



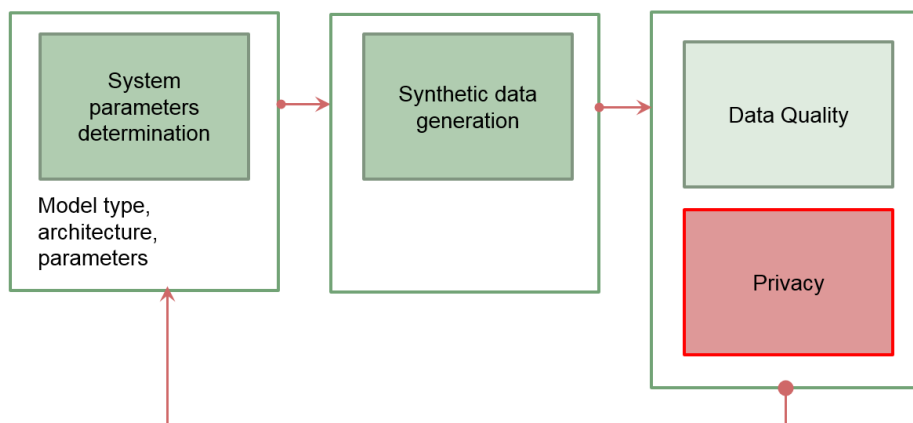
C. Joshi\*, I. Kaloskampis\*, D. Pugh\*, A. Noyvirt, L. Benedikt, L. Nolan

(\*) these authors contributed equally to this work

## Introduction

We develop a system which generates synthetic data to replace real data for processing and analysis. This is useful when the real data is sensitive (e.g. microdata, medical records, defence data).

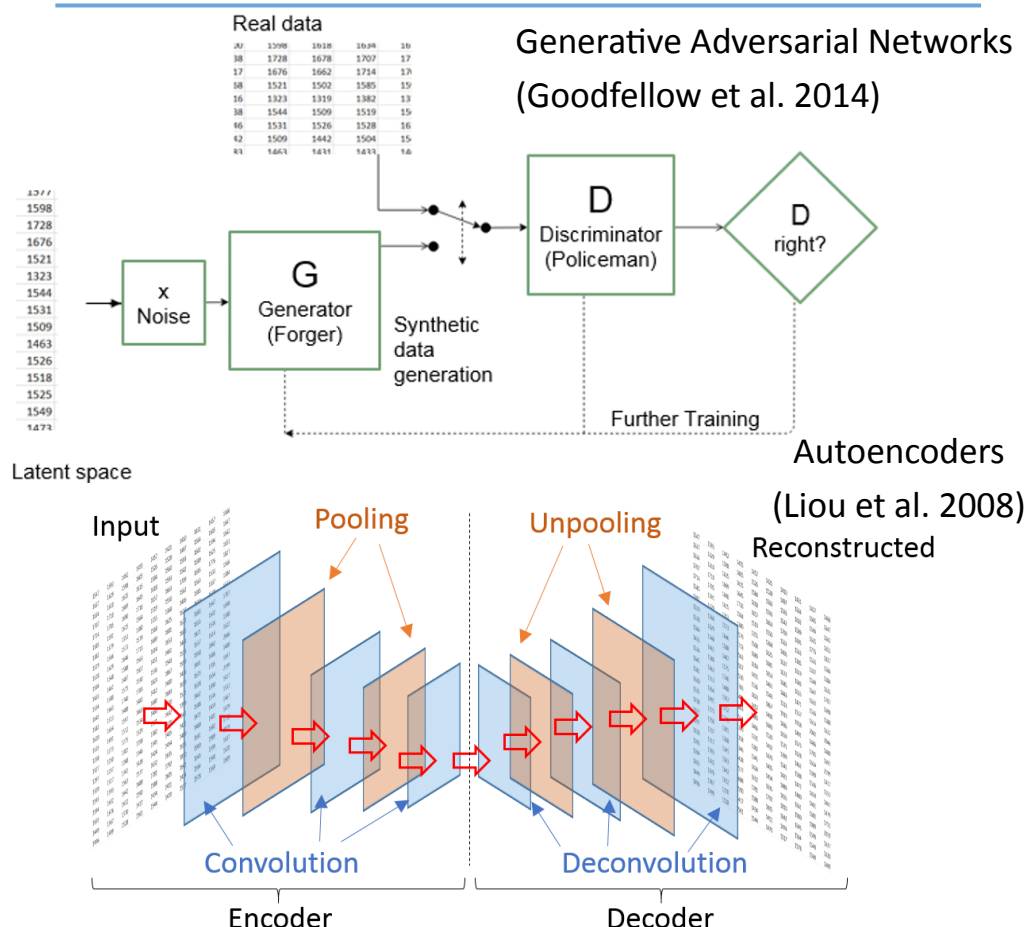
- ▶ Generative Adversarial Networks (GAN)
- ▶ (Variational) Auto-encoder
- ▶ Auto-regressive models
- ▶ Conventional generators (SMOTE etc.)



