# Private Federated Machine Learning
# The EPI Project
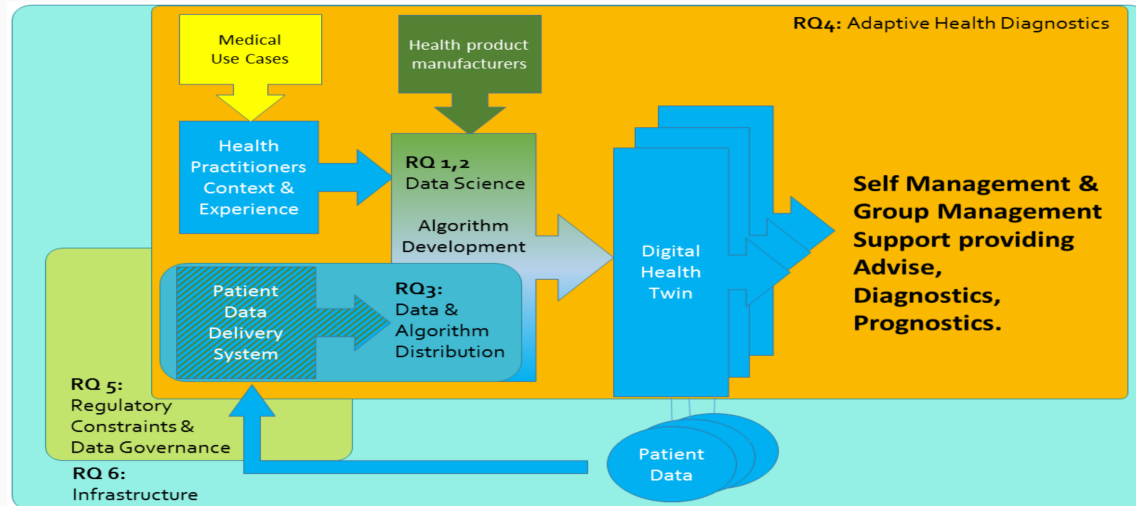
Saba Amiri

Adam Belloum, Sander Klous, Leon Gommans, Eric Nalisnick

UNIVERSITY OF AMSTERDAM
Informatics Institute

# Enabling Personalized Interventions
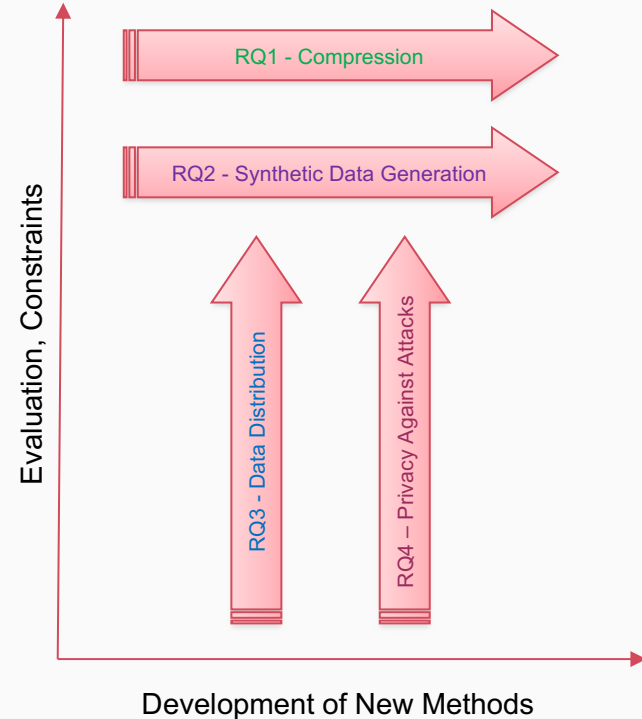
- EPI[*] project broadly aims to create a Digital Health Twin

  - The digital reflection of a person in terms of health related data and allows algorithms

  - Enables distributed processing of disparate relevant data, e.g. perform monitoring or predict outcomes of treatments



*https://enablingpersonalizedinterventions.nl/*

# Basic Setting

- RQ4-1 - How To Achieve Differential Privacy Through Compression?

