

# Reading: Generative AI for Data Anonymization

## Introduction

In today's data-driven world, organizations collect and store vast amounts of information, often containing sensitive personal data. Balancing the need to use this data for valuable insights with the ethical and legal obligation to protect individual privacy is a critical challenge. This is where data anonymization comes in.

## What is data anonymization?

Data anonymization is the process of modifying data to remove or obscure any information that could be used to directly or indirectly identify individuals. This can involve techniques like:

- Redaction: Replacing sensitive information with symbols like asterisks or hashes
- Generalization: Replacing specific values with broader categories, such as age groups instead of exact ages
- Pseudonymization: Replacing specific information with artificial identifiers

## Importance of data anonymization

Data anonymization is crucial for several reasons:

- Compliance with regulations: Many data privacy regulations, like GDPR and HIPAA, mandate the anonymization of personal data before sharing or using it for specific purposes.
- Protecting individual privacy: Anonymization safeguards individuals from the potential risks associated with data breaches or unauthorized access, such as identity theft or discrimination.
- Enabling data sharing and collaboration: Anonymized data can be shared more freely between organizations for research, analytics, and innovation while upholding individual privacy.

## Popular data anonymization strategies

Traditional data anonymization techniques often involve rule-based methods or suppression, which can be time-consuming, laborious, and potentially lead to data utility loss. Generative AI, however, offers a powerful alternative:

- **Generative models:** These AI models can learn the underlying patterns and relationships within a data set and use them to generate synthetic data that closely resembles the original data but without containing any personally identifiable information. The synthetic data allows organizations to achieve a high level of anonymization while preserving the statistical properties of the data for meaningful analysis.
- **Differential privacy:** This approach adds controlled noise to the data, making it difficult to identify individuals while still enabling accurate statistical inferences from the aggregated data.

## Generative AI for enhanced anonymization

As the amount of sensitive data collected by organizations continues to grow, balancing its use with the ethical and legal requirements of data privacy becomes increasingly crucial. Generative AI offers a promising solution through data anonymization. This technology uses AI models to learn underlying patterns within data sets. These models then generate synthetic data that closely resembles the original data while simultaneously removing any information that could potentially identify individuals. This allows valuable insights to be extracted from the data while upholding individual privacy. Generative AI, therefore, holds significant potential for fostering responsible data-driven innovation in a privacy-conscious world.

## Practical scenario: Anonymizing hospital data

Imagine a hospital wants to share anonymized patient data with researchers to study the effectiveness of a new treatment. The hospital's data includes patient names, ages, diagnoses, and treatment outcomes. Generative AI can be used to:

1. Train a generative model on the hospital data, excluding any direct identifiers like names.
2. The model can then generate synthetic patient records that preserve the statistical relationships between different data points (for example, age and diagnosis) but replace names and other identifying details with artificial values.
3. Researchers can then use the anonymized data to analyze treatment effectiveness without compromising patient privacy.

## GenAI to create Python code for anonymization

Consider a case where the data has attributes of name, age and address of different people. You need to run an anonymization process on this data in the following manner.

1. Replace the vowels in the name with the character '\*'.
2. Replace the exact age with a range, for example, 23 becomes 20s.
3. Replace the address with the last two characters, which represent the area code.

You can use GPT based model to create the code that performs this anonymization process. If using [ChatGPT](#), you can use the following prompt.

Create a Python code that performs the following anonymization tasks

1. Replace the vowels in the name with the character '\*'.
  2. Replace the exact age with a range, for example, 23 becomes 20s.
  3. Replace the address with the last two characters, which represent the area code.
- Test this code for at least 3 fictional entries

You can expect a response containing a code that will do the task. The code will be similar to what is shown below.

```
import re
def anonymize_data(name, age, address):
    # Task 1: Replace vowels in the name with '*'
    anonymized_name = re.sub('[aeiouAEIOU]', '*', name)

    # Task 2: Replace exact age with a range
    age_range = str(age // 10) + '0s'

    # Task 3: Replace address with last 2 characters (area code)
    area_code = address[-2:]

    return anonymized_name, age_range, area_code
# Test the anonymization function
entries = [
    ("John Doe", 23, "123 Main St, Seattle, WA"),
    ("Alice Smith", 35, "456 Elm St, Boston, MA"),
    ("Bob Johnson", 41, "789 Oak St, Manhattan, NY")
]
for entry in entries:
    anonymized_entry = anonymize_data(*entry)
    print("Anonymized Entry:", anonymized_entry)
```

If you execute this code on a Jupyter notebook interface, you can see the following output:

```
 import re

def anonymize_data(name, age, address):
    # Task 1: Replace vowels in the name with '*'
    anonymized_name = re.sub('[aeiouAEIOU]', '*', name)


    # Task 2: Replace exact age with a range
    age_range = str(age // 10) + '0s'

    # Task 3: Replace address with last 2 characters (area code)
    area_code = address[-2:]

    return anonymized_name, age_range, area_code

# Test the anonymization function
entries = [
    ("John Doe", 23, "123 Main St, Seattle, WA"),
    ("Alice Smith", 35, "456 Elm St, Boston, MA"),
    ("Bob Johnson", 41, "789 Oak St, Manhattan, NY")
]

for entry in entries:
    anonymized_entry = anonymize_data(*entry)
    print("Anonymized Entry:", anonymized_entry)
```

```
 Anonymized Entry: ('J*hn D**', '20s', 'WA')
Anonymized Entry: (*l*c* Sm*th', '30s', 'MA')
Anonymized Entry: ('B*b J*hns*n', '40s', 'NY')
```

This clearly shows the process has been successful in protecting the identity of the individuals.

## Author(s)

[Abhishek Gagneja](#)

© IBM Corporation. All rights reserved.

