



Data Without Barriers

Synthetic Data as a Catalyst for Responsible Innovation

Paul Tiwald

paul.tiwald@mostly.ai

Data Democratization

Data Access

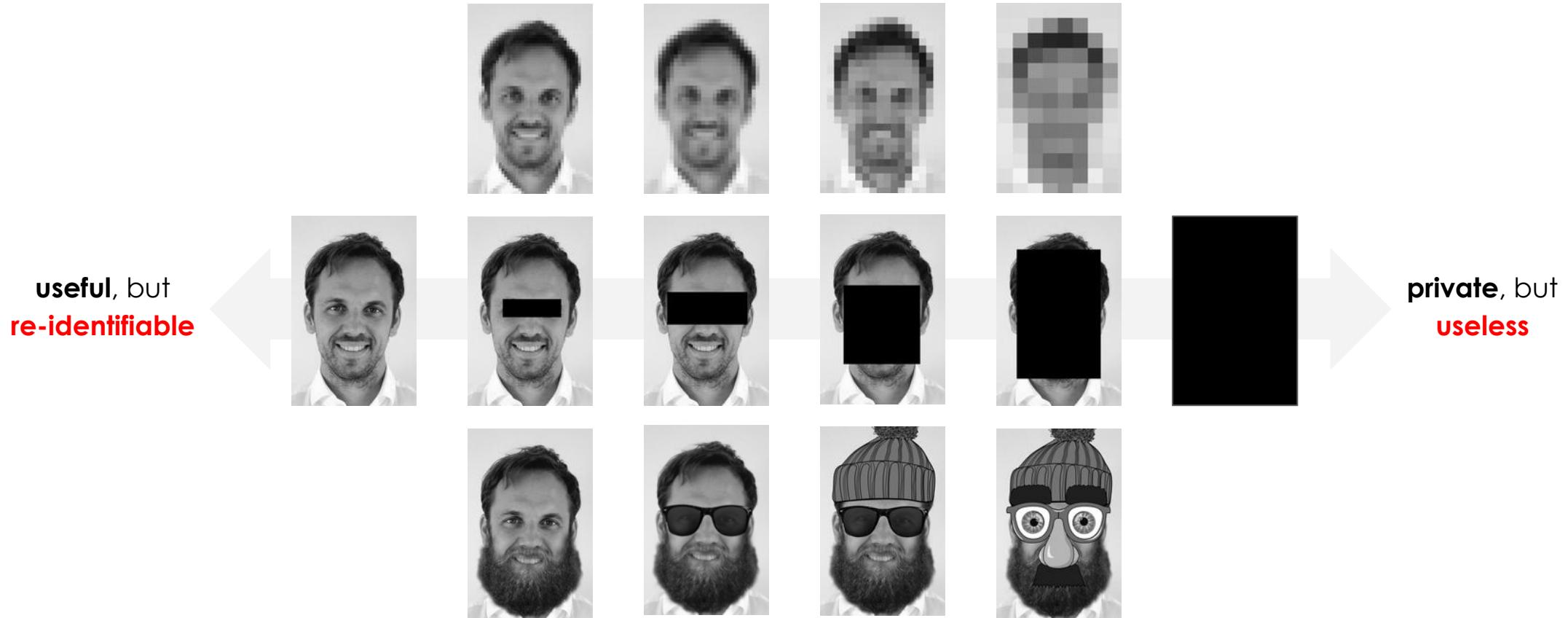
Data Insights

Synthetic Data

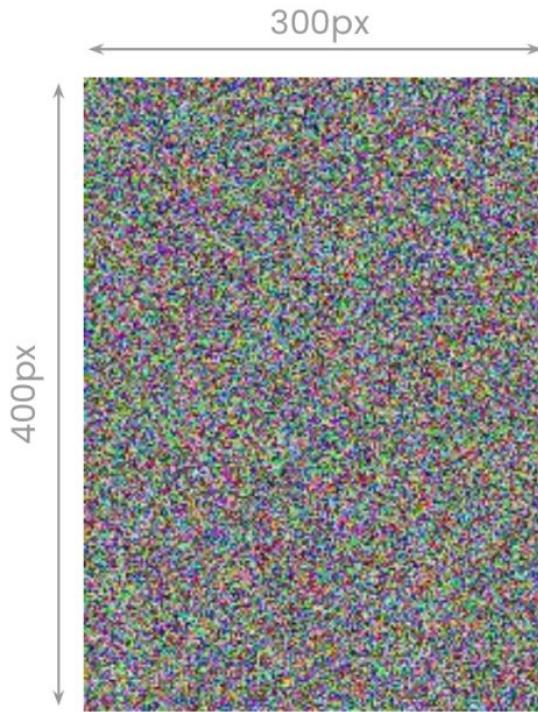
power of LLMs/assistants

... for everyone

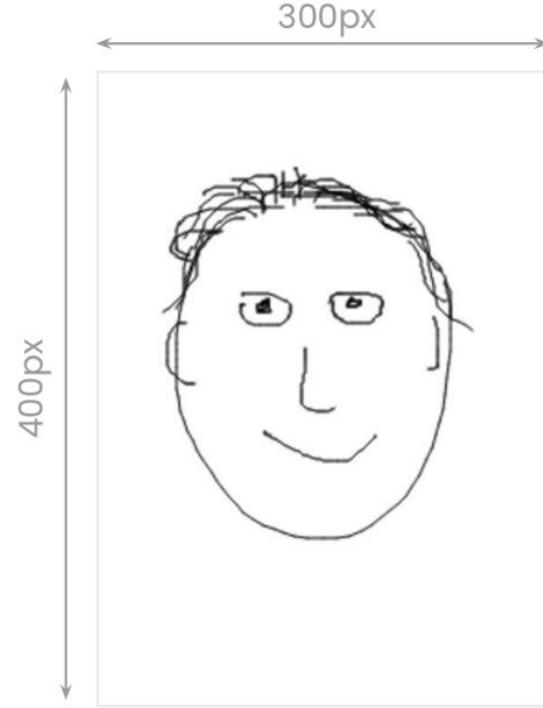
Why Synthetic? ... Real Data has its issues



What is Synthetic Data



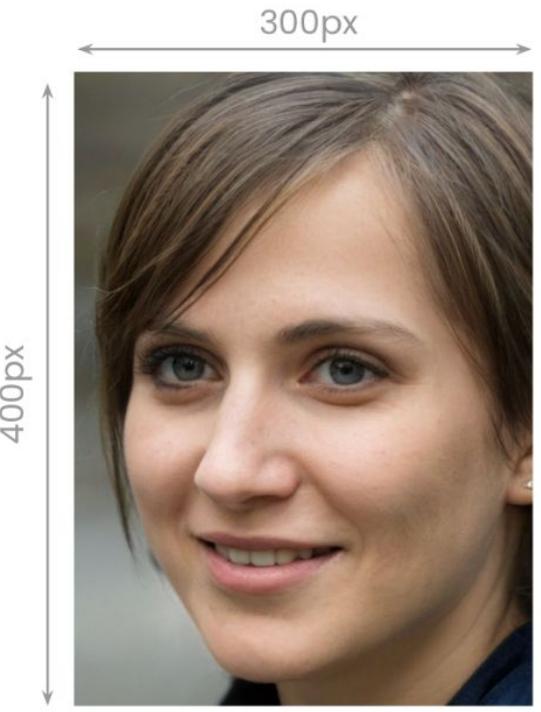
random data



self-generated data

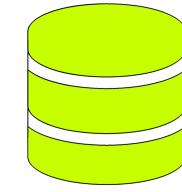
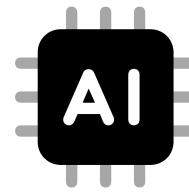
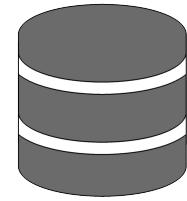


model-generated data
rule-based



AI-generated data
“data-based”

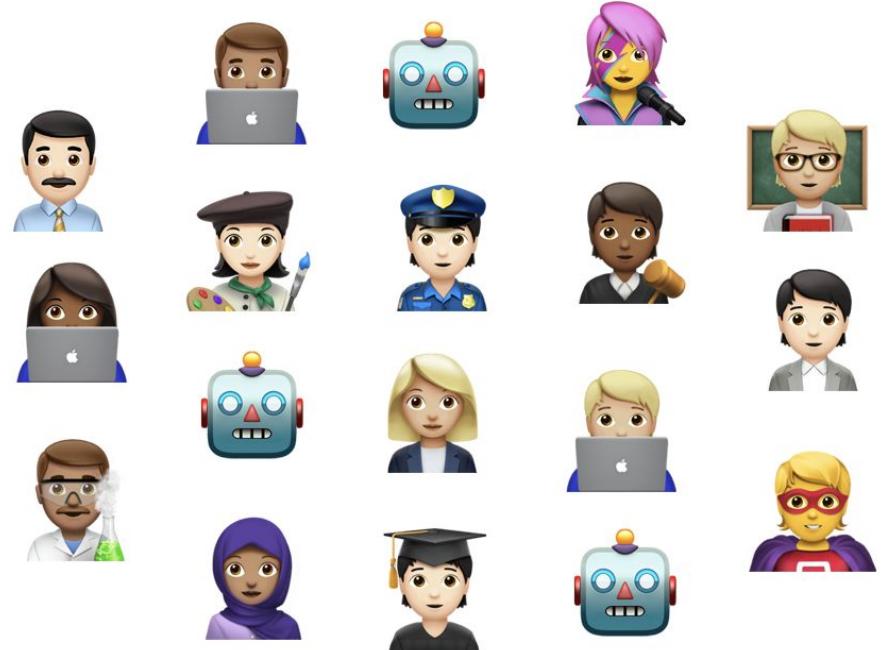
Synthetic Data = Generative AI



Actual Data
privacy-restricted
biased
incomplete

Generative Model

Synthetic Data
realistic
representative
anonymous
granular level



Data Consumers

people & algorithms

Tabular ARGN

Tabular ARGN - Auto-Regressive Generative Networks

arXiv > cs > arXiv:2501.12012

Search...

Help

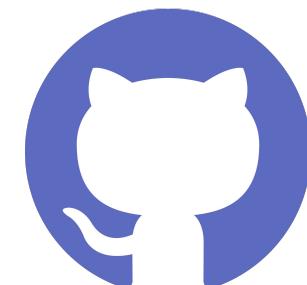
Computer Science > Machine Learning

[Submitted on 21 Jan 2025 (v1), last revised 6 Feb 2025 (this version, v2)]

TabularARGN: A Flexible and Efficient Auto-Regressive Framework for Generating High-Fidelity Synthetic Data

Paul Tiwald, Ivona Krchova, Andrey Sidorenko, Mariana Vargas Vieyra, Mario Scriminaci, Michael Platzer

Synthetic data generation for tabular datasets must balance fidelity, efficiency, and versatility to meet the demands of real-world applications. We introduce the Tabular Auto-Regressive Generative Network (TabularARGN), a flexible framework designed to handle mixed-type, multivariate, and sequential datasets. By training on all possible conditional probabilities, TabularARGN supports advanced features such as fairness-aware generation, imputation, and conditional generation on any subset of columns. The framework achieves state-of-the-art synthetic data quality while significantly reducing training and inference times, making it ideal for large-scale datasets with diverse structures. Evaluated across established benchmarks, including realistic datasets with complex relationships, TabularARGN demonstrates its capability to synthesize high-quality data efficiently. By unifying flexibility and performance, this framework paves the way for practical synthetic data generation across industries.

