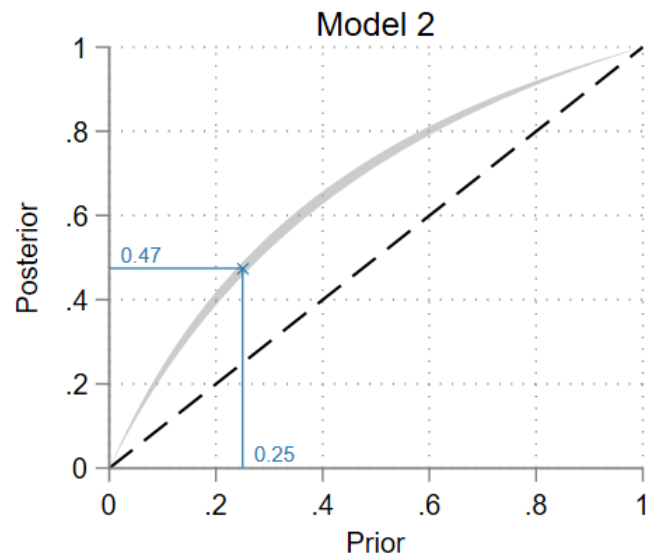
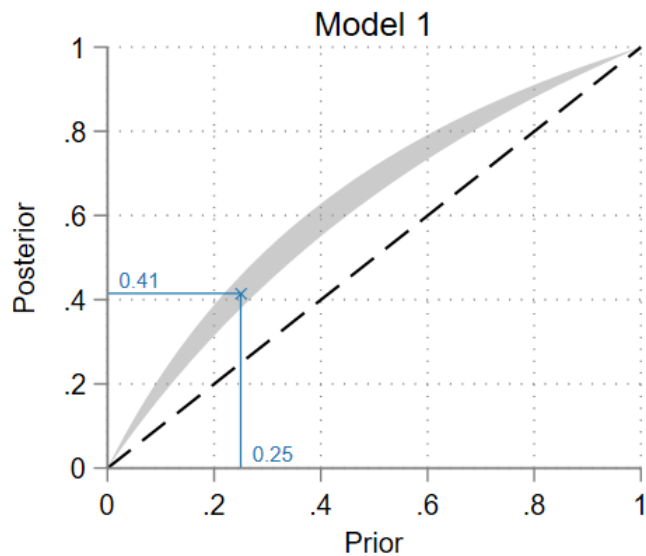


