# Generating Synthetic Data for the NHS

**Synthetic Data Summit -**
**30th November 2023**


**Dr. Jonny Pearson**
Lead Data Scientist
**Digital Analytics and Research Team**
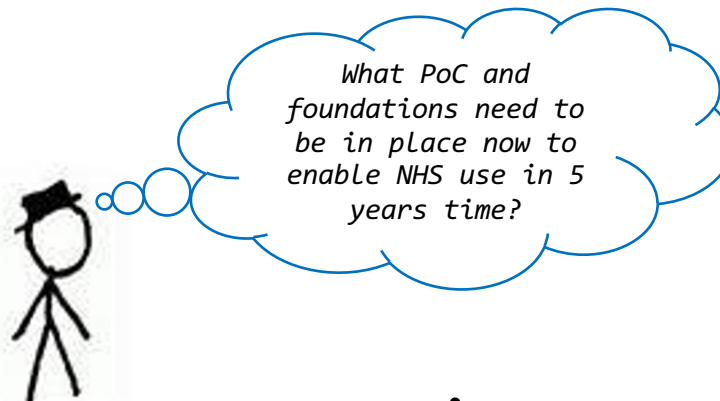NHS England
jonathanpearson@nhs.net

# Who We Are and Where We Fit

## Academia (and industry research)

Cutting edge models and thinking around getting value from large and complex data sources.

- Work often siloed and caught in long term projects
- Application often focussed on edge cases and ideal circumstances
- Sometime lacking real world or domain specific application

## DART Innovation

*What PoC and foundations need to be in place now to enable NHS use in 5 years time?*

Work with academia and NHS Ops to develop **both** push and pull (in different time scales)

- Short tangible outputs that clearly build towards wider context
- Take risks with expectation of high rate of failure
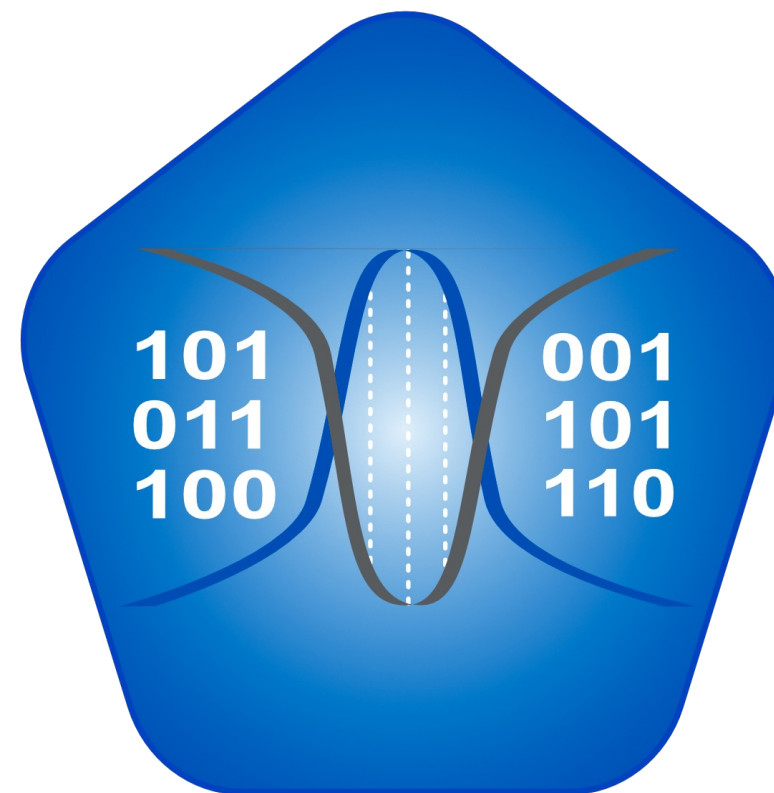- Include development cycle and tech transition plan

## NHS Operations and Decision Making

Need often driven by short-term priorities reducing desire for R&D.

