# Synthetic data as a novel anonymization technique: generation and evaluation approaches

Marta Batlle López, Data Scientist (Roche)

16 May 2024

# The core problems of sharing healthcare data

It is **challenging to access and share valuable healthcare data**

But high-quality data is still needed to **drive innovation**

In response, many **anonymisation techniques and privacy-enhancing technologies** were born
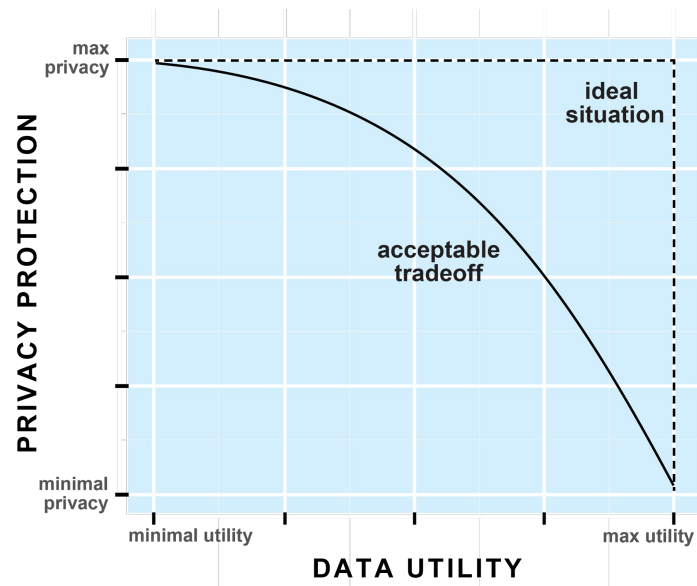
# How to choose between privacy enhancing technologies?

The quality and privacy trade-off

**Many** anonymisation techniques and privacy-enhancing technologies **were born** to tackle the problem of data sharing

Some can **compromise the quality of the data or pose technical challenges**

We look for that technology or technique that allows us to find an **acceptable trade-off** between the quality of the data and the level of privacy protection
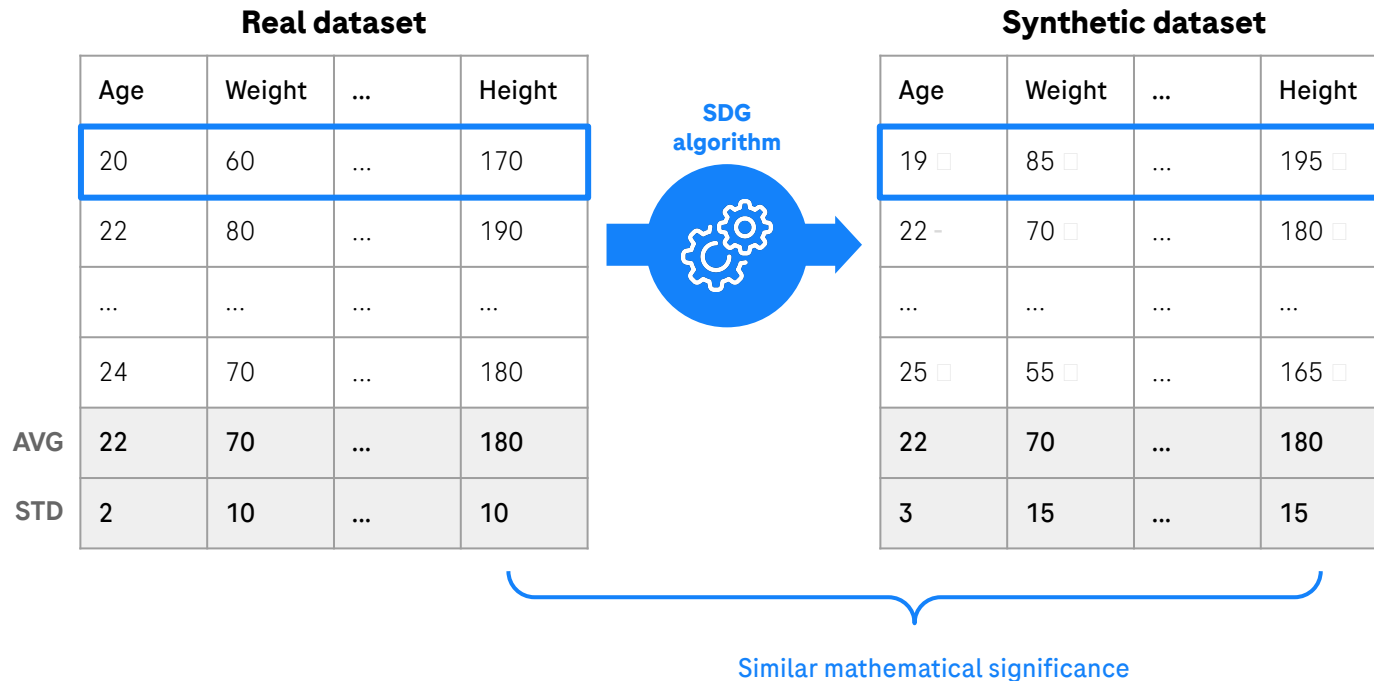