- RQ4-2 - How to generate differentially-private synthetic tabular data in a distributed setting?

- RQ4-3 - What is the effect of non-i.i.d data distribution on the performance of differentially private machine learning models?

- RQ4-4 - How can we measure the privacy level of DP machine learning methods from the perspective of privacy attacks?
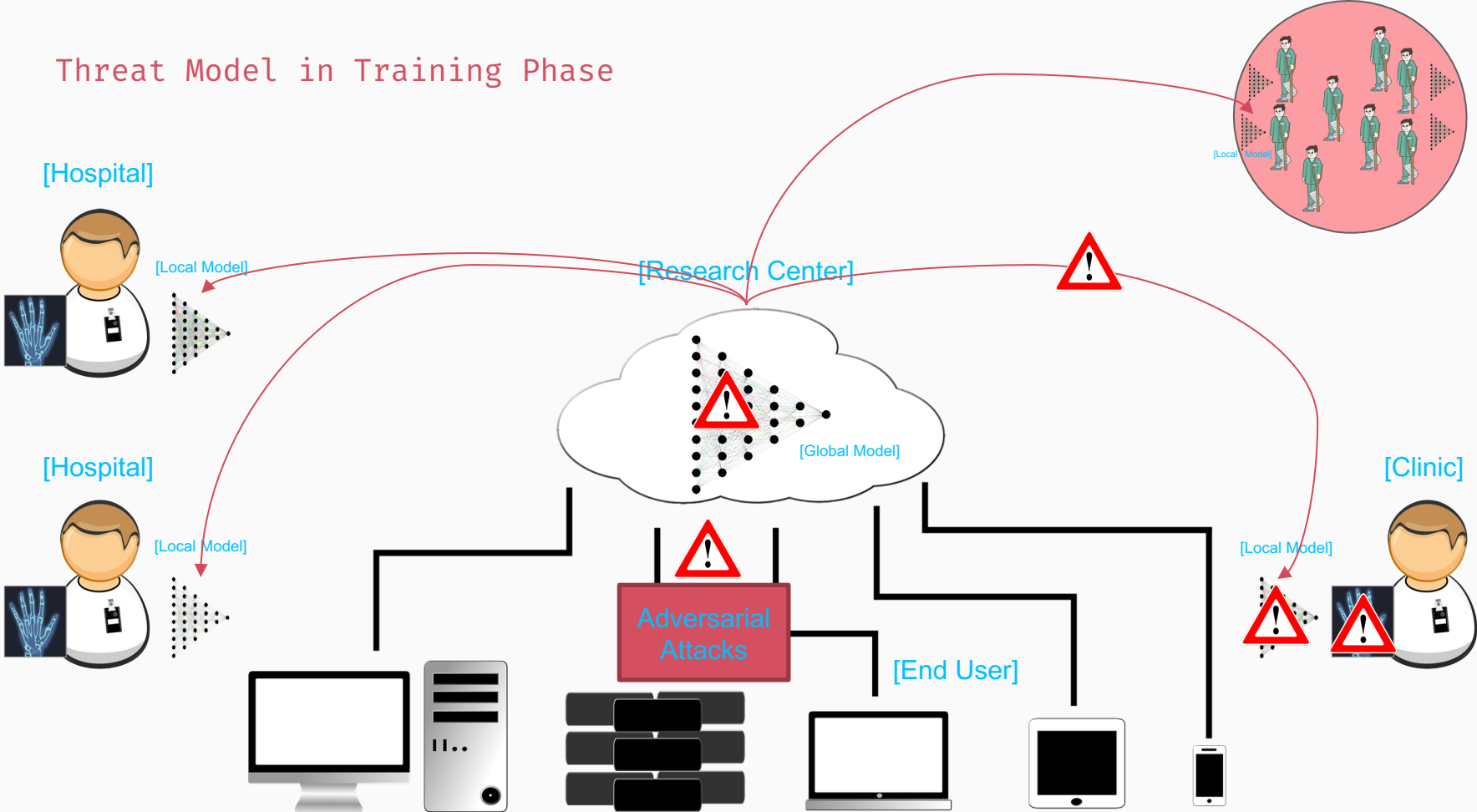