- Evidence-based decision making from robust data insights
- Live modelling and visualisation of data to support daily operations
- Linking data across a complex landscape

## What's Coming

- **Fidelity - Simple is Often Better**

- **The Generation Landscape**

- **Our Approach**

- **Evaluating the Data**

- **What to care about (for healthcare)**
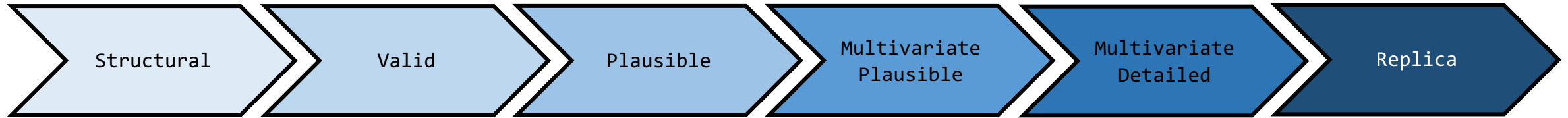


## Synthetic Data

# Setting the Scene

Synthetic data generation is not a silver bullet and often not an easy alternative, but it does have huge potential for datasets that are too low quality to use, too sensitive to share or, just doesn't exist.



**Raw Data**

**Reality**

Data Quality and Bias Issues

Model and Transformations

Representation

Fidelity, Fairness, Privacy Metrics

**Generated Data**

Use-Case and Amendments

# Fidelity - Simple is Often Better

**Range of Fidelity** *(how similar the generated data is to the ground truth)*

| Structural | Valid | Plausible | Multivariate Plausible | Multivariate Detailed | Replica |
|---|---|---|---|---|---|

*Source:* **Office for National Statistics**

**Range of Use-cases**



**End-to-end software testing**
e.g. Interoperability of architectures and systems to pass health data in FHIR formats

**Tool Demonstration**
e.g. New geospatial tool for showing impact of service planning on travel distance

**Faster Innovation**
e.g. Internal or external development of patient safety report classifier

**Novel Linkage**
e.g. generation of patient cancer pathways

**Evaluation of Solutions**
e.g. test clinical risk score prediction on rare patients.

**Addressing Bias and Quality**
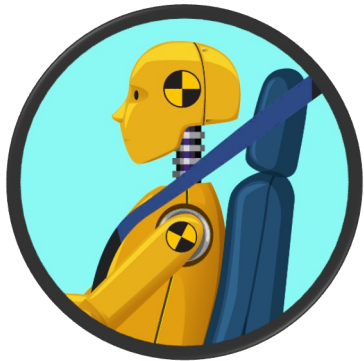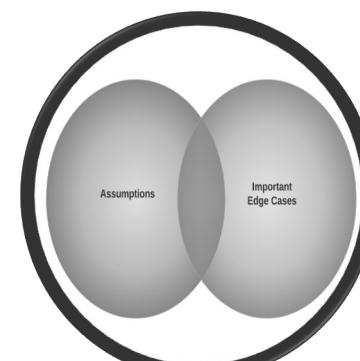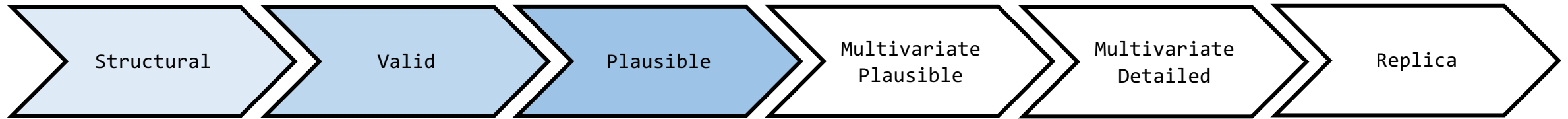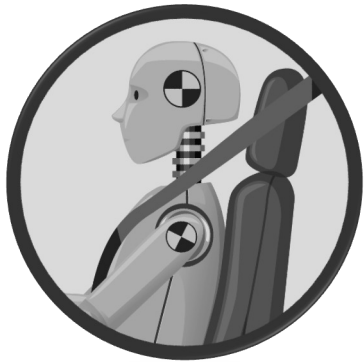e.g. creating a de-biased data set to highlight the impact that bias is having on the real data

174

# Fidelity - Simple is Often Better

**Range of Fidelity** *(how similar the generated data is to the ground truth)*

| Structural | Valid | Plausible | Multivariate Plausible | Multivariate Detailed | Replica |
|---|---|---|---|---|---|

*Source:* **Office for National Statistics**

## Range of Use-cases



**End-to-end software testing**

e.g. Interoperability of architectures and systems to pass health data in FHIR formats