Differential privacy auditing

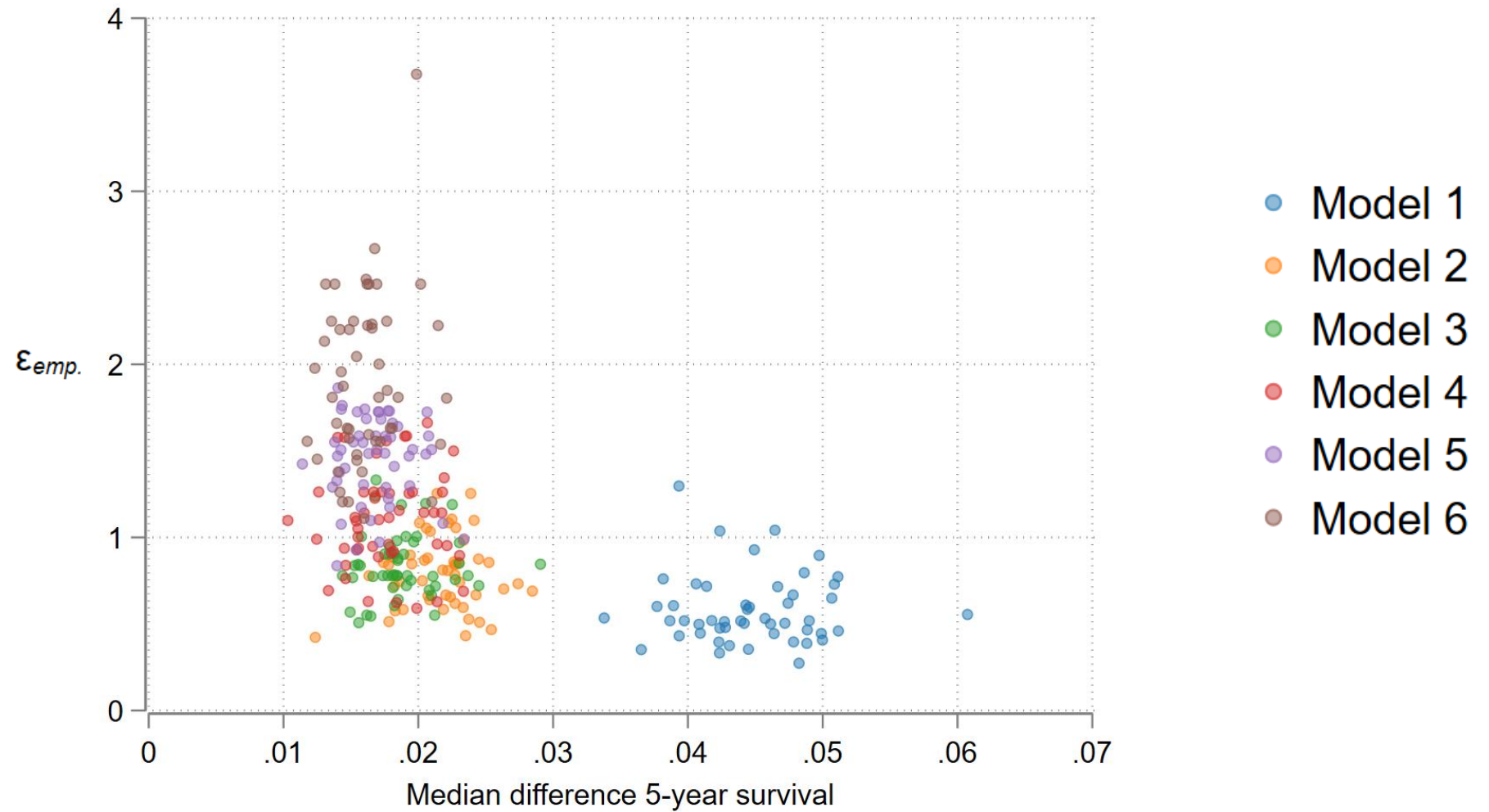


Maximum posterior probability

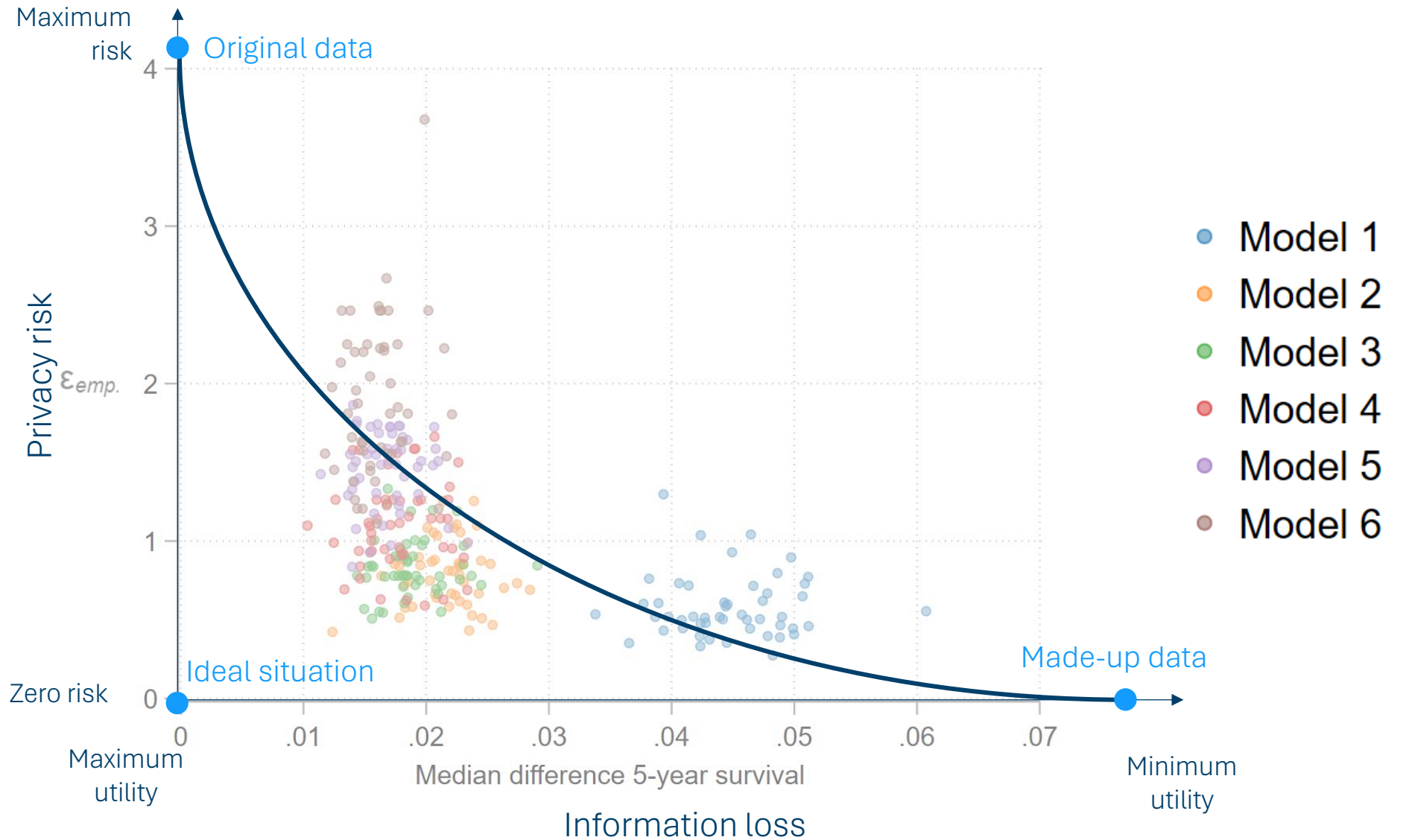