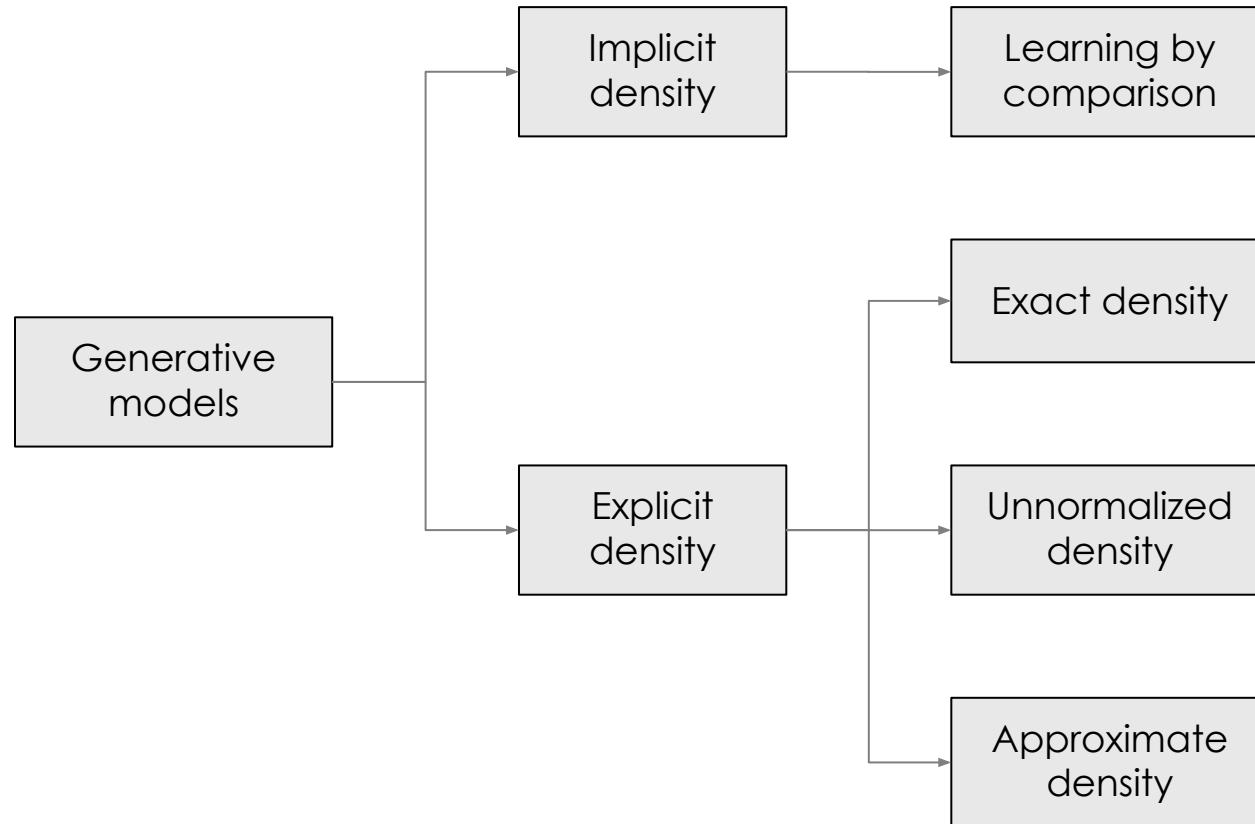


Tabular ARGN is implemented in the [Synthetic Data SDK](#):

<https://github.com/mostly-ai/mostlyai>

`pip install -U mostlyai[local]`

Taxonomy of deep generative models



- only care about generation, not $p(x)$
 - instead of maximizing the density, compare real vs generated sample (classification problem)
 - examples: GAN, GMMN
-
- directly learn density $p(x)$
 - examples: autoregressive models (Transformer, RNNs), flow-based models
-
- learn unnormalized density $E(x) \propto p(x)$
 - examples: EBM
-
- learn approximation (e.g. lower bound) of density $L(x) \leq p(x)$
 - examples: VAE, diffusion models

Flat Model

Fixed (column) Order Training Phase

patients data set:

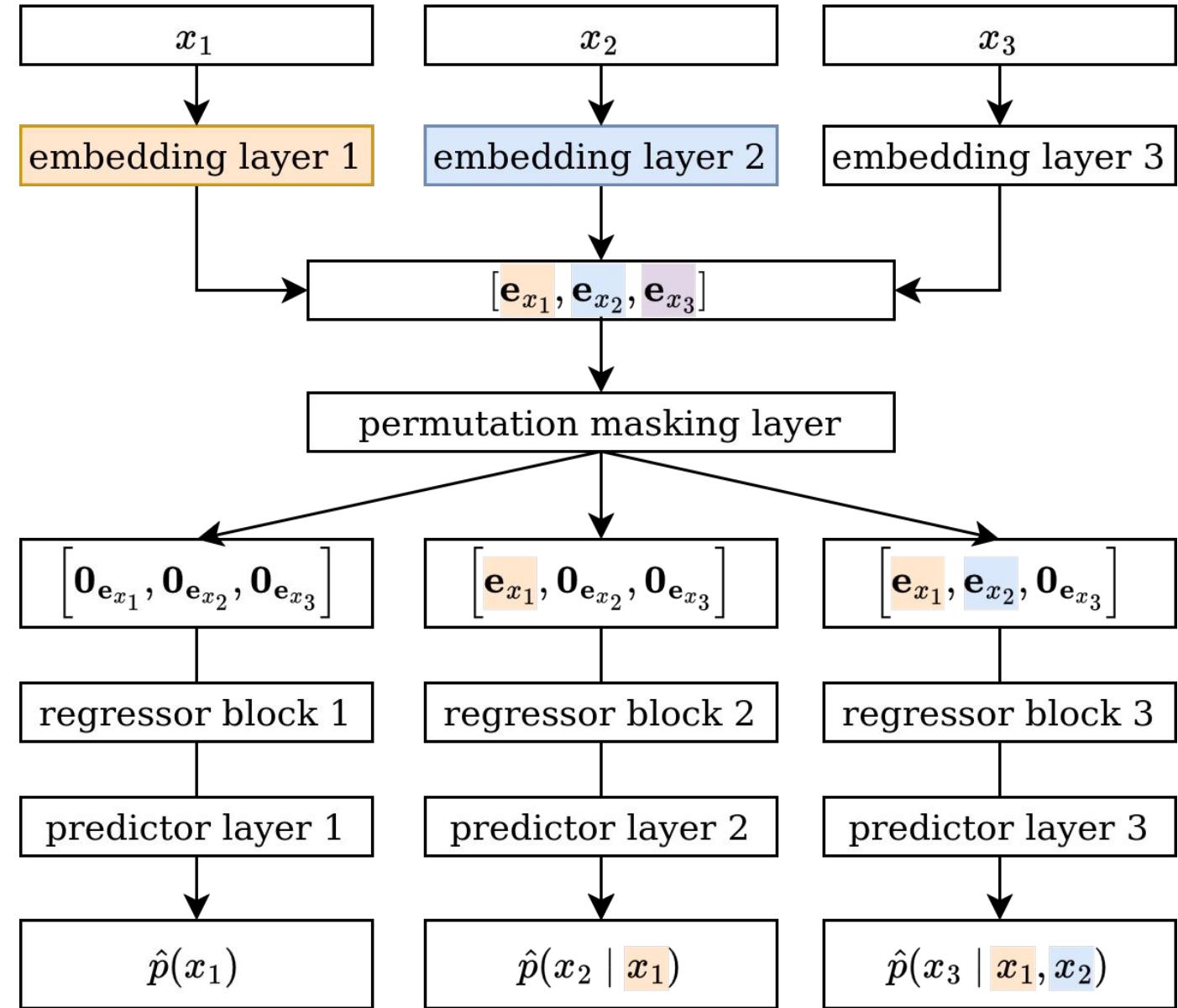
x_1 - age

x_2 - gender

x_3 - blood type

loss function:

$$\max_{\theta} \sum_{i=1}^D \log p_{\theta}(x_i | x_{<i})$$



Flat Model

Any (column) Order Training Phase

patients data set:

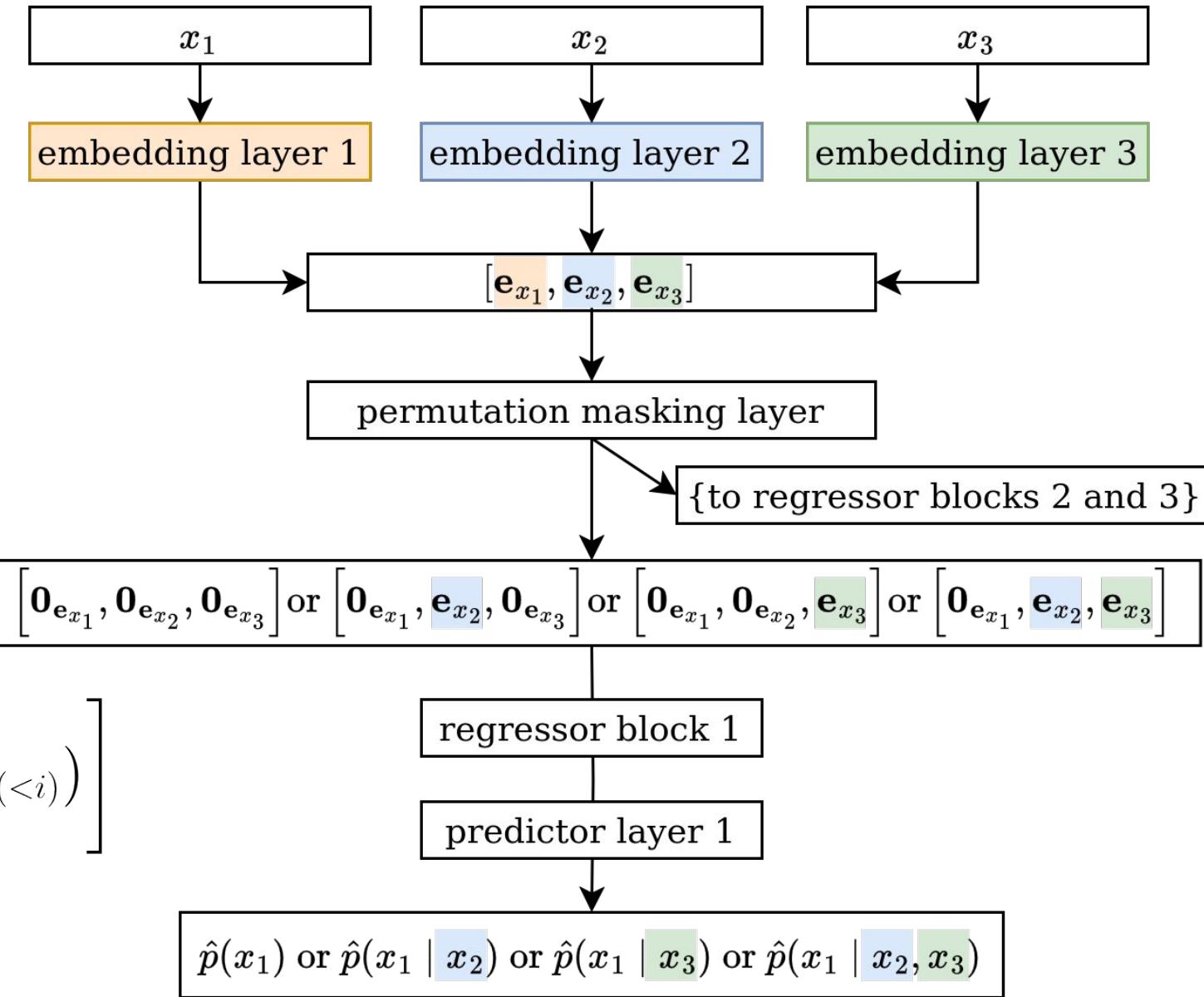
x_1 - age

x_2 - gender

x_3 - blood type

loss function:

$$\max_{\theta} \mathbb{E}_{\sigma \in \text{Uniform}(S_D)} \left[\sum_{i=1}^D \log p_{\theta}(x_{\sigma(i)} \mid x_{\sigma(<i)}) \right]$$