**Tool Demonstration**

e.g. New geospatial tool for showing impact of service planning on travel distance
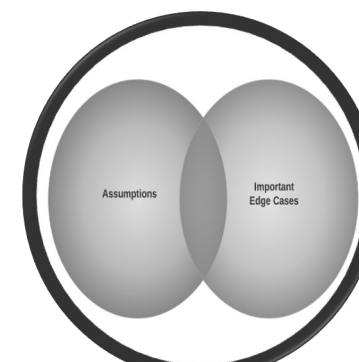


**Faster Innovation**

e.g. Internal or external development of patient safety report classifier



**Novel Linkage**

e.g. generation of patient cancer pathways



**Evaluation of Solutions**

e.g. test clinical risk score prediction on rare patients.
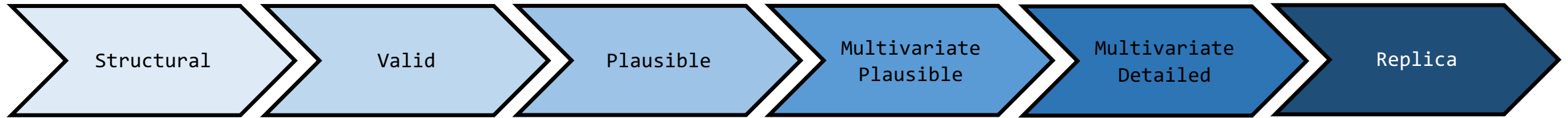


**Addressing Bias and Quality**

e.g. creating a de-biased data set to highlight the impact that bias is having on the real data

# Fidelity - Simple is Often Better

**NHS**

*Range of Fidelity* (how similar the generated data is to the ground truth)

| Structural | Valid | Plausible | Multivariate Plausible | Multivariate Detailed | Replica |

*Source:* **Office for National Statistics**

## *Range of Use-cases*



**End-to-end software testing**
e.g. Interoperability of architectures and systems to pass health data in FHIR formats

**Tool Demonstration**
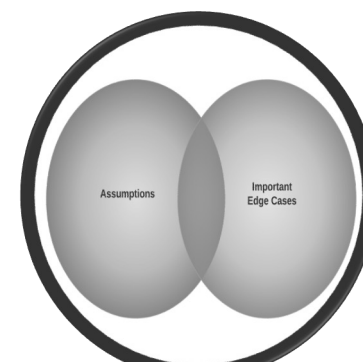e.g. New geospatial tool for showing impact of service planning on travel distance

**Faster Innovation**
e.g. Internal or external development of patient safety report classifier

**Novel Linkage**
e.g. generation of patient cancer pathways

**Evaluation of Solutions**
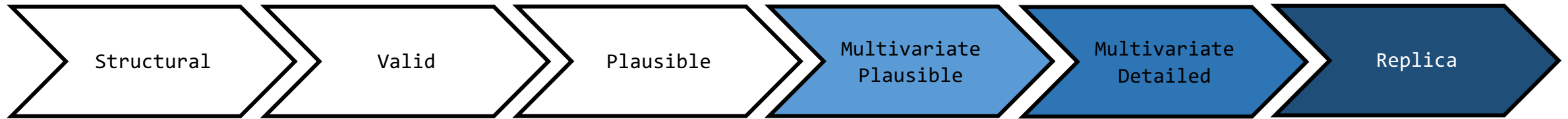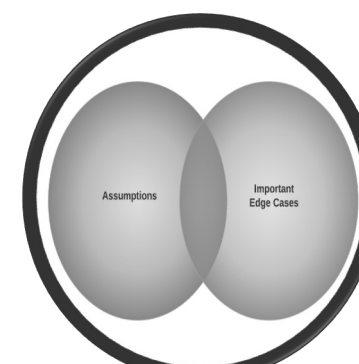e.g. test clinical risk score prediction on rare patients.