RQ1 - Compression

RQ2 - Synthetic Data Generation

RQ3 - Data Distribution

RQ4 – Privacy Against Attacks

Evaluation, Constraints

Development of New Methods

# Basic Setting

- Medical use-cases (EPI[*])

- Private, distributed, large datasets

- Common goal: train a machine learning model on these datasets while preserving privacy of the individuals in the datasets

# Basic Setting

- Medical use-cases (EPI[*])

- Private, distributed, large datasets

- Common goal: train a machine learning model on these datasets while preserving privacy of the individuals in the datasets

- Initial solution: accumulate data, train a centralized model

- Poses challenges, e.g. privacy, communication, etc.

Threat Model in Training Phase

[Hospital]

[Local Model]

[Research Center]

[Global Model]

[Hospital]

[Local Model]

[Clinic]

[Local Model]

Adversarial
Attacks

[End User]

UNIVERSITY OF AMSTERDAM
Informatics Institute

# Impact of non-i.i.d Distribution on Federated Learning

# The Problem w/ Federated Learning

- **Privacy**
  - FL solves the problem of data sharing
  - The training process is vulnerable
  - The model could leak information after being trained
- **Data distribution**
  - i.i.d assumption about data
  - 4 main types of imbalance in the data
    - Feature
    - Label
    - Temporal
    - Node
  - Has disparate impact on performance, fairness

# The Problem w/ Federated Learning

- **Adult dataset**

| | age | workclass | fnlwgt | education | education-num | marital-status | occupation | relationship | race | sex | capital-gain | capital-loss | hours-per-week | native-country | label |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 17 | Private | 124130 | Some-college | 9 | Separated | Protective-serv | Not-in-family | White | Male | 30 | 0 | 40 | Haiti | <=50K |
| 1 | 26 | Private | 168914 | HS-grad | 10 | Married-civ-spouse | Handlers-cleaners | Husband | Asian-Pac-Islander | Female | 21 | 1 | 39 | Yugoslavia | <=50K |
| 2 | 33 | Self-emp-not-inc | 218757 | HS-grad | 11 | Married-civ-spouse | Machine-op-inspct | Not-in-family | White | Male | 29 | 0 | 24 | United-States | >50K |
| 3 | 62 | Self-emp-not-inc | 558635 | Bachelors | 9 | Never-married | Prof-specialty | Wife | White | Male | 51 | 1 | 40 | United-States | <=50K |
| 4 | 27 | ? | 143612 | Masters | 13 | Separated | Priv-house-serv | Unmarried | White | Male | 89 | -2 | 40 | United-States | <=50K |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 995 | 44 | Private | 179779 | HS-grad | 9 | Never-married | Adm-clerical | Husband | White | Male | 2 | -3 | 40 | United-States | <=50K |
| 996 | 28 | Self-emp-not-inc | 180882 | Bachelors | 11 | Married-civ-spouse | Adm-clerical | Other-relative | Black | Female | 43 | 5 | 40 | United-States | <=50K |
| 997 | 15 | Private | 166548 | Bachelors | 6 | Married-civ-spouse | Protective-serv | Other-relative | White | Female | 23 | 7 | 38 | United-States | <=50K |
| 998 | 19 | Private | 158057 | Doctorate | 8 | Never-married | Other-service | Not-in-family | White | Male | 9 | -1 | 40 | United-States | >50K |
| 999 | 19 | Private | 119228 | Bachelors | 13 | Divorced | Other-service | Unmarried | White | Male | 69 | 5 | 40 | United-States | <=50K |

```
1 age                  2 workclass                3 education               4 education-num        5 marital-status
Min    17.             Private       22 696        HS-grad      10 501       Min    1.              Married-civ-spouse      14 976
1st Qu 28.             Self-emp-not-inc 2541       Some-college 7291         1st Qu 9.              Never-married           10 683
Median 37.             Local-gov      2093         Bachelors    5355         Median 10.             Divorced                4443
Mean   38.5816         ?              1836         Masters      1723         Mean   10.0807         Separated               1025
3rd Qu 48.             State-gov      1298         Assoc-voc    1382         3rd Qu 12.             Widowed                 993
Max    90.             Self-emp-inc   1116         11th         1175         Max    16.             Married-spouse-absent   418
                       (Other)        981          (Other)      5134                                Married-AF-spouse       23

6 occupation             7 relationship            8 race                           9 sex              10 capital-gain
Prof-specialty   4140     Husband        13 193    White            27 816          Male   21 790      1st Qu 0.
Craft-repair     4099     Not-in-family  8305      Black            3124            Female 10 771      3rd Qu 0.
Exec-managerial  4066     Own-child      5068      Asian-Pac-Islander 1039                             Median 0.
Adm-clerical     3770     Unmarried      3446      Amer-Indian-Eskimo 311                              Min    0.
Sales            3650     Wife           1568      Other            271                                Mean   1077.65
Other-service    3295     Other-relative 981                                                          Max    99 999.
(Other)          9541

11 capital-loss          12 hours-per-week         13 native-country        14 income
1st Qu 0.                Min    1.                 United-States 29 170      <=50K 24 720
3rd Qu 0.                1st Qu 40.                Mexico        643          >50K  7841
Median 0.                Median 40.                ?             583
Min    0.                Mean   40.4375            Philippines   198
Mean   87.3038           3rd Qu 45.                Germany       137
Max    4356.             Max    99.                Canada        121
                                                   (Other)       1709
```