Flat Model

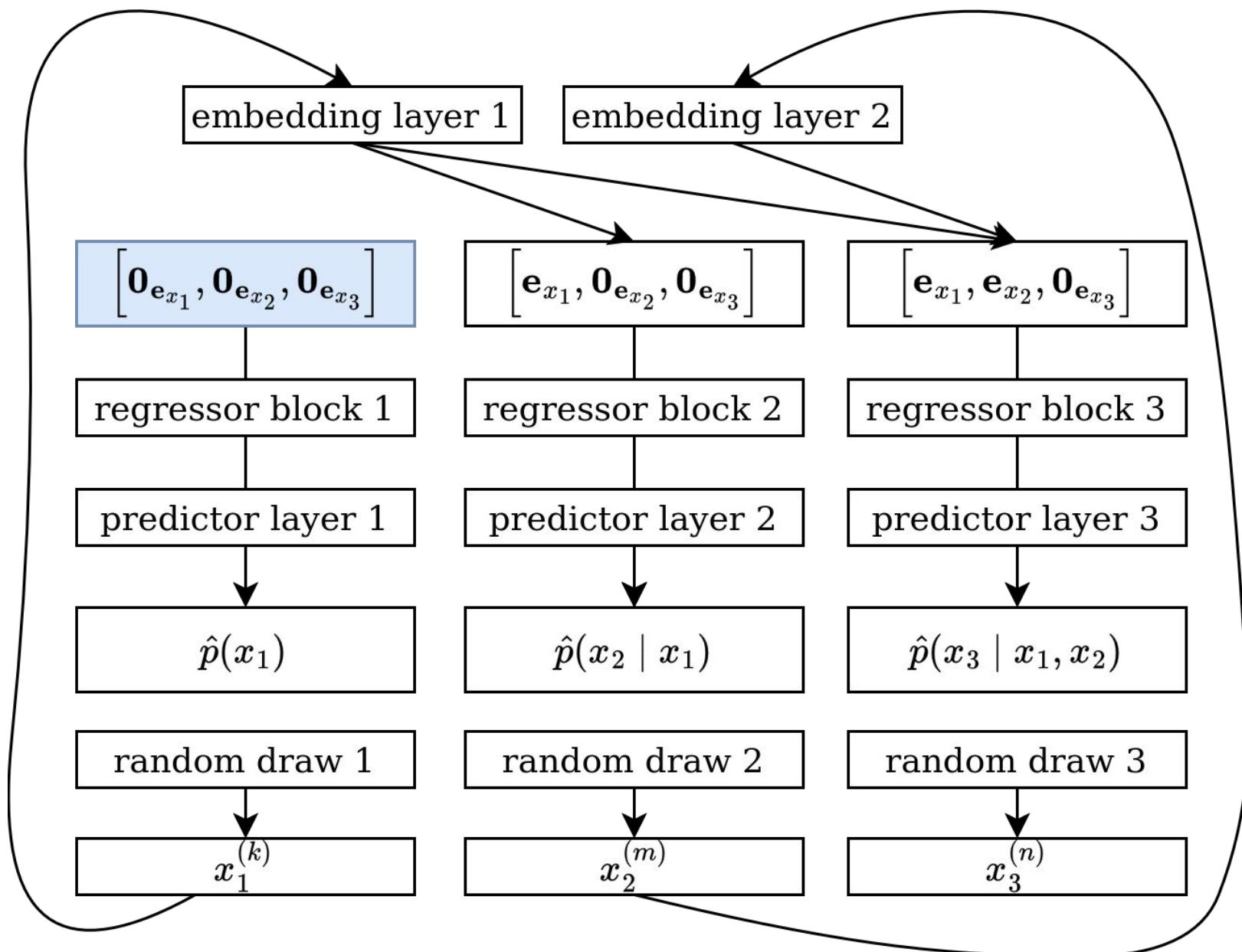
Generation Phase

patients data set:

x_1 - age

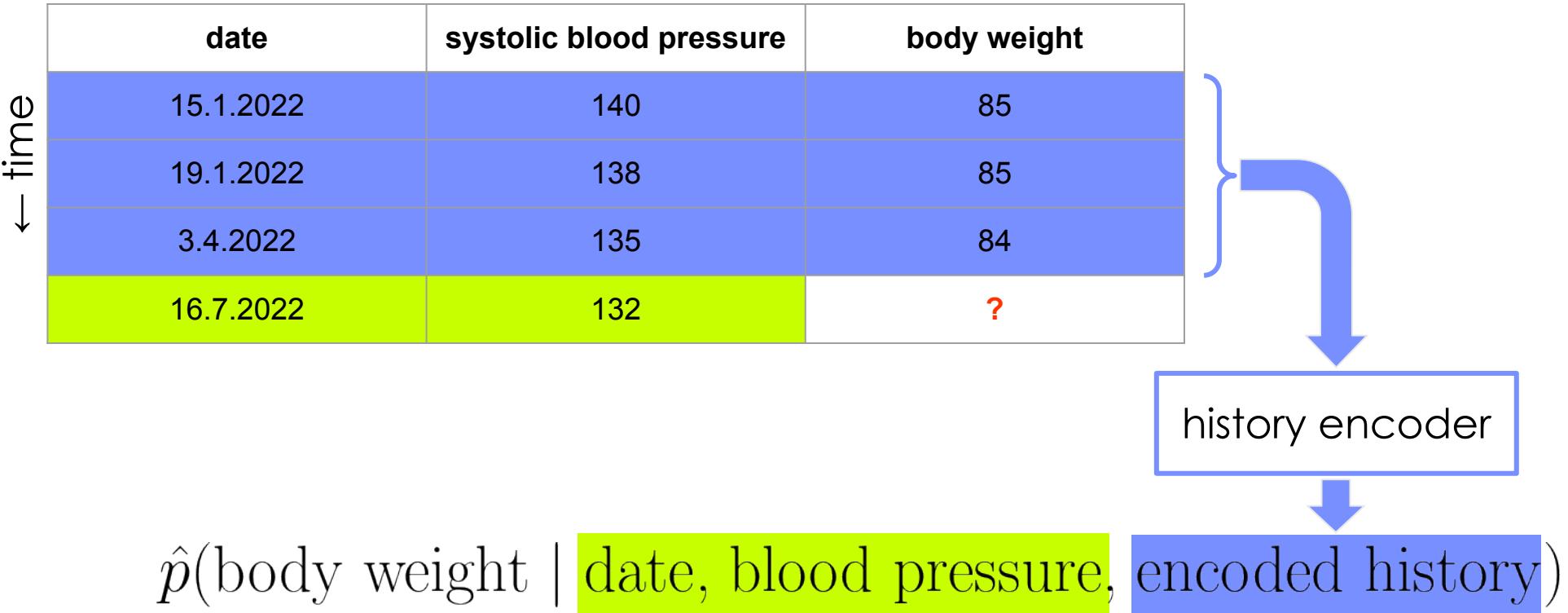
x_2 - gender

x_3 - blood type



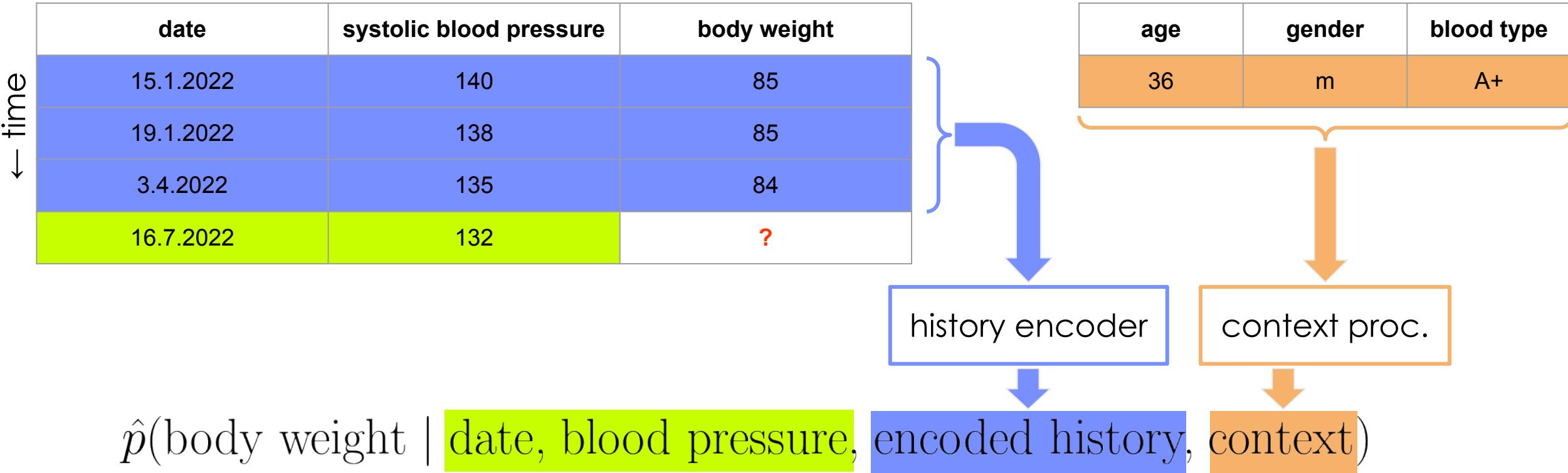
Sequential Model - doctor visits

auto-regressive along the column and the time dimensions

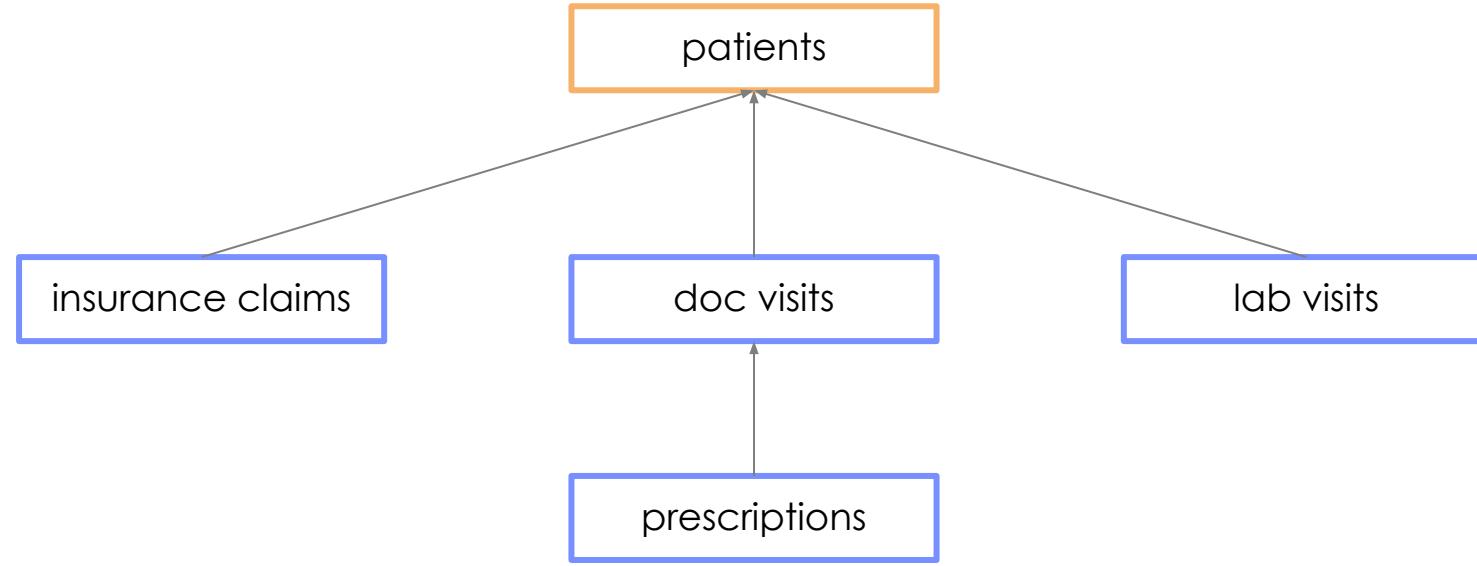


Sequential Model with context

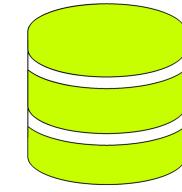
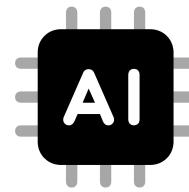
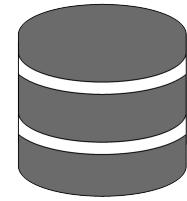
auto-regressive along the column, time, and table dimensions