3

# How to choose between privacy enhancing technologies?

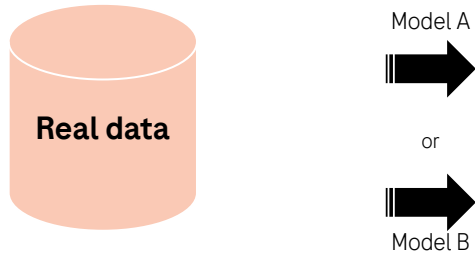| Technology | Pros | Cons |
|---|---|---|
| Homomorphic Encryption | High security & confidentiality (data usable even if encrypted) | Computationally intensive |
| Differential Privacy | Strong privacy guarantees (adds noise to data in a controlled way) | Data utility (may reduce accuracy of analysis) |
| Data Anonymization | Low risk data sharing (removes all PII data) | Data utility might be highly reduced |
| Pseudonymization | Easy data sharing (removes direct identifiers) | Risk of re-identification |
| Federated Learning | Collaborative learning (trains models on distributed data) | Complex coordination, potentially lower accuracy (aggregates) |
| Secure Multi-Party Computation (SMPC) | Strong Privacy for Joint Computations (Parties compute without revealing data) | Computationally intensive |
| Synthetic data | Private data sharing (generate artificial data based on real distributions) | Potentially lower utility of the data |

StyleGAN2 (Karras et al.)
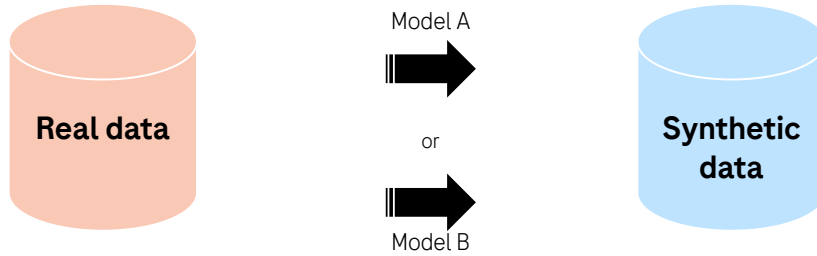
# What is synthetic data?

**Real dataset**

| Age | Weight | ... | Height |
|-----|--------|-----|--------|
| 20 | 60 | ... | 170 |
| 22 | 80 | ... | 190 |
| ... | ... | ... | ... |
| 24 | 70 | ... | 180 |

| | Age | Weight | ... | Height |
|-----|-----|--------|-----|--------|
| **AVG** | 22 | 70 | ... | 180 |
| **STD** | 2 | 10 | ... | 10 |

**SDG algorithm**

**Synthetic dataset**

| Age | Weight | ... | Height |
|-----|--------|-----|--------|
| 19 | 85 | ... | 195 |
| 22 | 70 | ... | 180 |
| ... | ... | ... | ... |
| 25 | 55 | ... | 165 |

| Age | Weight | ... | Height |
|-----|--------|-----|--------|
| 22 | 70 | ... | 180 |
| 3 | 15 | ... | 15 |

Similar mathematical significance

6

# Synthetic Data Generation process

**Real data**

# Synthetic Data Generation process



**Real data**

Model A

or

Model B

# Synthetic Data Generation process

Real data

Model A

or

Model B

Synthetic data

# Synthetic Data Generation process

Real data

Model A

or

Model B

Synthetic
data

Quality and privacy tests

# Synthetic Data Generation process

Real data

Model A

or

Model B

Synthetic data

Quality and privacy tests

**Quality and privacy report**

# Synthetic Data Generation process



Real data

Model A

or

Model B

Synthetic data

Privacy and quality tests

**Quality and privacy report**

# Generic utility evaluation metrics

How do we assess the quality of the synthetic data?