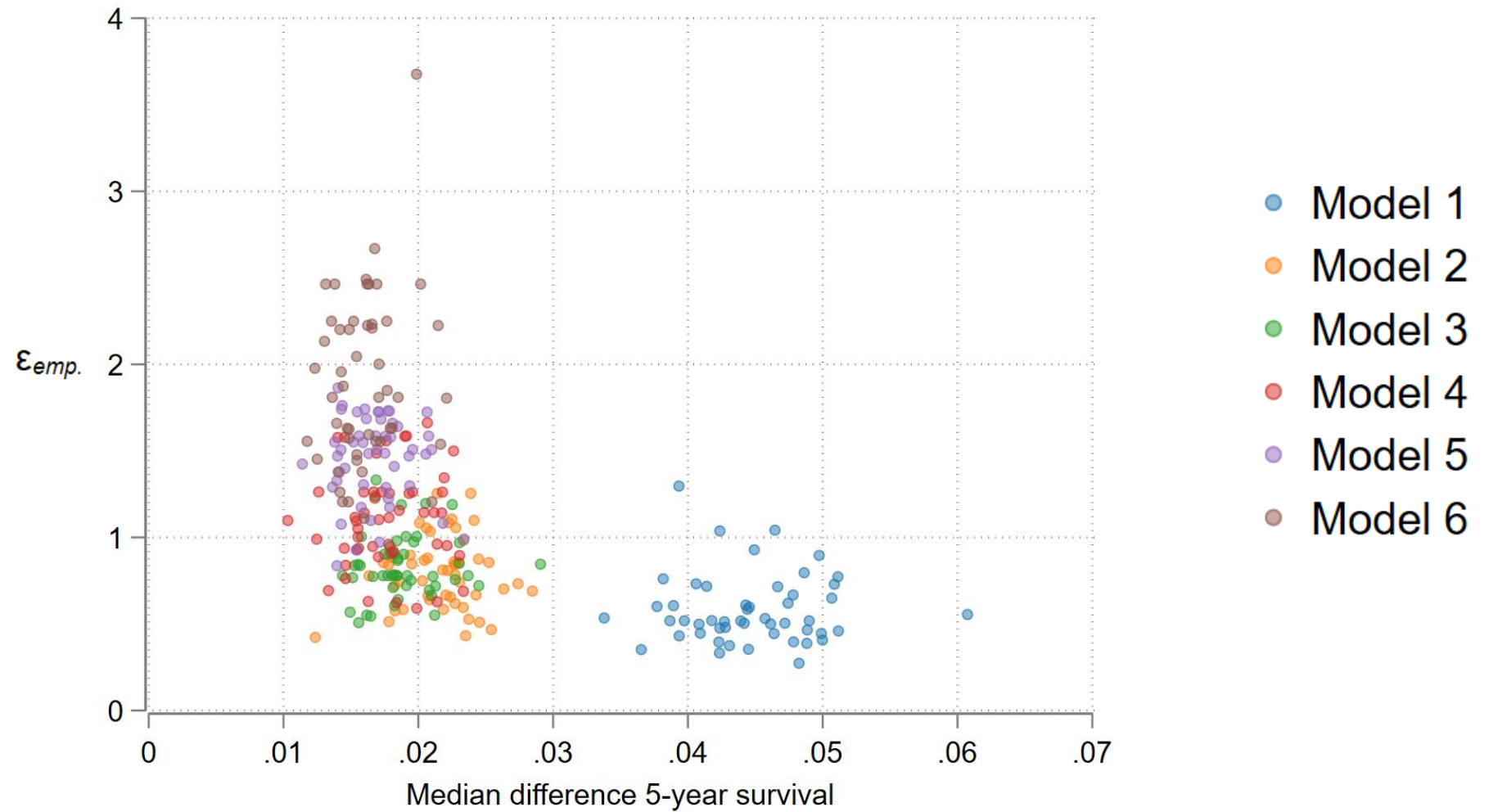
Privacy-utility trade-off



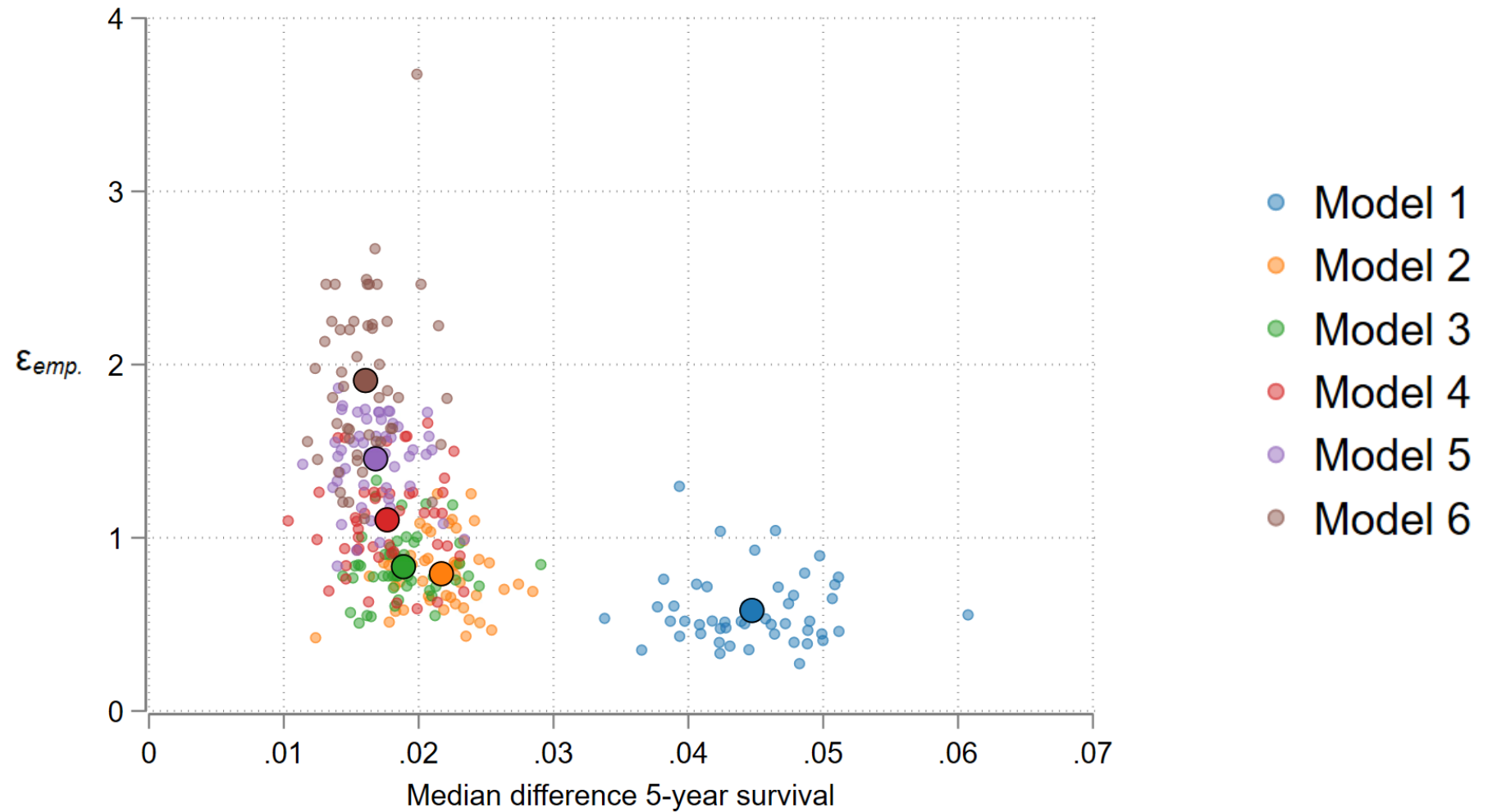