Flexible context allows for synthesis of multi-table setups



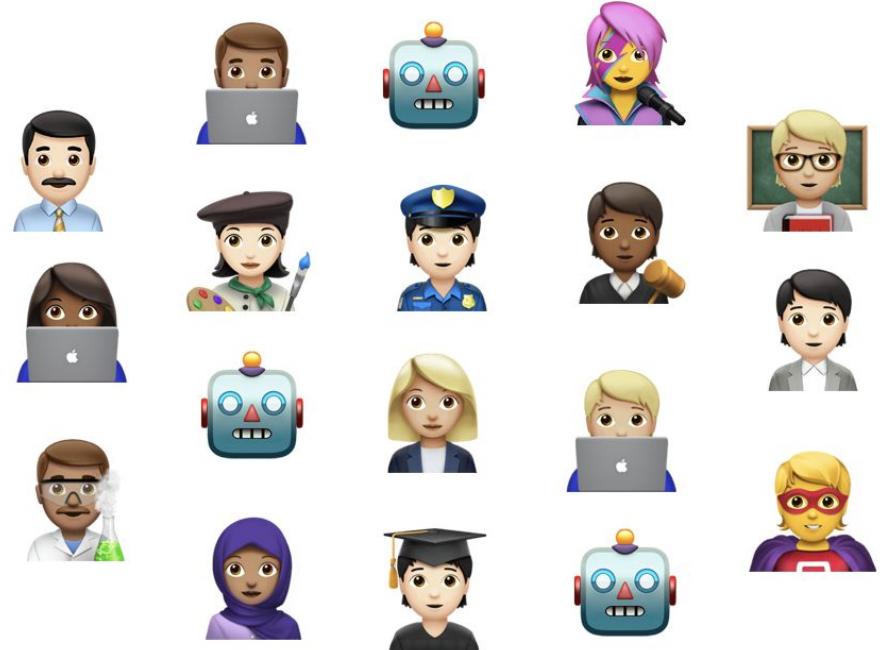
Synthetic Data = Generative AI



Actual Data
privacy-restricted
biased
incomplete

Generative Model

Synthetic Data
realistic
representative
anonymous
granular level



Data Consumers

people & algorithms

The main use case: Privacy and Data Access

Privacy/Data Access - the main use case

- reducing the “time-to-data”

Privacy/Data Access - the main use case

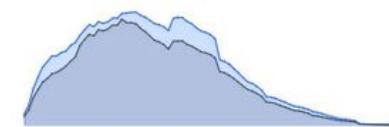
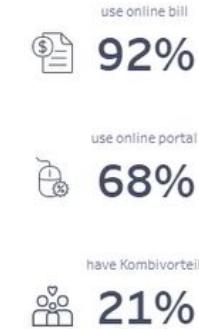
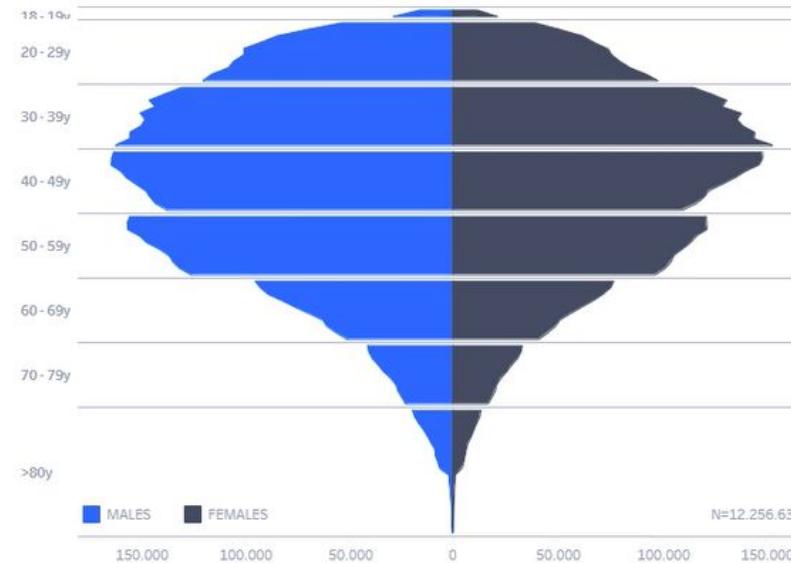
- reducing the “time-to-data”
- breaking down data silos within organizations
(e.g. synthetic-data products)

Know your customer

[OPEN CB EXPLORER](#)

Age distribution of B2C Postpaid Customers

Results of the synthesized data according to the "Title" field. Proportion of unknown and not clearly assignable contents excluded.



Privacy/Data Access - the main use case

- reducing the “time-to-data”
- breaking down data silos within organizations
(e.g. synthetic-data products)
- share data between subsidiaries in different countries
- share data between organizations (e.g. clean rooms)

Privacy/Data Access - the main use case

- reducing the “time-to-data”
- breaking down data silos within organizations
(e.g. synthetic-data products)
- share data between subsidiaries in different countries
- share data between organizations (e.g. clean rooms)
- open-data initiatives by public entities

make the information in EPCs accessible

EPC - Analysis
Synthetic EPC

Powered by


 Home
 Analysis
 User
 Help

Analysis

Analysis of EPCs

DPR412_classification	construction_year	degree_days	altitude	floors	net_area	heat_loss_surface
13	2017	2981	537	0	622.69	1664
8	1970	3396	840	0	164.15	579.5
8	2005	2426	197	1	129.18	487.23
4	1900	2765	354	1	74.37	228.65
13	1975	2591	265	0	2532.5	5743.9
5	1940	2778	283	0	6678.71	11930.02
6	1930	2741	345	0	299.21	1340.18
4	1900	3224	664	0	145.38	623.28
4	1992	2528	172	2	197.347	413.73
12	2022	2961	405	0	81.68	801.35

Rows per page
10
1-10 of 2519
<
<
>
>>

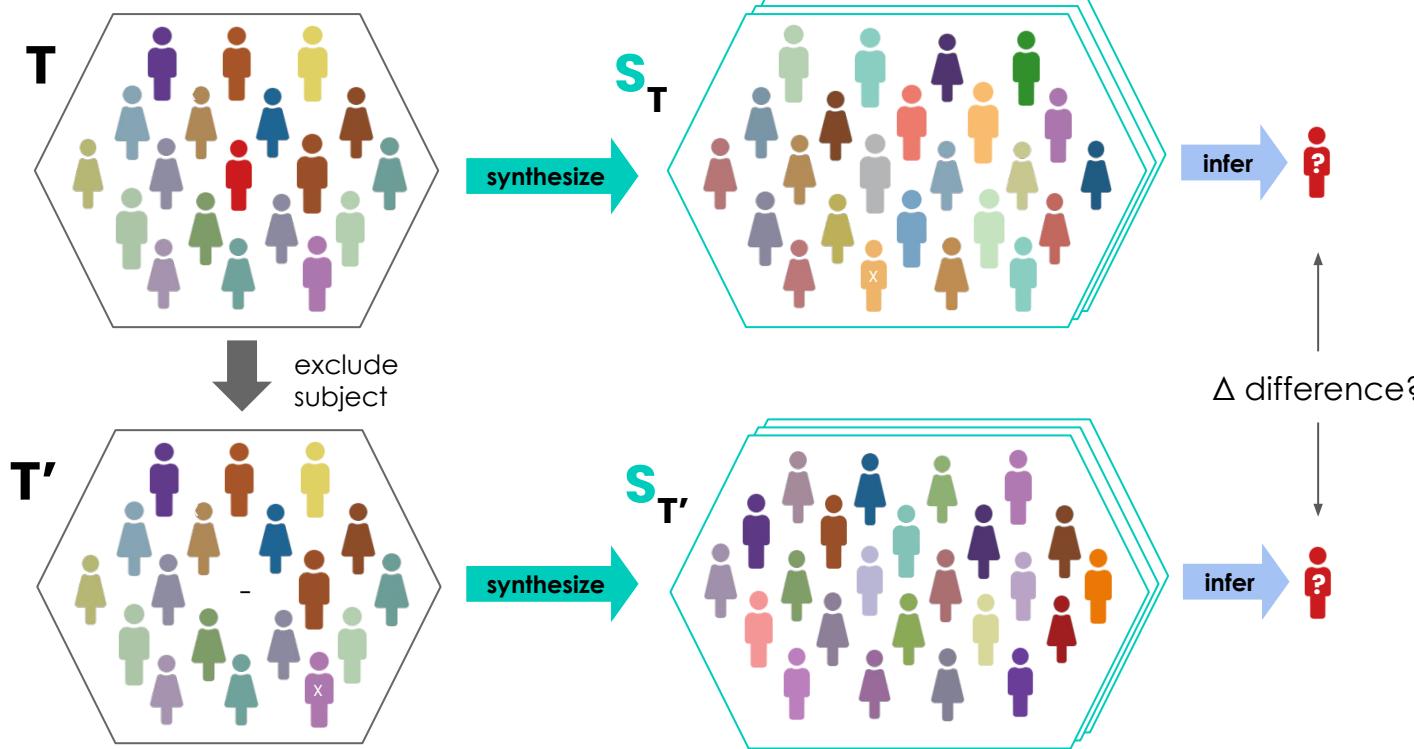


Horizon Europe research and innovation programme under grant agreement No 101069834. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or CINEA. Neither the European Union nor the granting authority can be held responsible for.

About MODERATE

In a nutshell
News
Contact

How to test Synthetic-Data Privacy



We must not be able to infer more about an individual, when that person is included in the database used for synthesis.

