# Improving Structured Data ML Process with GANs
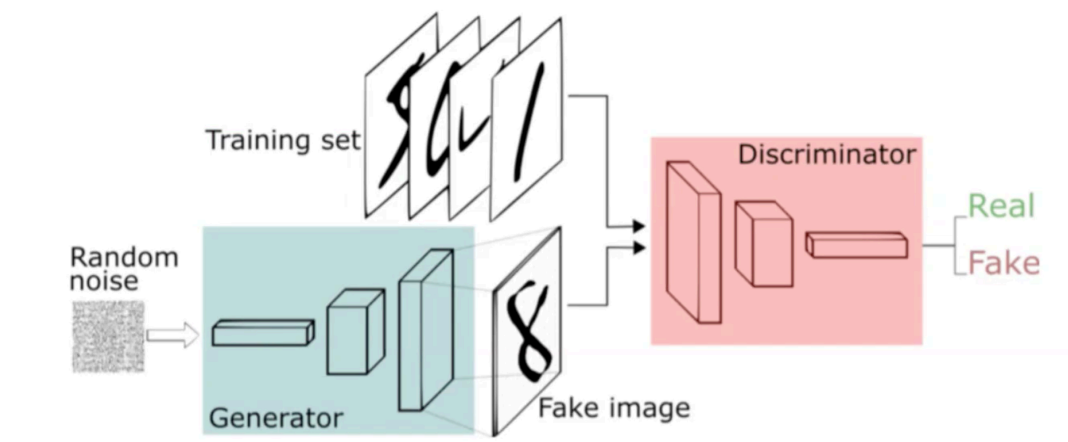
3/30/21     11:10 AM

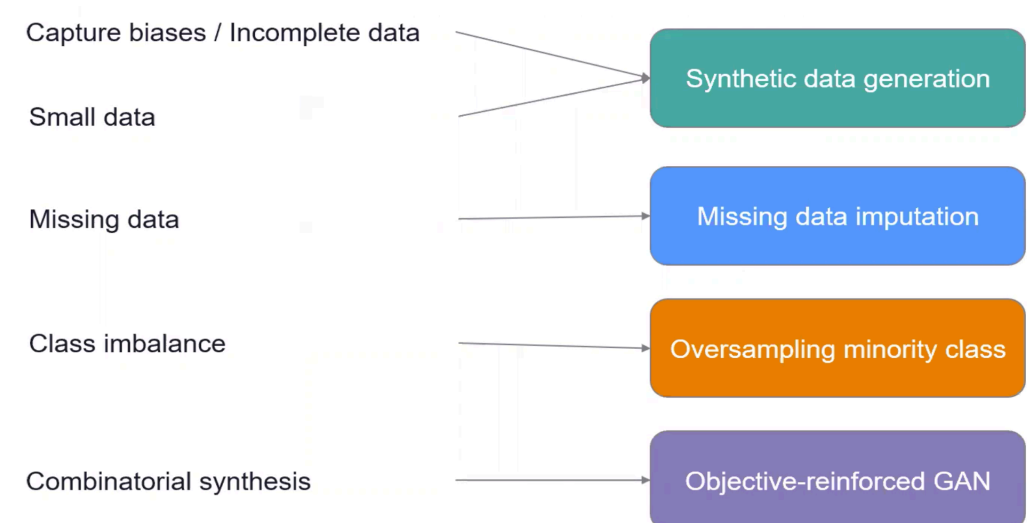*Speaker: Srinivas Chilukuri (ZS Associates)*
*Date: 03/30/2021*

## Structured datasets have unique challenges of their own
- Capture biases
- Missing data
- Small data
- Class imbalance

## GAN overview for image generation



## Applications of GANs



## Synthetic generation of tabular data
- Generator create values according to schema of columns
- Discriminator produces log probability