**Addressing Bias and Quality**
e.g. creating a de-biased data set to highlight the impact that bias is having on the real data

176

# Fidelity - Simple is Often Better

**Range of Fidelity** *(how similar the generated data is to the ground truth)*

| Structural | Valid | Plausible | Multivariate Plausible | Multivariate Detailed | Replica |

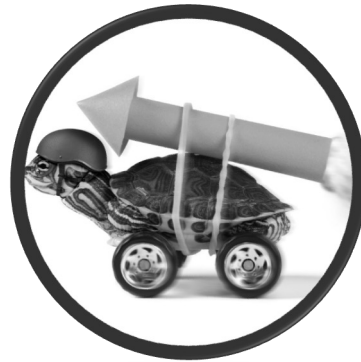*Source:* **Office for National Statistics**

## Range of Use-cases

**End-to-end software testing**
e.g. Interoperability of architectures and systems to pass health data in FHIR formats

**Tool Demonstration**
e.g. New geospatial tool for showing impact of service planning on travel distance
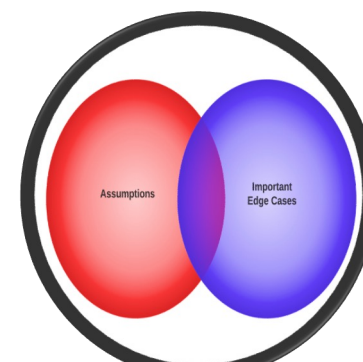
**Faster Innovation**
e.g. Internal or external development of patient safety report classifier

**Novel Linkage**
e.g. generation of patient cancer pathways

**Evaluation of Solutions**
e.g. test clinical risk score prediction on rare patients.

**Addressing Bias and Quality**
e.g. creating a de-biased data set to highlight the impact that bias is having on the real data

# Fidelity - Simple is Often Better

*Range of Fidelity* (how similar the generated data is to the ground truth)

| Structural | Valid | Plausible | Multivariate Plausible | Multivariate Detailed | Replica |
|---|---|---|---|---|---|

*Source:* **Office for National Statistics**

## Range of Use-cases



**End-to-end software testing**
e.g. Interoperability of architectures and systems to pass health data in FHIR formats



**Tool Demonstration**
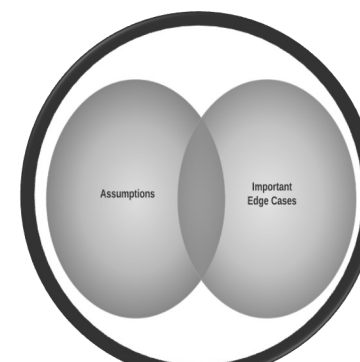e.g. New geospatial tool for showing impact of service planning on travel distance



**Faster Innovation**
e.g. Internal or external development of patient safety report classifier



**Novel Linkage**
e.g. generation of patient cancer pathways



**Evaluation of Solutions**
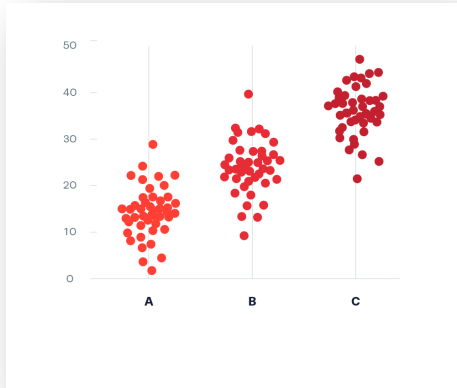e.g. test clinical risk score prediction on rare patients.



**Addressing Bias and Quality**
e.g. creating a de-biased data set to highlight the impact that bias is having on the real data

178

# Fidelity - Simple is Often Better

NHS

*Range of Fidelity* (how similar the generated data is to the ground truth)

| Structural | Valid | Plausible | Multivariate Plausible | Multivariate Detailed | Replica |
|---|---|---|---|---|---|

*Source:* **Office for National Statistics**

*Range of Use-cases*

| **End-to-end software testing** | **Tool Demonstration** | **Faster Innovation** | **Novel Linkage** | **Evaluation of Solutions** | **Addressing Bias and Quality** |
|---|---|---|---|---|---|
| e.g. Interoperability of architectures and systems to pass health data in FHIR formats | e.g. New geospatial tool for showing impact of service planning on travel distance | e.g. Internal or external development of patient safety report classifier | e.g. generation of patient cancer pathways | e.g. test clinical risk score prediction on rare patients. | e.g. creating a de-biased data set to highlight the impact that bias is having on the real data |

179

# Fidelity - Simple is Often Better

**Range of Fidelity** *(how similar the generated data is to the ground truth)*

| Structural | Valid | Plausible | Multivariate Plausible | Multivariate Detailed | Replica |
|---|---|---|---|---|---|