How to test Synthetic-Data Privacy - Empirically

Accuracy scores for 50 randomly chosen subjects, that were part of training

	NB	SVM	KNN	RF	LR	FRNN	ENS	DUM	RMEAN
Target T	42.8±5	43±6.6	41.6±9.3	49.8±9.4	38.3±3.7	49.1±9.8	45±6.8	32	44.2±3.8
Synthetic S_T	42.1±4.2	39.5±7.3	36.2±6.8	37.9±5.9	36.5±6	37±6.3	39.7±6.1	32	38.4±2

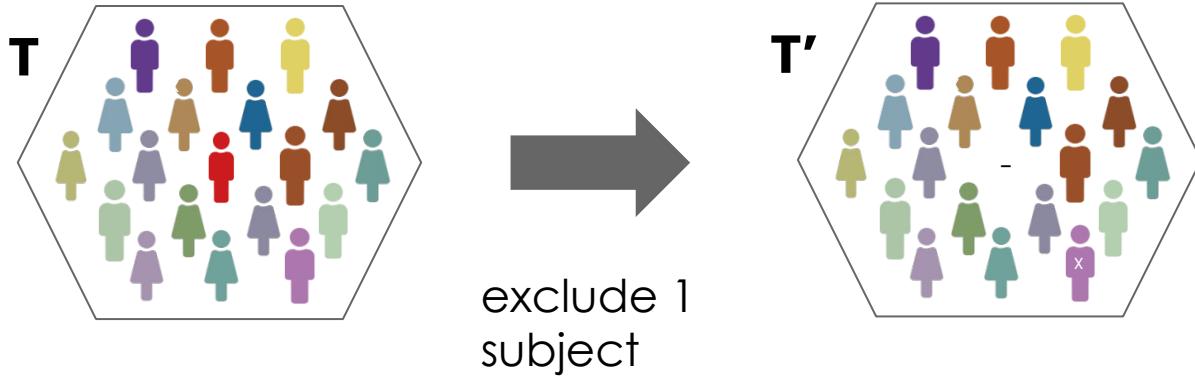
Accuracy scores for 50 randomly chosen subjects, that were NOT part of training

	NB	SVM	KNN	RF	LR	FRNN	ENS	DUM	RMEAN
Target T'	42.1±5	39.9±7.1	36.9±5.7	39.7±5.1	37.6±4.1	39.1±5	39.7±5.9	32	39.3±1.6
Synthetic $S_{T'}$	43.7±4.2	40.4±6.4	35.9±6.4	39.1±6.2	36.9±4.6	38±5.8	40.5±7	32	39.2±2.4

SBA Research, 2020

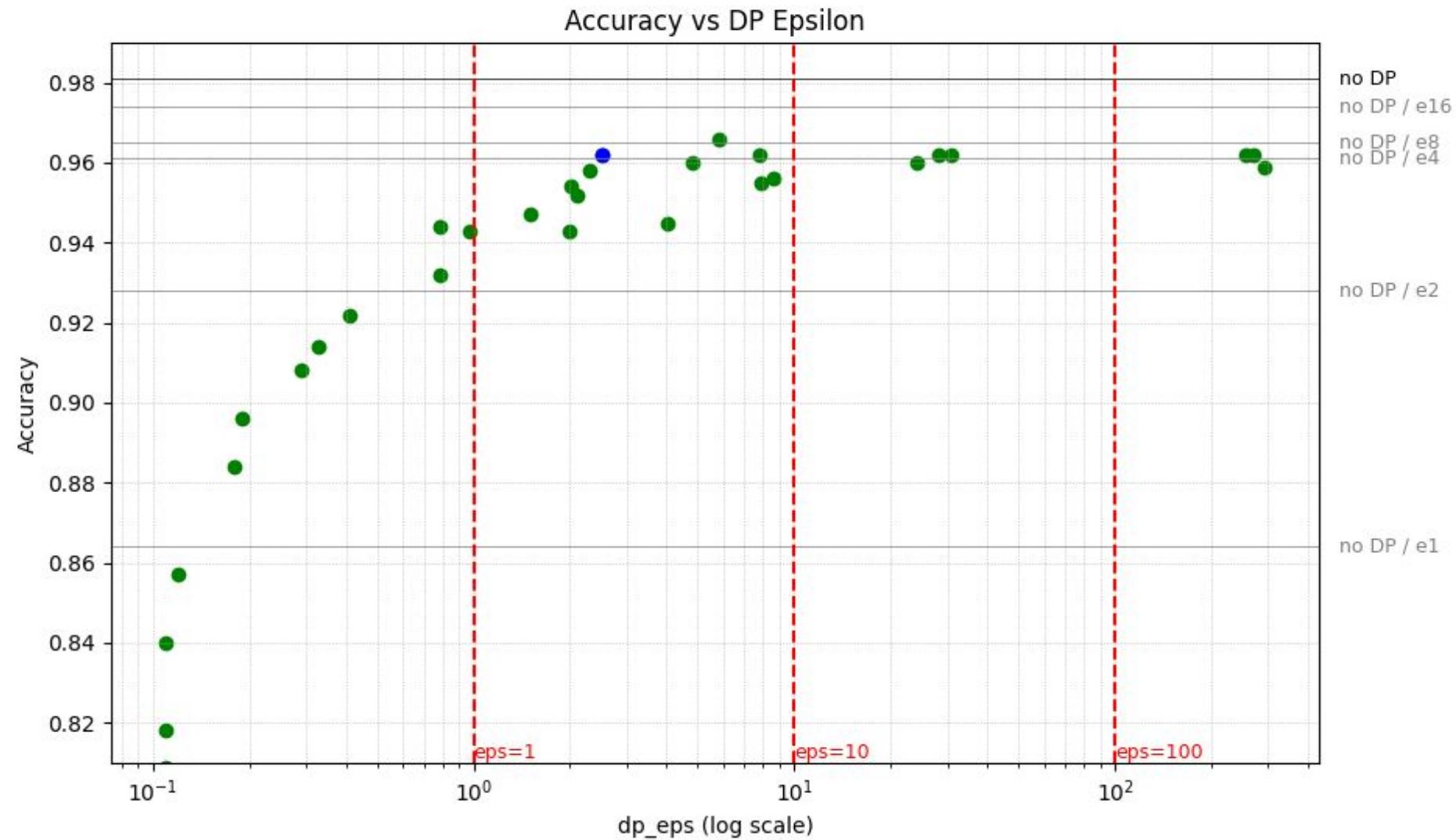
How to test Synthetic-Data Privacy - Mathematically

$$\Pr[\mathcal{A}(T)] \leq e^\epsilon \cdot \Pr[\mathcal{A}(T')]$$



- gold standard definition of privacy
- idea: limit influence of single individuals
- provides a mathematical guaranteed upper bound (ϵ) for the difference in outcomes of an algorithm A applied to the adjacent data sets T and T'

Differential Privacy



ML on synthetic data

Synthetic Switzerland

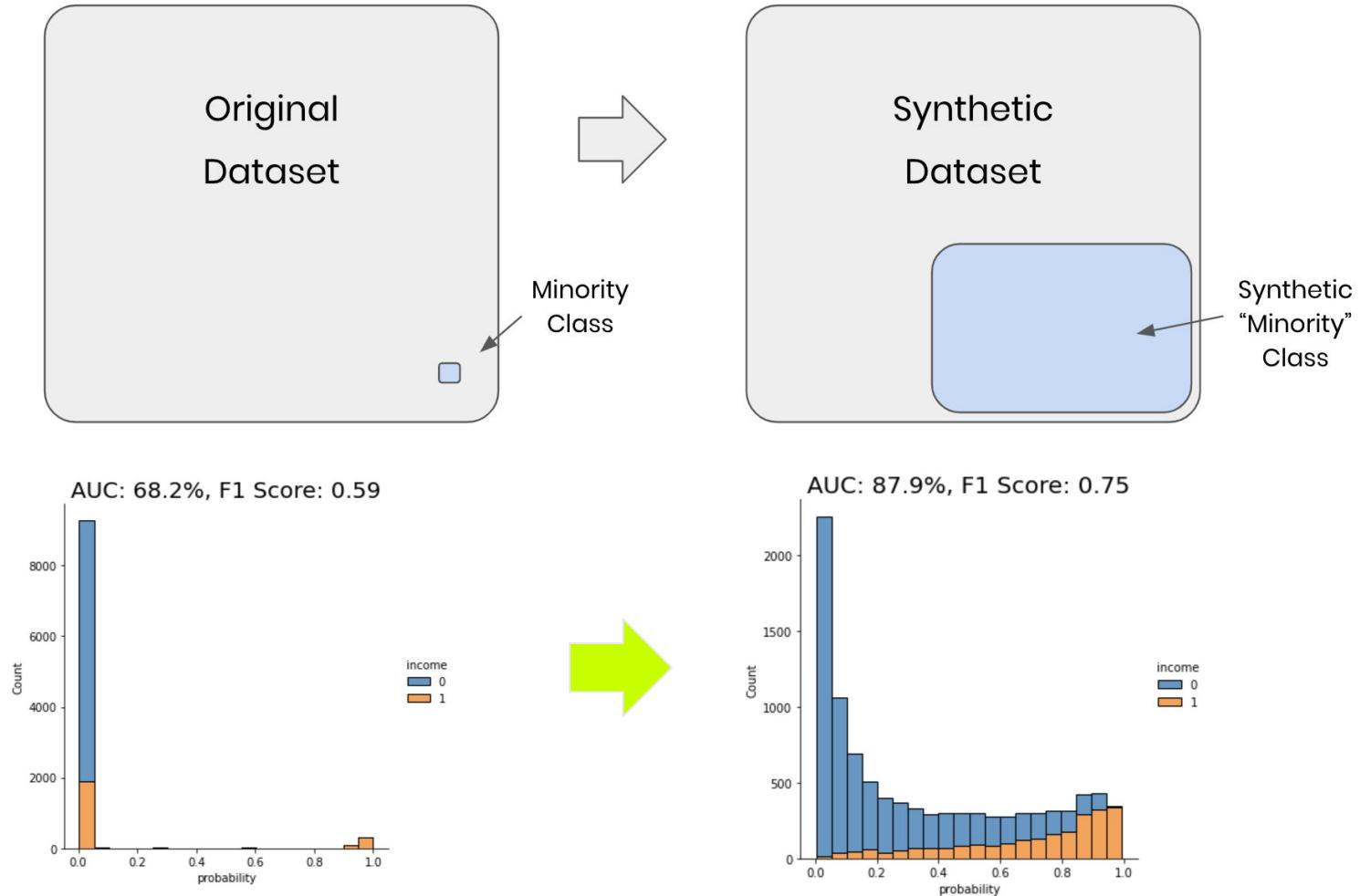


- Switzerland has around 9 Mio. Inhabitants
- Only 15% (1/9) provide explicit consent for marketing use
- 8/9 remain locked behind privacy
- Synthetic Switzerland mirrors population patterns without real identities
- We can now analyse and train models with 9/9 citizens instead of 1/9
- Zero personal data is used