UNIVERSITY OF AMSTERDAM
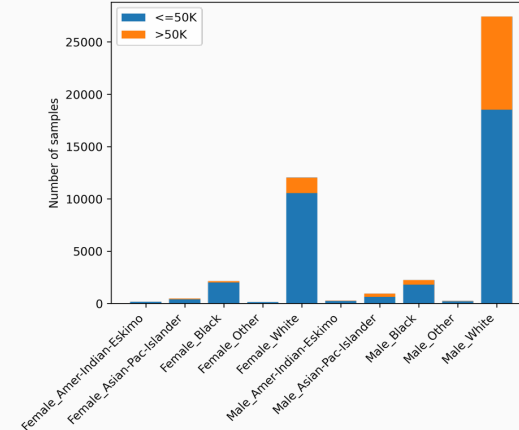Informatics Institute

# The Problem w/ Federated Learning
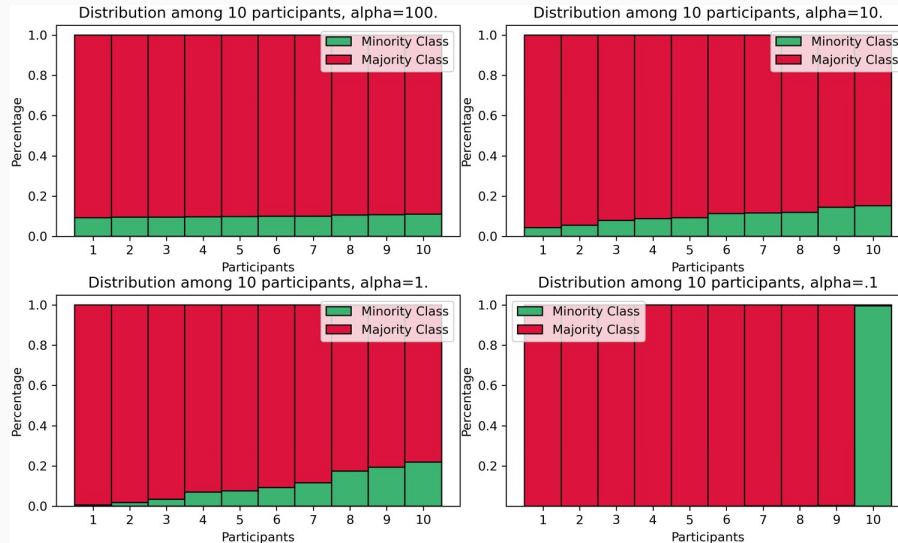
- **Adult dataset**

| | age | workclass | fnlwgt | education | education-num | marital-status | occupation | relationship | race | sex | capital-gain | capital-loss | hours-per-week | native-country | label |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 17 | Private | 124130 | Some-college | 9 | Separated | Protective-serv | Not-in-family | White | Male | 30 | 0 | 40 | Haiti | <=50K |
| 1 | 26 | Private | 168914 | HS-grad | 10 | Married-civ-spouse | Handlers-cleaners | Husband | Asian-Pac-Islander | Female | 21 | 1 | 39 | Yugoslavia | <=50K |
| 2 | 33 | Self-emp-not-inc | 218757 | HS-grad | 11 | Married-civ-spouse | Machine-op-inspct | Not-in-family | White | Male | 29 | 0 | 24 | United-States | >50K |
| 3 | 62 | Self-emp-not-inc | 558635 | Bachelors | 9 | Never-married | Prof-specialty | Wife | White | Male | 51 | 1 | 40 | United-States | <=50K |
| 4 | 27 | ? | 143612 | Masters | 13 | Separated | Priv-house-serv | Unmarried | White | Male | 89 | -2 | 40 | United-States | <=50K |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 995 | 44 | Private | 179779 | HS-grad | 9 | Never-married | Adm-clerical | Husband | White | Male | 2 | -3 | 40 | United-States | <=50K |
| 996 | 28 | Self-emp-not-inc | 180882 | Bachelors | 11 | Married-civ-spouse | Adm-clerical | Other-relative | Black | Female | 43 | 5 | 40 | United-States | <=50K |
| 997 | 15 | Private | 166548 | Bachelors | 6 | Married-civ-spouse | Protective-serv | Other-relative | White | Female | 23 | 7 | 38 | United-States | <=50K |
| 998 | 19 | Private | 158057 | Doctorate | 8 | Never-married | Other-service | Not-in-family | White | Male | 9 | -1 | 40 | United-States | >50K |
| 999 | 19 | Private | 119228 | Bachelors | 13 | Divorced | Other-service | Unmarried | White | Male | 69 | 5 | 40 | United-States | <=50K |

```
1 age                    2 workclass                      3 education                 4 education-num           5 marital-status
Min    17.               Private          22 696          HS-grad        10 501       Min     1.               Married-civ-spouse   14 976
1st Qu 28.               Self-emp-not-inc  2541           Some-college    7291        1st Qu 9.                Never-married        10 683
Median 37.               Local-gov         2093           Bachelors       5355        Median 10.               Divorced              4443