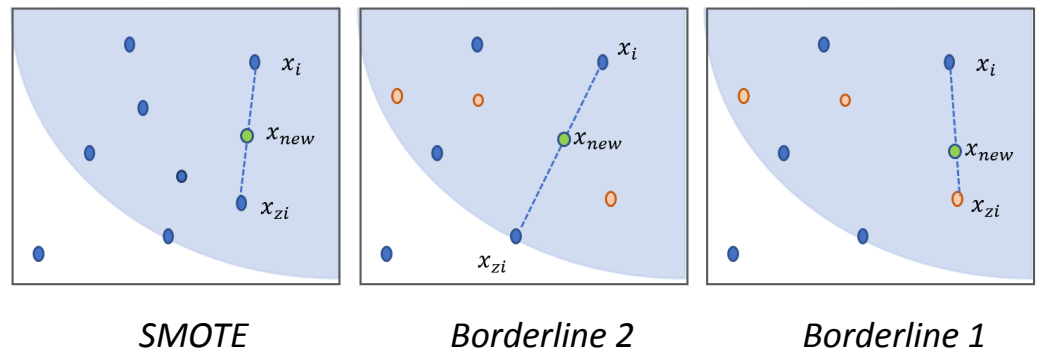
Our system utilises robust statistical, machine learning and state of the art deep learning techniques.

Safer, easier and faster sharing of sensitive data between research communities and ONS.

## Models

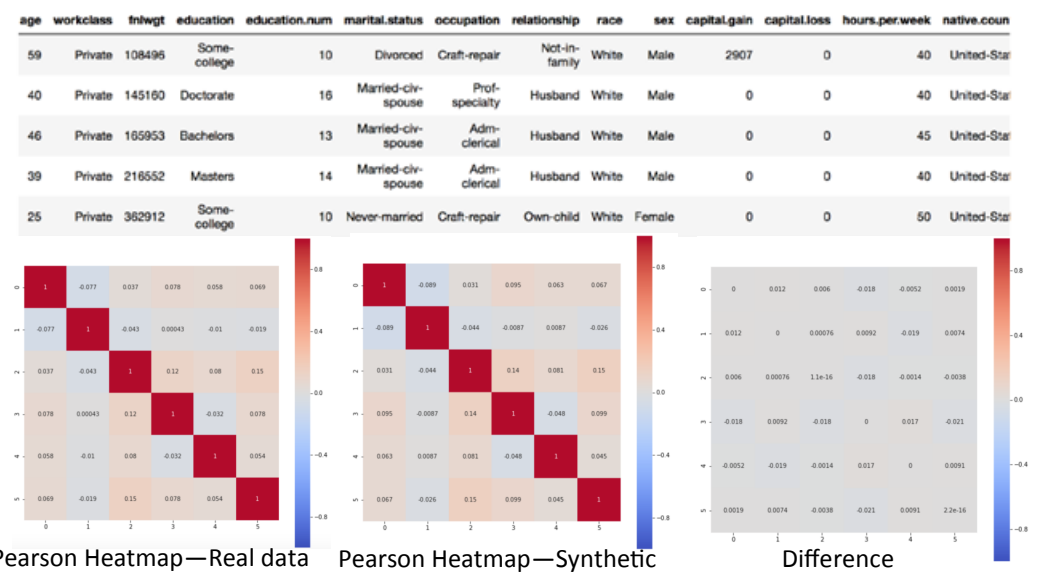


SMOTE (Chawla et al. 2002)

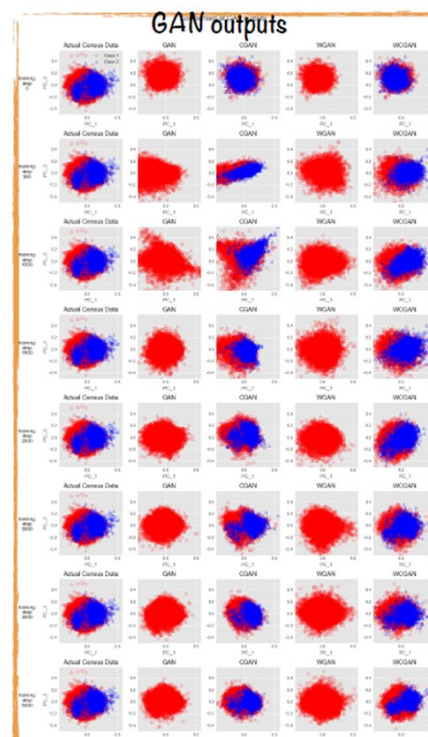


## Results

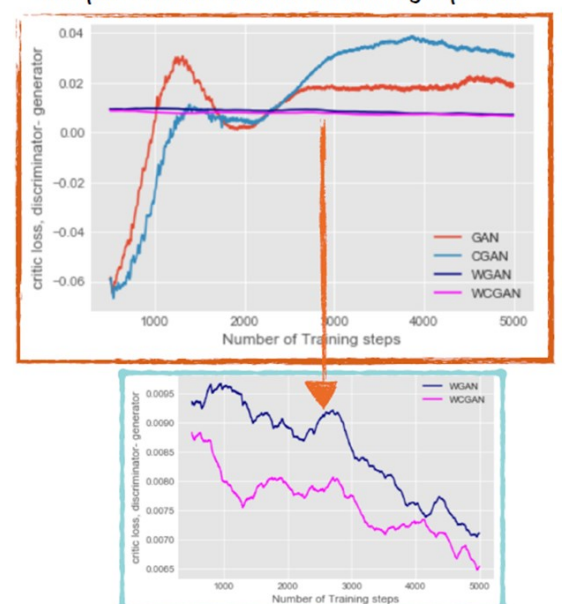
Adult Census Dataset (UCI Repository)



Model	Dataset	Application	Testing for...	Result
GAN	OCS Adult Census (numerical variables)	Synthesis	Data Quality – Pearson Correlation	0.085802
WGAN				0.064000
SMOTE				0.011192
SMOTE B1				0.011192
SMOTE B2				0.009105
ADASYN				0.012401
Autoencoder				0.008293
Variational Autoencoder			absolute mean difference between real & synthetic correlation matrices	0.211314



Compare the loss function as training improves



WGAN, WCGAN further training seems to improve the results in contrast to GAN and CGAN