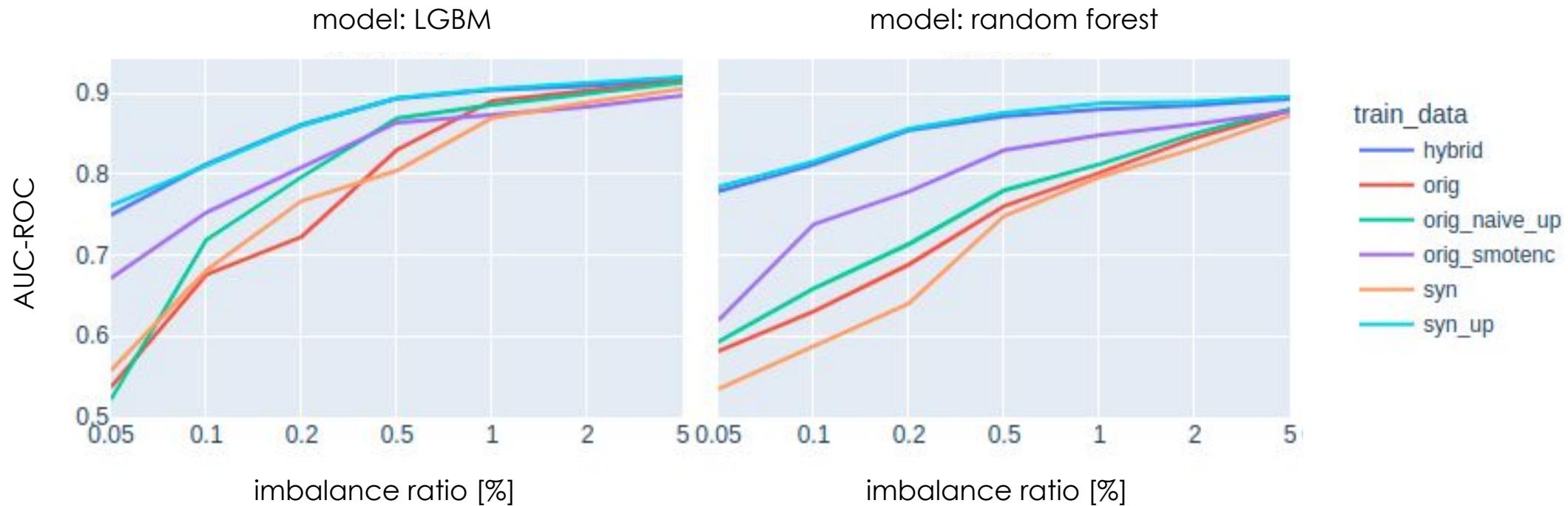
<https://www.youtube.com/watch?v=Rt5gXclc0jY>

Rebalancing of underrepresented Classes

 Open in Colab

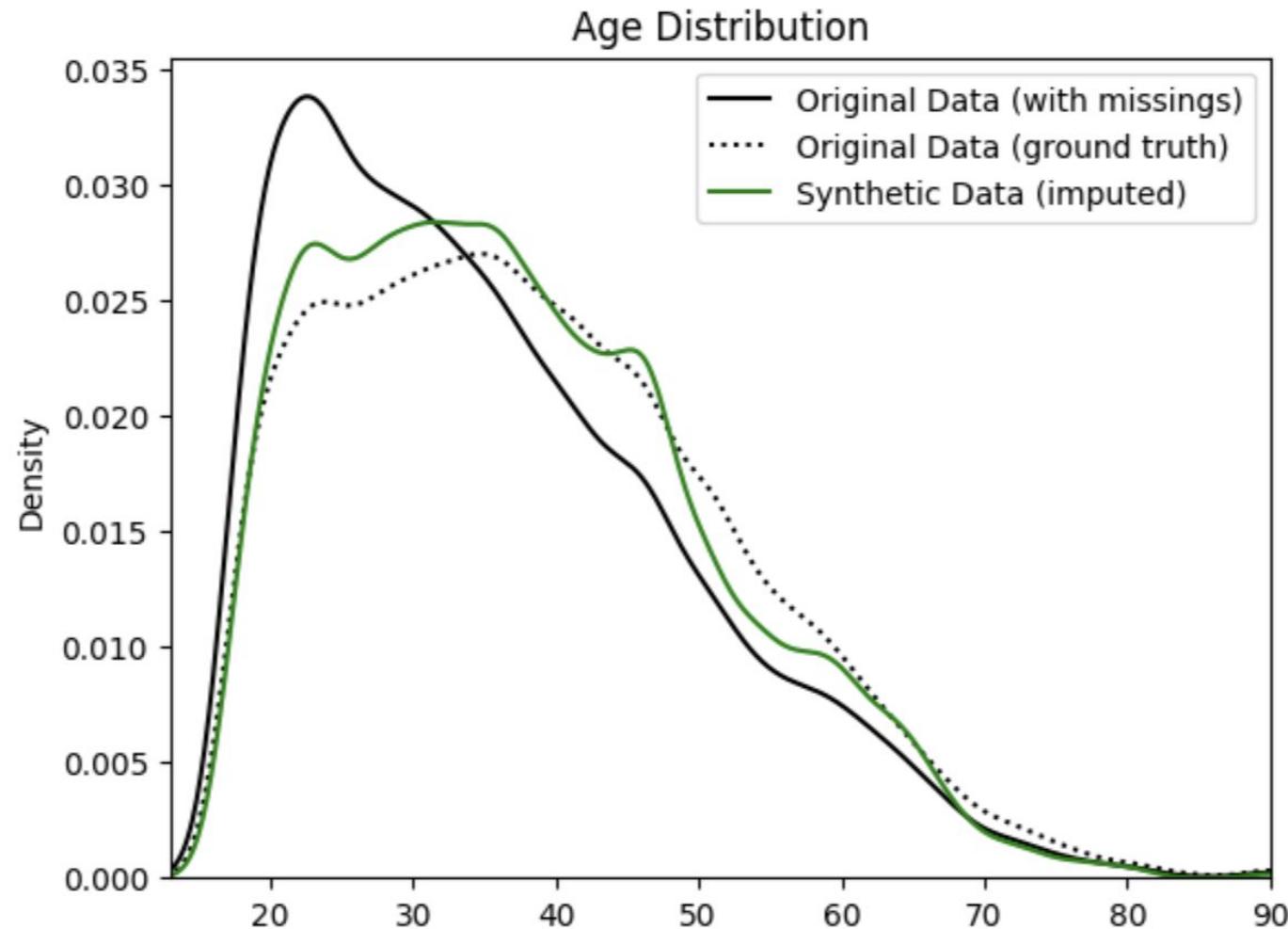


Rebalancing of underrepresented Classes



Imputation

Smart Imputation of Missing Data



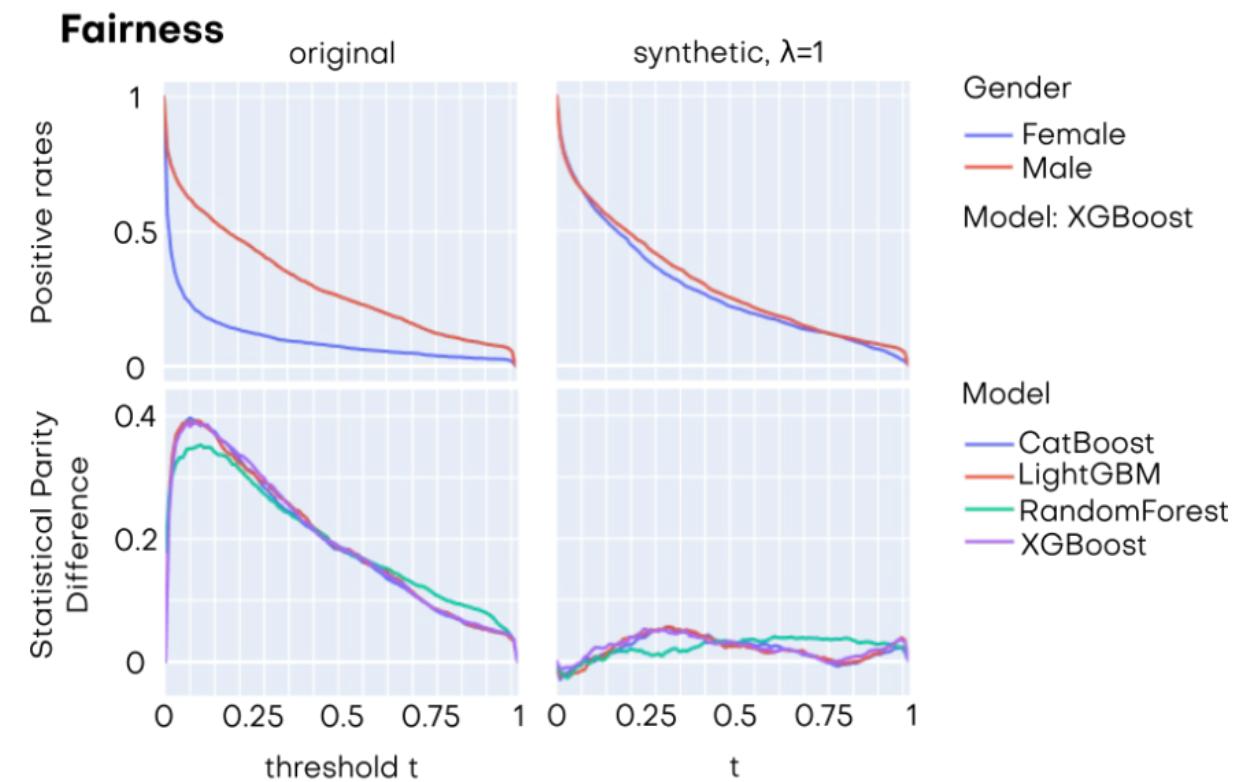
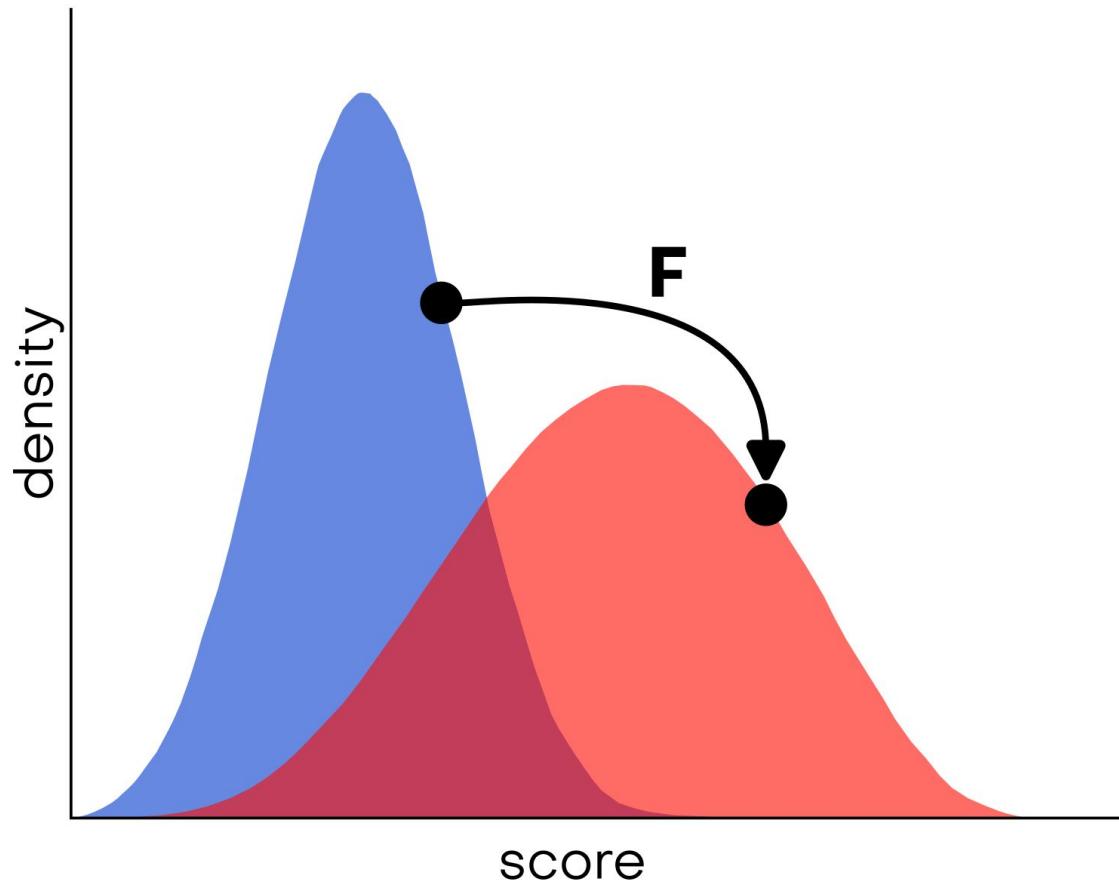
[Open in Colab](#)

