# REPORT-3

# (INTERNSHIP)

**PRESENTED TO: Mrs SHOBHARANI MA'AM (BEL)**

**FROM: BHOOMIKA K (AIML, MS RAMAIAH)**

# CONTENTS:

# WHAT IS TOOLWORK?

In the context of synthetic data generation, a "tool" typically refers to a software or a set of algorithms designed to create artificial data that mimics the characteristics of real-world data. The goal of synthetic data generation tools is to generate data that is statistically similar to the original data but does not contain any sensitive or private information. These tools are commonly used in situations where obtaining real data is challenging, expensive, or involves privacy concerns.

# USING TOOLWORK TO SYNTHETISE DATA.

Synthetic data tools work by creating artificial datasets that mimic the statistical properties, patterns, and structures of real-world data. These tools use various techniques and algorithms to generate data that is representative of the underlying characteristics of the domain for which the synthetic data is being created. Here's a general overview of how synthetic data tools work:

- Understanding the Real Data:

The process often starts with a thorough analysis of the real data that the synthetic data is intended to simulate. This involves understanding the statistical distributions, relationships, and patterns present in the original data.

- Modeling Data Generation:

Synthetic data tools employ mathematical models and algorithms to capture the key features of the real data. Different techniques may be used, such as rule-based systems, probabilistic models, or machine learning approaches like Generative Adversarial Networks (GANs) or Variational Autoencoders (VAEs).

- Parameterization:

The tools are typically parameterized based on the characteristics observed in the real data. This involves setting parameters that control the distribution of values, relationships between variables, and other properties of the synthetic data.

- Randomization and Variability:

To make the synthetic data realistic, tools often introduce randomization and variability within the specified constraints. This ensures that the synthetic dataset is not an exact replica of the original data but still captures its essential characteristics.

- Privacy-Preserving Techniques:

In cases where privacy is a concern, synthetic data tools may incorporate privacy-preserving techniques to ensure that individual records in the synthetic dataset cannot be linked back to specific individuals in the real data. This can involve techniques such as differential privacy or data anonymization.

- Validation and Quality Control:

Synthetic data generation is an iterative process. The quality of the synthetic data is assessed by comparing it to the real data using various metrics. This validation step helps refine the parameters and improve the fidelity of the synthetic dataset.

- Application-Specific Customization:

Depending on the use case, synthetic data tools may allow for customization to meet specific application requirements. For example, in machine learning, synthetic data may be generated with a focus on certain features or patterns relevant to the training of a particular model.

- Output Generation:

The final output of a synthetic data tool is a dataset that can be used for various purposes, such as testing and development, model training, or data augmentation. This synthetic dataset should share key statistical properties with real data while avoiding the disclosure of sensitive information.

It's important to note that the effectiveness of synthetic data depends on the quality of the models and algorithms used, as well as the accuracy of the parameterization. Additionally, rigorous validation processes are crucial to ensuring that the synthetic data adequately represents the characteristics of the real data.

# Best Synthetic Data Generation Tools.



## 01. MDClone

Due to many privacy considerations, evaluating actual patient data is frequently difficult in the healthcare industry. However, such issues are no longer an issue. MDClone is a synthetic data generator designed exclusively for healthcare professionals to generate as much clinical data as you require from real patient profiles.

MDClone provides a systematic way to access healthcare data for research, synthesis, and analytics while avoiding the disruption of sensitive data. It can produce synthetic data from any sort of organized or unstructured patient-oriented data without revealing the patient's identity.

## 02. MOSTLY AI

MOSTLY.AI provides the most accurate synthetic data. It lets you unlock, share, update, and simulate data. MOSTLY.AI employs the most advanced artificial intelligence or AI model to generate fake data that looks and feels just like actual data. You will be able to keep valuable, granular-level information while ensuring that no individual is ever exposed.

MOSTLY.AI supports a wide range of data types, including structured data, text, pictures, and time series data. You can use it in a wide range of sectors and use cases. This versatility makes it suitable for a vast array of industries and applications.