## Domain specific tabular GANs
- E-commerce:
  - Generate low-dimensional dense representation of e-commerce orders
  - To generate plausible orders data for new products
- Insurance: (tackle data privacy issue)
  - Generate accessible insurance datasets for actuarial studies
  - Tweaked CTGAN for imbalance categorical levels
- Medical (medGAN)
  - Generate privacy preserving, accessible healthcare datasets
  - Use of recurrent Encoders/Decoders for handling sequences
  - Does not handle lab values, clinical notes, diagnostic images, …
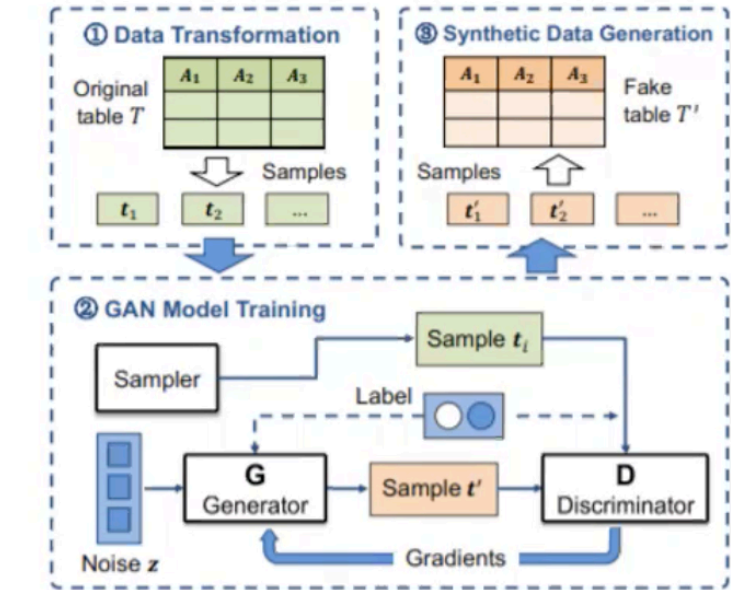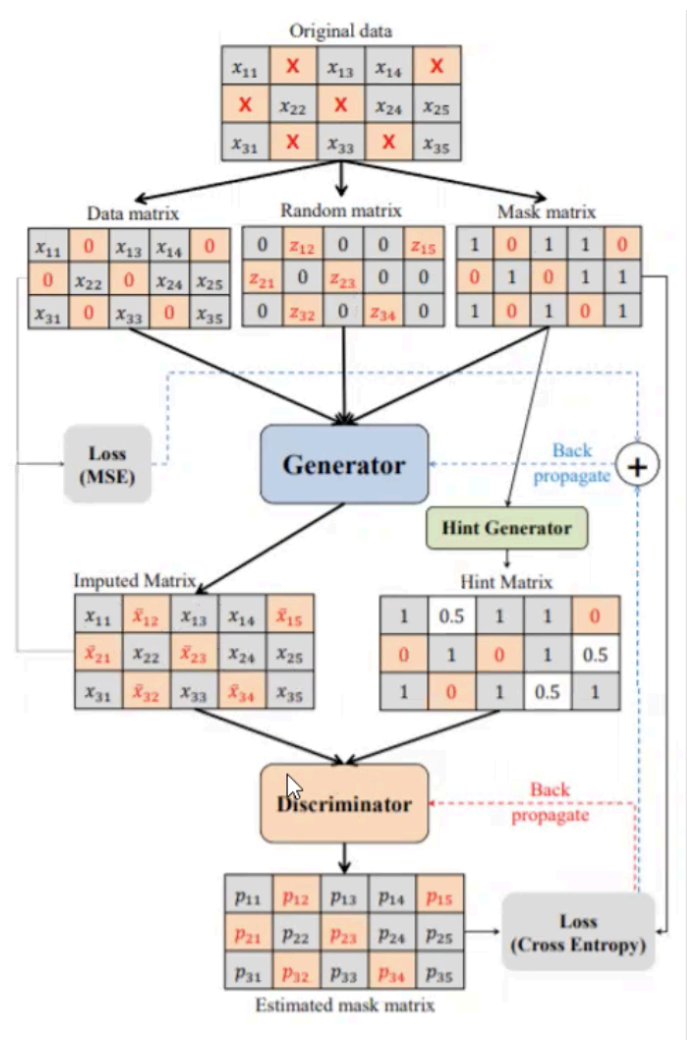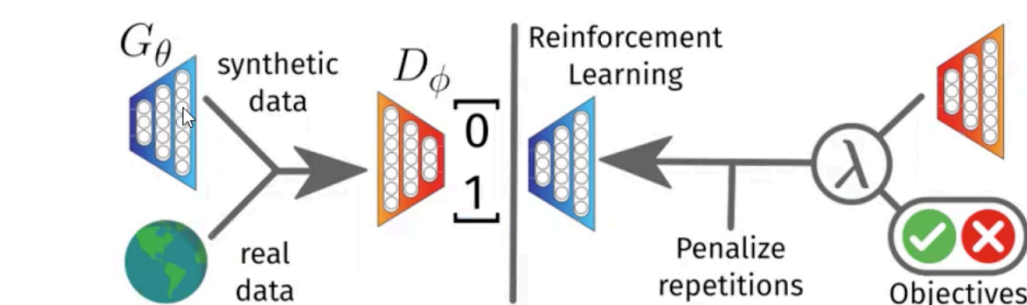- Relational datasets



Figure 2: Overview of data synthesis using GAN. (1) It transforms each record in a relational table into a sample $t \in \mathbb{R}^d$. (2) It takes the samples as input to train a deep generative model $G$ using the adversarial training framework in GAN. (3) It utilizes the trained $G$ to generate a set of synthetic samples, which are then transformed back into fake records.

## Missing data imputation
- Through masking data points in original dataset



## Generate data based on an objective



Figure 1: Schema for ORGAN. *Left:* $D$ is trained as a classifier receiving as input a mix of real data and generated data by $G$. *Right:* $G$ is trained by RL where the reward is a combination of $D$ and the objectives, and is passed back to the policy function via Monte Carlo sampling. We penalize non-unique sequences.

## Future directions
1. **How to evaluate tabular GANs?**
   Sample similarity
   - Basic statistics
   - Column correlations
   - Mirror column associations
   - PCA variance correlation
   Machine learning efficacy
   - Evaluate how downstream ML model performance varies with real data, synthetic data, real data + synthetic
   Privacy evaluations
   - Check how many samples are replicated from the training (real) dataset
   Human evaluations
   - On the lines of Turing test - to see whether humans can distinguish between real and synthetic data
   - In domain specific implementations, this will help with understanding which patterns GANs are able to learn
2. **Open challenges**
   - Data hungry - Adaptive Discriminator Augmentation
   - Mode collapse for imbalanced categorical data
   - Domain specific data still requires significant work e.g. HER data is still not realistic