Mean   38.5816           ?                 1836           Masters         1723        Mean   10.0807           Separated             1025
3rd Qu 48.               State-gov         1298           Assoc-voc       1382        3rd Qu 12.               Widowed                993
Max    90.               Self-emp-inc      1116           11th            1175        Max    16.               Married-spouse-absent  418
                         (Other)            981           (Other)         5134                                 Married-AF-spouse       23

6 occupation                   7 relationship                  8 race                            9 sex                10 capital-gain
Prof-specialty   4140          Husband         13 193          White            27 816           Male    21 790       1st Qu 0.
Craft-repair     4099          Not-in-family    8305           Black             3124            Female  10 771       3rd Qu 0.
Exec-managerial  4066          Own-child        5068           Asian-Pac-Islander 1039                                Median 0.
Adm-clerical     3770          Unmarried        3446           Amer-Indian-Eskimo 311                                 Min    0.
Sales            3650          Wife             1568           Other              271                                 Mean   1077.65
Other-service    3295          Other-relative   981                                                                   Max    99 999.
(Other)          9541

11 capital-loss                12 hours-per-week               13 native-country               14 income
1st Qu 0.                      Min     1.                      United-States  29 170           <=50K 24 720
3rd Qu 0.                      1st Qu 40.                      Mexico           643            >50K   7841
Median 0.                      Median 40.                      ?                583
Min    0.                      Mean    40.4375                 Philippines      198
Mean   87.3038                 3rd Qu 45.                      Germany          137
Max    4356.                   Max     99.                     Canada           121
                                                               (Other)         1709
```
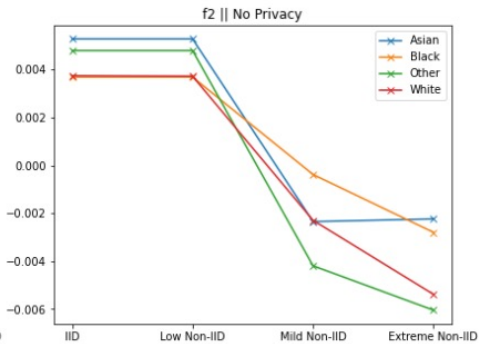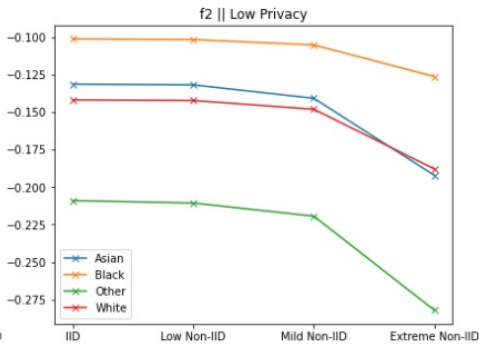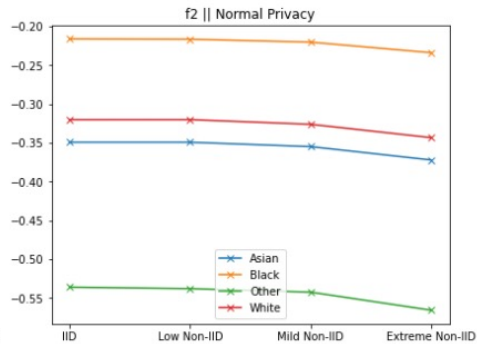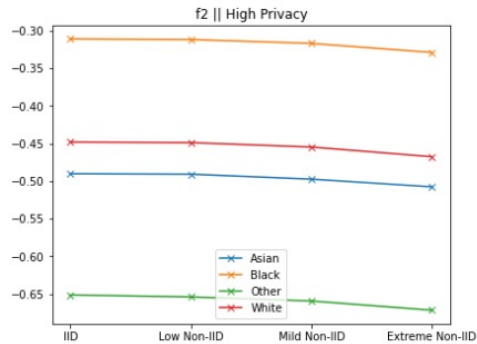
Number of samples of adult dataset per label per target class