## 03. Hazy

Hazy sets itself apart from the competition by offering models capable of generating top-quality synthetic data while incorporating a differential privacy mechanism. Whether your data is tabular, sequential, or spread across multiple tables in a relational database, Hazy has you covered.

Hazy's innovative data modeling approach empowers you to accelerate analytics workflows without the inherent risks associated with collecting real customer data. With Hazy, you can confidently develop and test your analytics solutions while safeguarding sensitive information.

## 04. Ydata

YData offers a data-centric platform that accelerates development and maximizes the ROI of your AI solutions. With YData, You can improve the quality of your training datasets and make them more robust and effective. Data scientists can use automated data quality analysis and cutting-edge synthetic data generation techniques to improve the performance of your dataset.

When it comes to data quality, YData goes the extra mile. It provides high-quality synthesized data and assures that it is free from bias or any personally identifiable information, which protects your privacy and compliance.

## 05. BizDataX

Whether you work as a test data engineer, bank professional, security officer, or business or data analyst, BizDataX gives you the tools you need to use synthetic data generation to protect personally identifiable information (PII) in your pre-production environment.

You can feel confident that you are in compliance with GDPR rules when you use BizDataX. The platform includes comprehensive data masking algorithms to ensure that sensitive data is secured throughout your testing and analysis procedures.

## 06. Sogeti

Sogeti is a cognitive-based tool that can help you with generating fake data. It is known as one of the most effective synthetic data generation tools, particularly for engineering, research, quality assurance, and testing.

You will benefit from Sogeti's Artificial Data Amplifier (ADA) technology, which has the unique capacity to read and reason with data of any type. It is a synthetic structured data generator that also creates unstructured data. ADA uses deep learning techniques to recreate its recognition capabilities, distinguishing it from its competitors.

## 07. Gretel

Gretel.ai is a new synthetic data generation tool for creating synthetic data. Gretel is a self-proclaimed "Privacy Engineering as a Service" that builds statistically similar datasets without using any sensitive customer data from the original source.

Gretel's ML method compares real-time information by employing a sequence-to-sequence model to enable prediction while generating fresh data and training the data for synthesis. Gretel also

employs differential privacy, which ensures that no original data is memorized or re-identified in the system.

## 08. Tonic

Tonic.ai offers an automated and anonymous data creation method for testing and development needs. With Tonic's technology, you can rest assured that your data remains anonymous through the use of database de-identification. This process separates PII from real data and prioritizes your client's privacy.

Tonic's powerful AI system categorizes distinct tables across databases using the Generative Adversarial Network, or GAN model. The platform preserves behaviors and dependencies within the data and allows the data science team to work with equally valuable data by eliminating hours of manual work.

## 09. CVEDIA

CVEDIA is an excellent alternative for a powerful computer vision cross-industry platform. The platform can generate synthetic data to power its AI and machine learning algorithms, and it does it effectively. CVEDIA's patented simulation engine, SynCity, allows it to generate high-quality synthetic data, which is extremely useful for testing and training models based on neural network architectures.

CVEDIA has you covered whether you work in security, manufacturing, or aerospace. Through NVIDIA's Metropolis initiative, the platform provides a holistic solution that addresses both your hardware and software requirements.

## 10. OneView

OneView is a scalable, cost-effective synthetic data solution for accelerating remote sensing imaging analytics. The platform provides synthetic data solutions and generates virtual synthetic datasets for you to employ in the training of machine learning models.

With OneView, you can avoid the time-consuming process of collecting, categorizing, and evaluating real-world photos from drones, aircraft, and satellites. The platform can generate datasets customized to suit your individual needs for any environment, object, or sensor.

## 11. Faker

Faker is a Python library that provides a simple and flexible way to generate fake data. It is widely used for creating synthetic datasets for various purposes, including testing, development, and training machine learning models. The primary goal of Faker is to generate data that looks real but doesn't correspond to any actual individuals or entities.