Fair synthetic data

Fair Synthetic Data

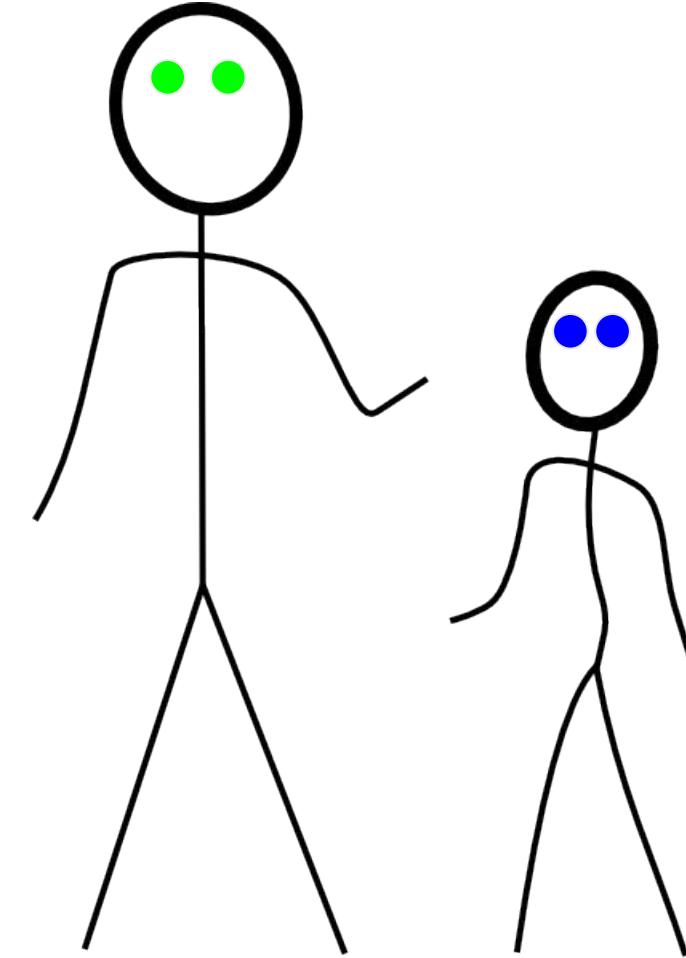
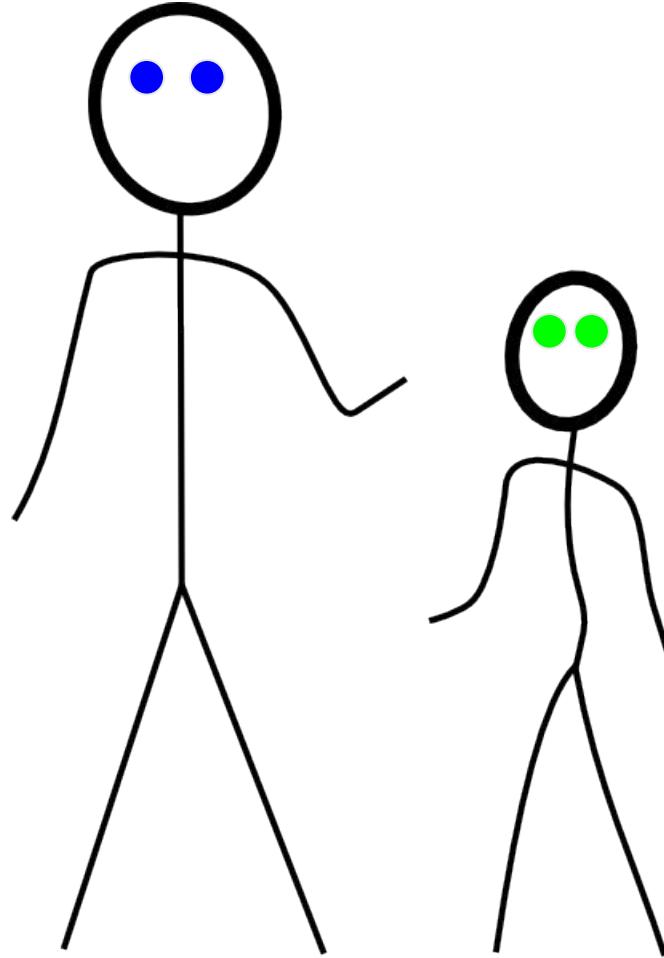
... subject to Strong Statistical Parity

<https://openreview.net/pdf?id=HbU5QuPZj6>

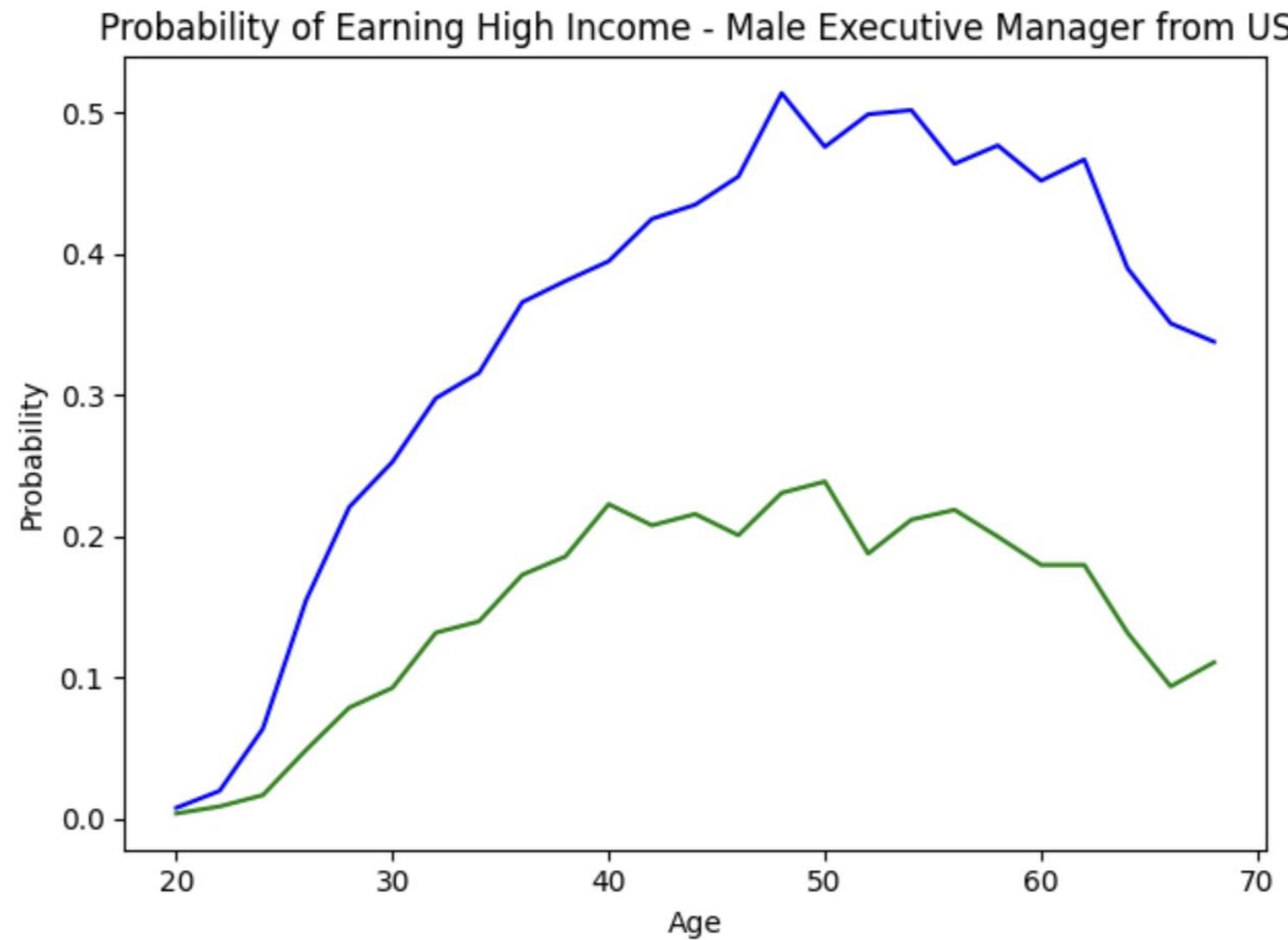


“Simulation”

Simulation with Synthetic Data



Simulation through Flexible Conditional Generation



[Open in Colab](#)

Simulation with Synthetic Data



International Journal of Research in
Marketing

Volume 39, Issue 4, December 2022, Pages 988-1018

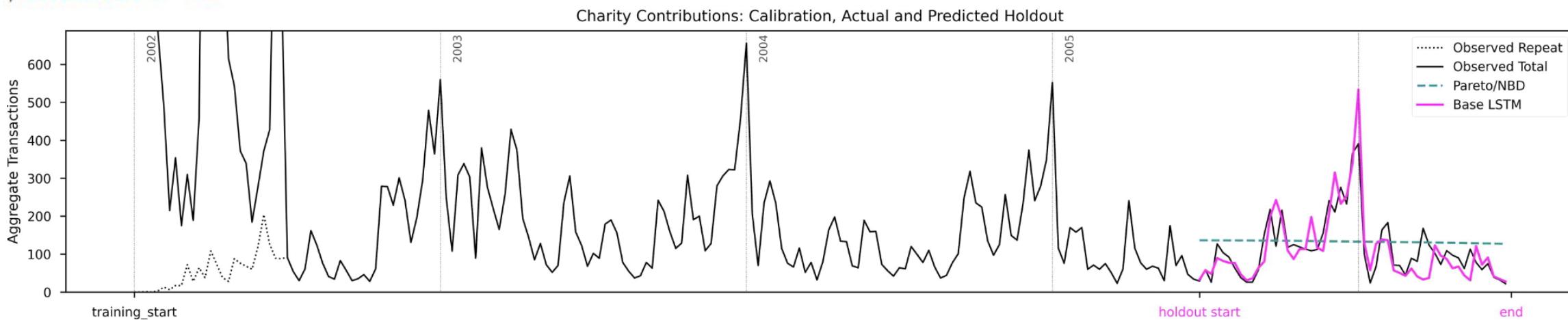


Full length article

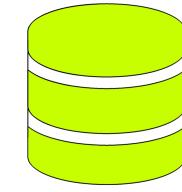
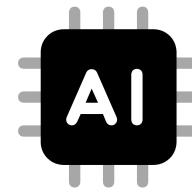
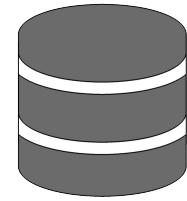
Customer base analysis with recurrent neural networks

Jan Valentin^a , Thomas Reutterer^a , Michael Platzer^b
, Klaudius Kalcher^b

[Customer base analysis with recurrent neural networks](#)