Privacy-utility trade-off



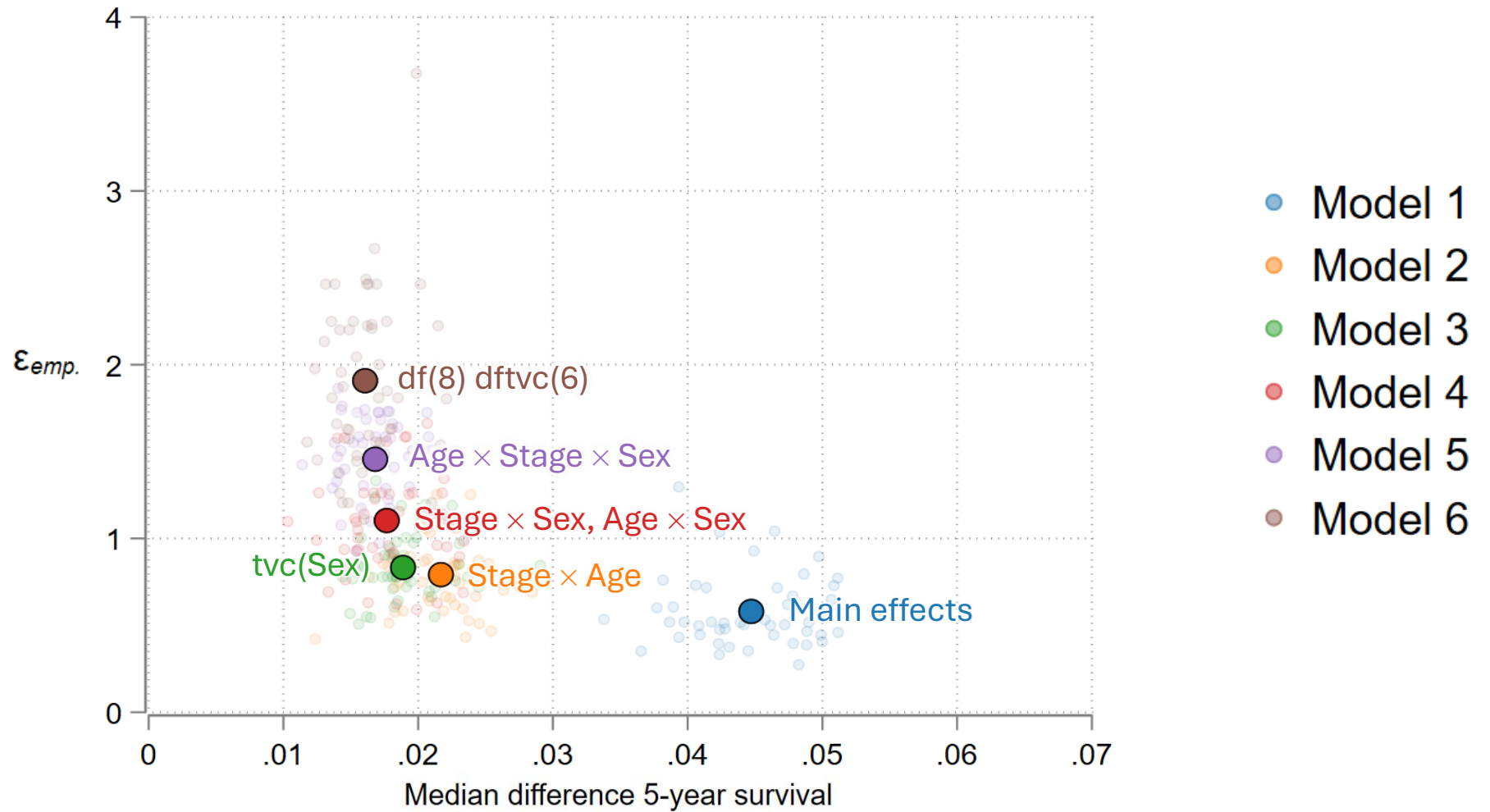
Privacy-utility trade-off



Privacy-utility trade-off



Privacy-utility trade-off



Thank you!

Questions?



GitHub

Cancer

Registry of Norway



A part of the Norwegian Institute of Public Health