# CODE REPRESENTATION.

## Faker:

Here are some key features and aspects of Faker:

### Data Types:

Faker supports a wide range of data types, allowing you to generate fake data for names, addresses, phone numbers, dates, times, email addresses, and more.

Examples include name.first_name(), address.street_address(), phone_number(), date_of_birth(), etc.

### Localization:

Faker allows you to generate data that is localized for different regions and languages. This is particularly useful when you need data that resembles a specific cultural or linguistic context.

You can set the locale using Faker.seed() or by creating an instance of the Faker class with a specific locale.

### Customization:

You can customize the generated data by providing your own patterns or by using specific methods provided by Faker. This allows you to tailor the data to your specific needs.

For example, you can use pystr_format to create custom patterns for generating strings.

### Reproducibility:

Faker allows you to set a seed for the random number generator, ensuring that the same dataset can be regenerated when needed. This can be useful for reproducibility in testing or research scenarios.

Faker.seed() or using a seed when creating an instance of the Faker class achieves this.

### Usage in Python:

To use Faker in Python, you need to install it first using a package manager like pip: pip install faker

After installation, you can import the library and start using it in your Python script or Jupyter notebook.

```python
from faker import Faker

fake = Faker()

# Generate fake data
print(fake.name())
print(fake.address())
print(fake.email())
```

Integration with Other Libraries:

Faker can be easily integrated with other Python libraries and frameworks. For example, you might use it in conjunction with Pandas to create synthetic datasets for data analysis or machine learning.

Overall, Faker is a versatile and user-friendly tool for generating synthetic data quickly. It is widely adopted in the Python community and is particularly useful in scenarios where you need realistic-looking data for development and testing purposes.

Code:

```python
# import the Faker library
from faker import Faker

# create an instance of the Faker class
fake = Faker()

# generate synthetic data
for _ in range(5):  # adjust the number based on how many records you want
    # use Faker methods to generate different types of data
    name = fake.name()
    address = fake.address()
    email = fake.email()
    phone_number = fake.phone_number()

    # print or use the generated data as needed
    print(f"Name: {name}, Address: {address}, Email: {email}, Phone: {phone_number} \n")
```

Output:

```
Name: Andrew Cooper, Address: 95704 Amber Manors
Lake Kevin, MO 42927, Email: kristinethomas@example.net, Phone: 001-834-938-1255x3840

Name: Erika Carter, Address: 0970 Howard Bridge Apt. 910
Hoffmanview, CA 15264, Email: briannacooper@example.net, Phone: (398)785-8855x510

Name: Jesus Stewart, Address: 645 Bonnie Unions Apt. 420
Bruceton, AS 80133, Email: nicholas61@example.net, Phone: 782.732.9302x6041

Name: Lauren Thomas, Address: 30797 Gibson Crossroad
Kylebury, WI 34133, Email: ivega@example.org, Phone: 001-915-798-2039x8171

Name: Beverly Taylor, Address: 7949 Clark Common Suite 429
Hillborough, KS 57118, Email: stephen90@example.org, Phone: +1-215-712-6668x157
```

## Mostly AI

"Mostly AI" is a company that specializes in synthetic data generation. Synthetic data is artificially generated data that mimics the characteristics of real-world data but does not contain any personally identifiable information (PII). The use of synthetic data can be valuable for tasks such as training machine learning models, testing applications, and preserving privacy.

Mostly AI provides a platform for generating high-quality synthetic data with a focus on maintaining the statistical properties and privacy of the original data. The company's solutions are designed to help organizations overcome challenges related to data privacy regulations and security concerns while still enabling effective data-driven decision-making.

Their technology often involves advanced techniques to ensure that the synthetic data generated is representative of the original dataset without disclosing sensitive information. This can be particularly useful in industries where data privacy and compliance with regulations (such as GDPR) are critical considerations.

# <u>REFERENCES</u>

**I STUDIED VARIOUS ARTICLES TO COME UP WITH THIS REPORT AND VISITED WEBSITES ONLINE.**

- **Wikipedia.**

- **aimultiple.**

- **Youtube.**

- **questionpro.com**

- **chatgpt**