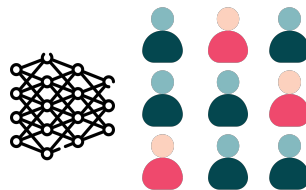**Feature distributions**



**Feature correlations**



**Distance metrics**



80%

75%

**ML performance**



**Distinguishability**

# Measuring privacy risks

How do we assess the privacy of the synthetic data?

An effective anonymisation prevents an attacker from:

1. **Singling out an individual** in a dataset
2. **Linking two records** within a dataset (or between two separate datasets)
3. **Inferring any information about individuals** in such dataset

Original sample

| Age | Gender | BMI | Hypertension |
|-----|--------|------|--------------|
| 58  | F      | 32.4 | True         |
| 66  | M      | 25.4 | False        |
| 35  | M      | 27.8 | True         |

Synthesis

Synthetic sample

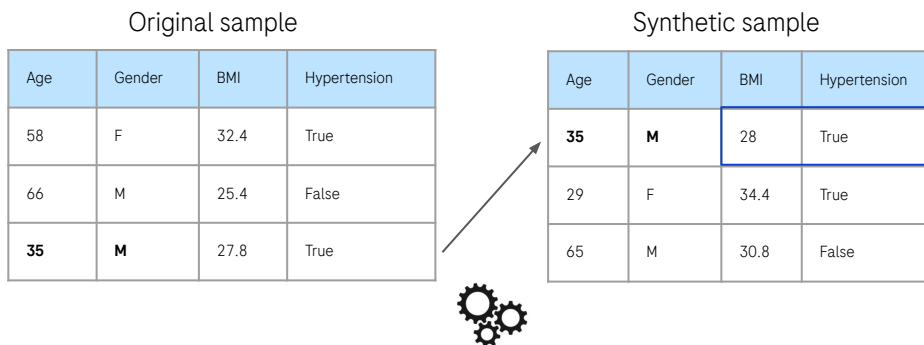| Age | Gender | BMI  | Hypertension |
|-----|--------|------|--------------|
| 35  | M      | 28   | True         |
| 29  | F      | 34.4 | True         |
| 65  | M      | 30.8 | False        |

# Measuring privacy risks
How do we assess the privacy of the synthetic data?

**Identity Disclosure** involves the risk of inferring the true identity of an individual from synthetic data
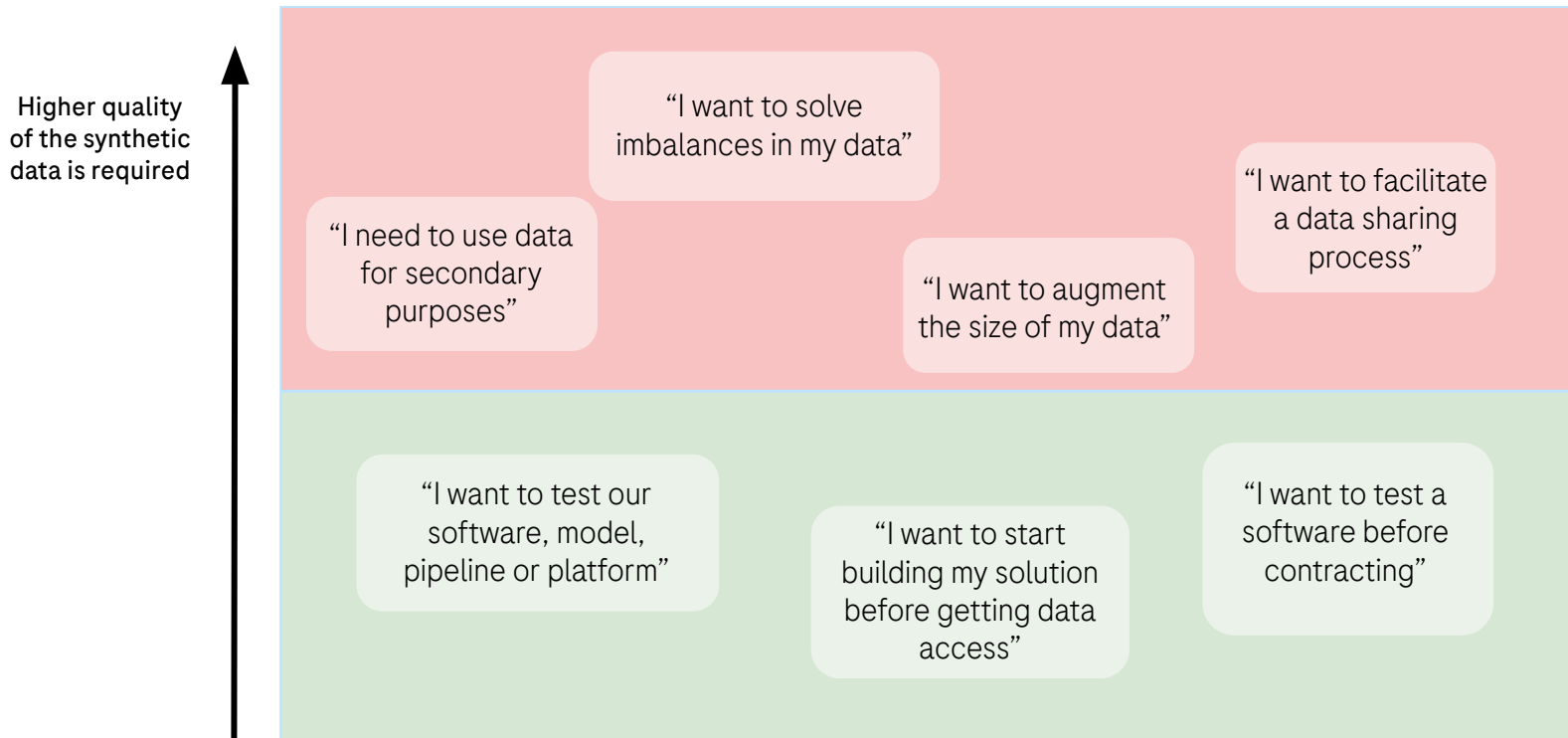
Adversary gains information on the individual through the **matched record**.