*Source:* **Office for National Statistics**

## Range of Use-cases



**End-to-end software testing**
e.g. Interoperability of architectures and systems to pass health data in FHIR formats

**Tool Demonstration**
e.g. New geospatial tool for showing impact of service planning on travel distance

**Faster Innovation**
e.g. Internal or external development of patient safety report classifier

**Novel Linkage**
e.g. generation of patient cancer pathways

**Evaluation of Solutions**
e.g. test clinical risk score prediction on rare patients.

**Addressing Bias and Quality**
e.g. creating a de-biased data set to highlight the impact that bias is having on the real data

180

# The Generation Landscape – By Technique



**Adding Noise / Data Erosion**

Adding Noise/Jitter

Suppression and relocation

Generalisation

[A Review of Anonymization for Healthcare Data](#)

**Statistical/Probabilistic models**

Sampling from independent marginals

Sampling from joint probabilities

*SynthPop*
*Faker*
*Simlacrum*
*CPRD Syntehtic Data Generation*

**Simulations**

Digital Twins

Clinical Practice guidelines (CPGs)

Agent based simulations

[Synthea](#), [simhospital](#)

**Perturbations of the manifold**

Synthetic Minority Over-Sampling Technique ([SMOTE](#))

Variational Autoencoders ([VAE](#))

[TVAE Synthetic Patient Generation](#)

**Iterative Comparisons**

Generative Adversarial Networks ([GAN](#))

GPT-3-based architecture

[CT-GAN](#)
PATE-GAN, ADS-GAN, DECAF
*[Van Der Schaar Lab](#)*
[SynGatorTron](#)

181

# The Generation Landscape – By Modality

# Our Approach - SynthVAE

Dom Danks and David Brind joined our team as a PhD Data Science Interns developing a variational autoencoder with differential privacy (SynthVAE).

## Fidelity



$$loss = \| x - \hat{x} \|^2 + KL[\ N(\mu_x, \sigma_x), N(0, I)\ ] = \| x - d(z) \|^2 + KL[\ N(\mu_x, \sigma_x), N(0, I)\ ]$$



The NHS AI Lab Skunkworks team published a case study of their use of our original tool in a user-friendly end-to-end process using QuantumBlack's Kedro.

## Privacy

The privacy is the most difficult element to quantise and demonstrate confidence in.

We investigated adding differential privacy through it's impact on privacy metrics from SDMetrics.



Figure 2: The DP-SGD algorithm. Credit: Abadi et al. (2016).

## Fairness



Fairness can be considered before, during or after model application.

We have explored incorporating manually adjusting a learnt DAG representation through the ability to meet different fairness metrics



Table 1: Table 1 from review by Mehrabi et al. [9] showing the breakdown of fairness metrics and segregating them into similar groups

# Our Approach – NHSSynth

In 2023, Harry Wilde has taken on the code with the task of bring everything together into a public facing python package called NHSSynth.

# Evaluating the Data - Quality



## Profile Comparisons
Are variable relationships maintained?

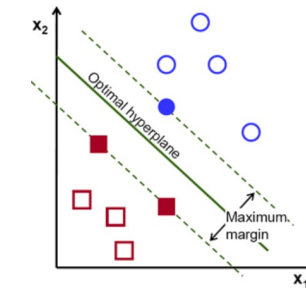e.g. Pearson's / **Similarity Metrics**

## Distribution Comparisons
Are variable profiles consistent?
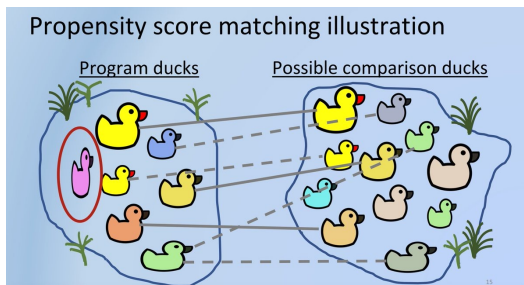
e.g. KS test, Chi-Squared

## Detection Metrics
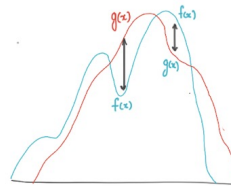Can a classifier differentiate between real and synthetic data?