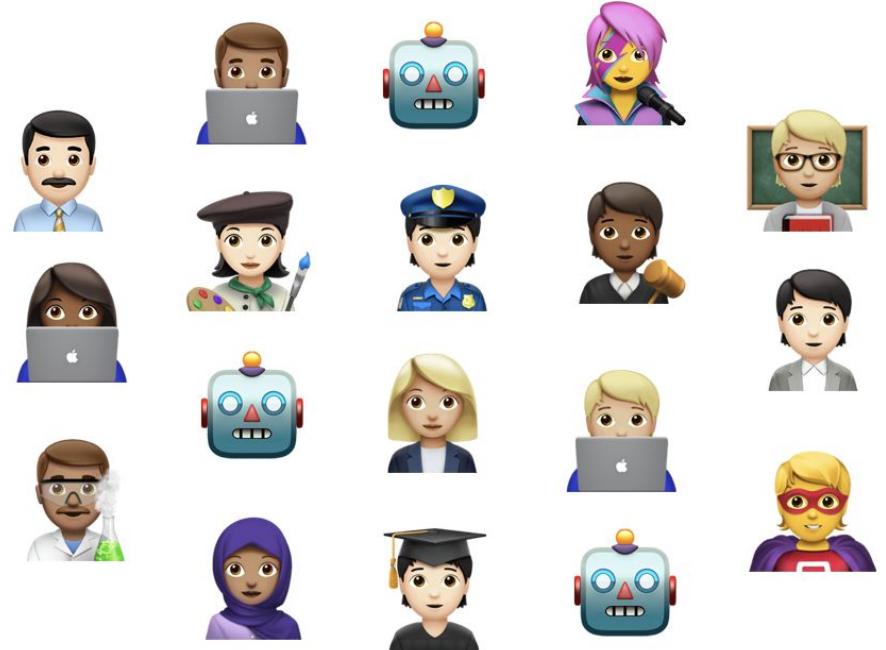
Synthetic Data = Generative AI



Actual Data
privacy-restricted
biased
incomplete

Generative Model

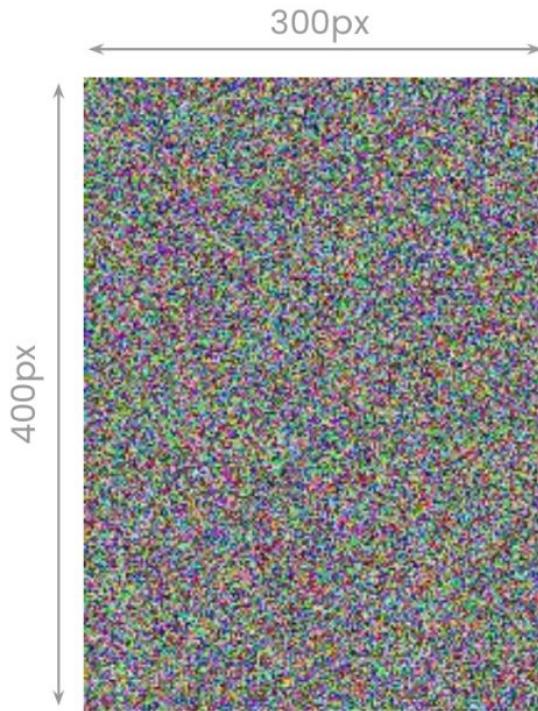
Synthetic Data
realistic
representative
anonymous
granular level



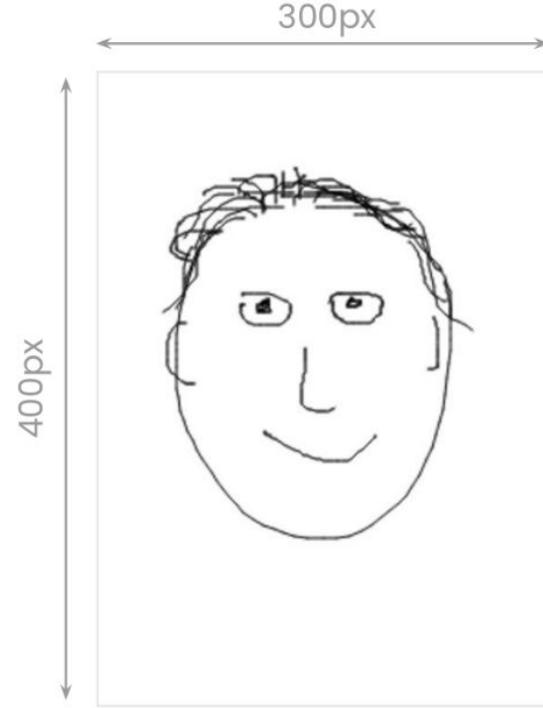
Data Consumers

people & algorithms

What is Synthetic Data



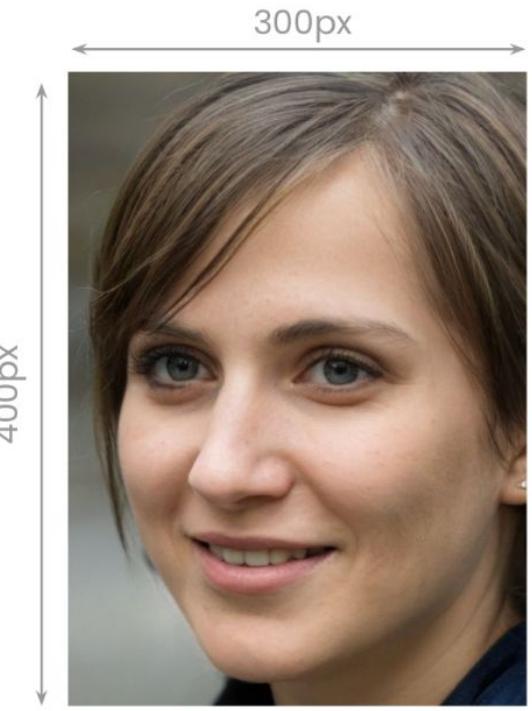
random data



self-generated data

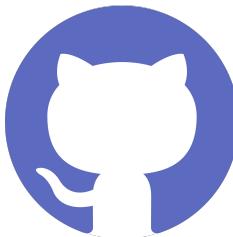
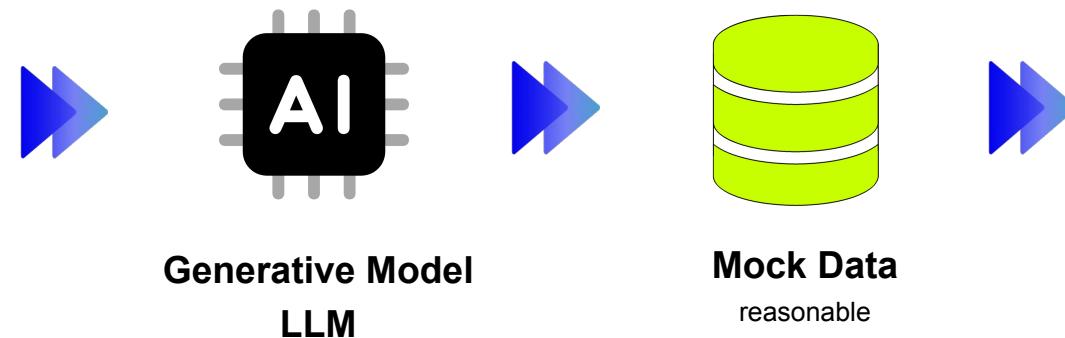


model-generated data
rule-based



AI-generated data
"data-based"

Get reasonable mock data out of “nothing”



<https://github.com/mostly-ai/mostlyai-mock>

`pip install -U mostlyai-mock`

Data Democratization

Data Access

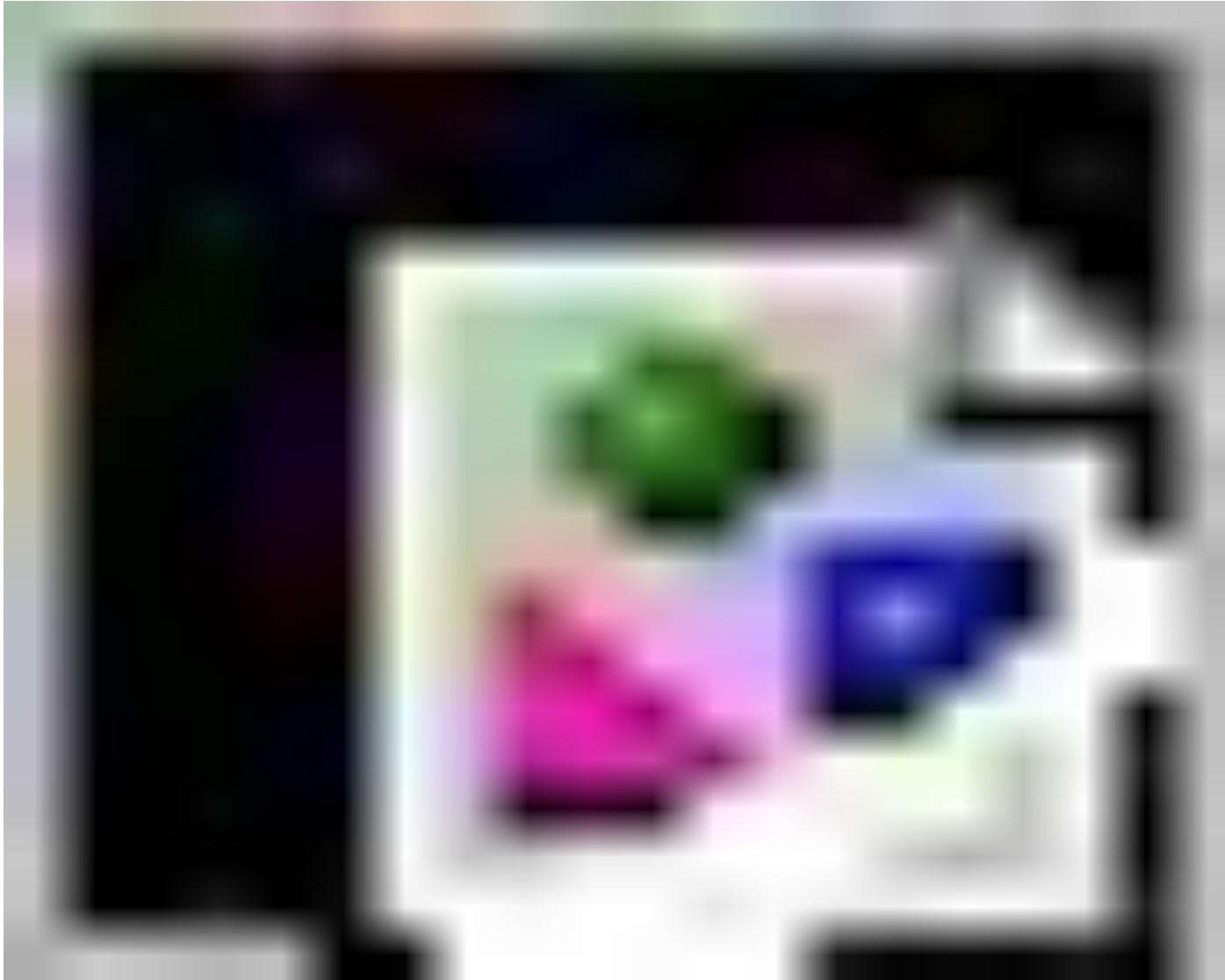
Data Insights

Synthetic Data

power of LLMs/assistants

... for everyone

Use natural language to get insights



app.mostly.ai

MOSTLY•AI

MOSTLY•AI

THE MOSTLY AI PRIZE



CASH PRIZE OF

100k USD

<https://www.mostlyaiprize.com/>