**Membership Disclosure** refers to the risk that an individual's presence in a dataset can be disclosed through the synthetic data.

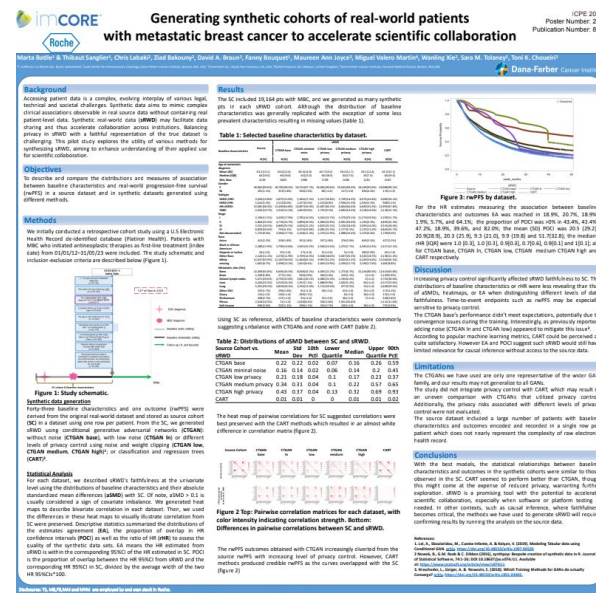Adversary gains information on the individual through the **membership in the dataset.**

Original sample

| Age | Gender | BMI | Hypertension |
|-----|--------|------|--------------|
| 58 | F | 32.4 | True |
| 66 | M | 25.4 | False |
| **35** | **M** | 27.8 | True |

Synthetic sample

| Age | Gender | BMI | Hypertension |
|-----|--------|------|--------------|
| **35** | **M** | 28 | True |
| 29 | F | 34.4 | True |
| 65 | M | 30.8 | False |

# Applications of synthetic data



Higher quality of the synthetic data is required

"I want to solve imbalances in my data"

"I want to facilitate a data sharing process"

"I need to use data for secondary purposes"

"I want to augment the size of my data"

"I want to test our software, model, pipeline or platform"

"I want to start building my solution before getting data access"

"I want to test a software before contracting"

# A real example of synthetic data sharing pilot

- Pilot on data sharing to accelerate collaboration with Dana Farber Institute (Harvard Medical School)

- Synthesis of sRWD from existing _analytical dataset_ using different methods CTGAN or CART

  - ~10k MBC patients

  - Over 100 variables

# A real example of synthetic data sharing pilot

- A total of 9,770 pts with MBC were included in the SC and as many synthetic pts were generated in each sRWD cohort.
- Distributions of continuous and categorical variables were closely replicated
- Measures of association between baseline characteristics and outcomes were largely preserved
- CART outperformed CTGAN
- Dimensionality of dataset has a big impact on utility

# Thank you

Roche

**Doing now what patients need next**