e.g. **logistic**, SVM

## Variance Metrics
Is the range of variation of data points consistent?

Propensity score matching illustration

Program ducks    Possible comparison ducks

**Voas-Williamson** or propensity scores

## Aggregate difference Metrics
How much work is required to align data points to expectations?

e.g. **KL divergence, Gower distance** or Wasserstein metric

## Off-Manifold and Latent Space Checks
Are there unexpected features in the latent space?

e.g. PCA, t-SNE

**End use case comparisons and task performance**

# Evaluating the Data - Privacy

"Synthetic Data – Anonymisation Groundhog Day" by Stadler, Oprisanu, and Troncoso proposed using shadow modelling to apply generalised membership inference and attribute inference attacks on any synthetic data model.

## Tapas: a Toolbox for Adversarial Privacy Auditing of Synthetic Data

Project is under active development.  A Python library for evaluating the privacy of synthetic data from an adversarial perspective.

# Evaluating the Data - Fairness

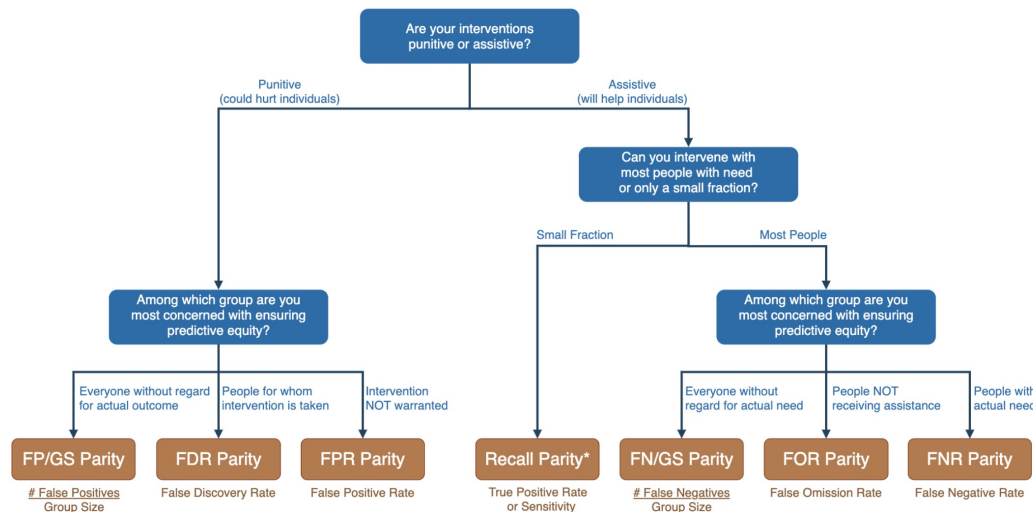It is desirable to be able to "de-bias" originally unfair training data or at least understand if the synthetic process has changed bias in the data.

**Aequitas:** an open-source bias audit toolkit for machine learning developers, analysts, and policymakers to audit machine learning models for discrimination and bias, and to make informed and equitable decisions around developing and deploying predictive risk-assessment tools. Developed by the Centre for Data Science and Public Policy at University of Chicago.



FAIRNESS TREE
(Zoomed in)



Bias and Fairness Audit Report

# Success / Failures

NHS

## Success

## Failures (yet to succeed)

+ Data generation using Bayesian Network, VAE, GAN & Transformer approaches

+ Development of end-to-end framework allowing experiments beyond single generation

+ IG Buy-In for data erosion approach

+ Partial success on data transformers for categorical, date and numerical variables

- Clarity around privacy threshold for different data types and sensitivities (characterising metrics as confidence levels)

- Clear privacy metrics

- Longitudinal Data generation

- Multi-table data generation

- Multimodal Data generation

- Implementing causal modelling into the generation algorithm

**Further examples required in this area**

**Further research required in this area**

# Main Take-Away

The balance for utility/privacy is key for the actual generation of synthetic data but two components which are fundamental to a successful project are:

## Explainability

*A project which is clear in how it has handled the data, where the data is limited and what level of risk is present is more useful that high quality or absolute privacy.*

## Adoption

*Putting user need and route to deployment above technical opportunity*

**Final Point:**

For healthcare a golden thread needs to be discussed throughout the work that allows the end decision to be set in context of the level or quality/risk created by the synthetic generation