# Our Research: Impact of non-i.i.d data on private FL

# Our Research: Impact of non-i.i.d data on private FL

# Differentially Private Synthetic Data Generation

# Our Research: Differentially Private Synthetic Data Generation

- Adding DP to ML models is costly
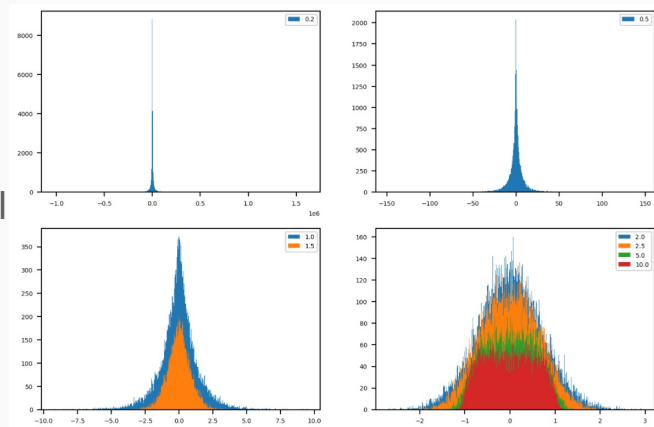
- Alternatively, we can make the data "privacy preserving"

- How?

  - Use a differentially private generative model to estimate the distribution of the data

  - Train the model on real data

  - Use model to generate a synthetic dataset

  - Due to post-processing theorem, any model trained on our synthetic data is at least differentially private with the same level as our generative model

# Our Research: Differentially Private Synthetic Data Generation

- Generate privacy preserving synthetic data from original data

- Differentially private with an acceptable privacy budget

- On tabular data

- Preserve statistical properties

- Maintain machine learning efficacy

- Distributed environment

- No i.i.d assumptions about data distribution

# Our Research: Differentially Private Synthetic Data Generation

- Generate privacy preserving synthetic data from original data

- Differentially private with an acceptable privacy budget

- On tabular data

- Preserve statistical properties

- Maintain machine learning efficacy

- Distributed environment

- No i.i.d assumptions about data distribution

- **Data quality: Semantic integrity**

# Our Research: Differentially Private Synthetic Data Generation

- Why semantic integrity?



Percentage of women and men with female condition

# Our Research: Differentially Private Synthetic Data Generation

- Proposed model's performance



MLP F-Scores

# Our Research: Differentially Private Synthetic Data Generation

- Proposed model's quality

**Synthetic data generation for data with long tailed distributions**

# Synthetic data generation for data with long tailed distributions

- Long tailed data

  - Have a generative model that is able to capture the tail behavior of long-tailed distributions

# Synthetic data generation for data with long tailed distributions

- Long tailed data

  - Have a generative model that is able to capture the tail behavior of long-tailed distributions

- Initial approach: GANs with a differentiable generalized Gaussian base distribution

# Synthetic data generation for data with long tailed distributions

- Second approach: normalizing flows

  - A normalizing flow

    - describes the transformation of a probability density through a sequence of invertible mappings.

    - Transforms a simple distribution into a complex one by applying a sequence of invertible transformation functions.

    - Flowing through a chain of transformations, we repeatedly substitute the variable for the new one according to the change of variables theorem and eventually obtain a probability distribution (i.e. normalized) of the final target variable

    - Normalizing flows can exactly estimate the density function

    - There is theory on capabilities of NFs on capturing the tail behavior of long-tailed distributions



$$\mathbf{z}_0 \xrightarrow{f_1(\mathbf{z}_0)} \mathbf{z}_1 \cdots \mathbf{z}_{i-1} \xrightarrow{f_i(\mathbf{z}_{i-1})} \mathbf{z}_i \xrightarrow{f_{i+1}(\mathbf{z}_i)} \cdots \mathbf{z}_K = \mathbf{x}$$

$$\mathbf{z}_0 \sim p_0(\mathbf{z}_0) \qquad \mathbf{z}_i \sim p_i(\mathbf{z}_i) \qquad \mathbf{z}_K \sim p_K(\mathbf{z}_K)$$
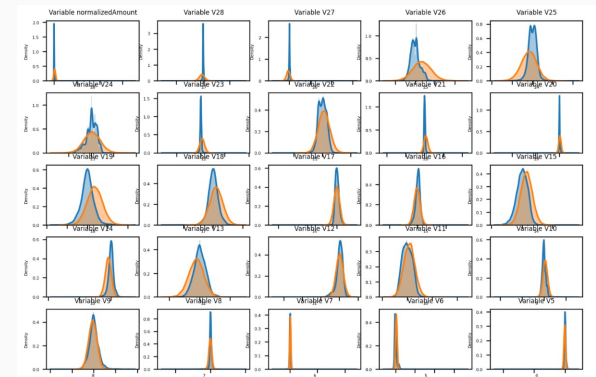
# Synthetic data generation for data with long tailed distributions

- Second approach: normalizing flows

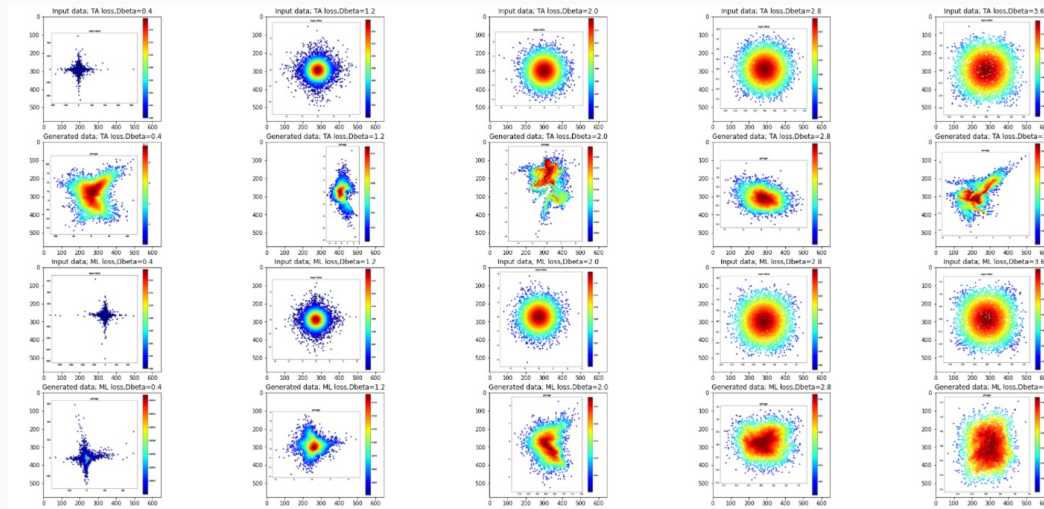  - How does a flow-based model compare to a GAN?

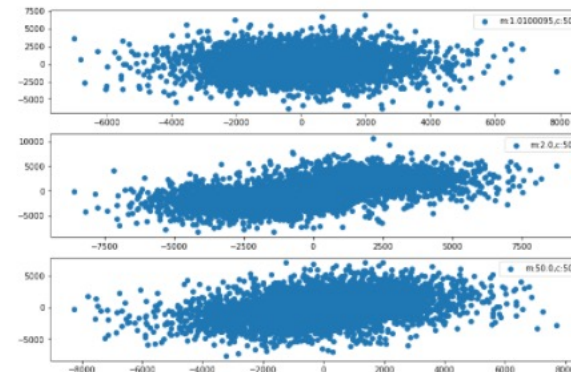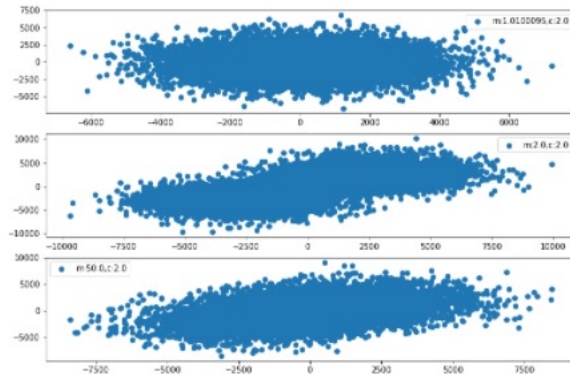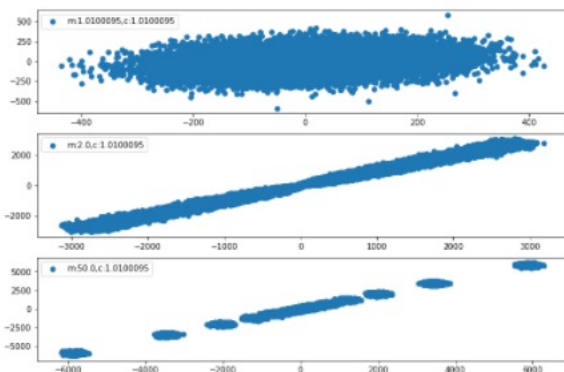- FLOW                                                    GAN

# Synthetic data generation for data with long tailed distributions

- Second approach: normalizing flows

  - Next step: dual training with ML-based training for the flow model and a loss function utilizing tail-adaptive alpha divergence for the base parameters



UNIVERSITY OF AMSTERDAM
Informatics Institute

# Synthetic data generation for data with long tailed distributions

- Last step: flexible mixture base distribution

  - Smooth contraction/expansion of the base mixture distribution to help the flow-based model capture the tail properties of the target distribution

# Thank You!