

This document describes the improved data processing flow for generating a privacy-compliant synthetic dataset using Microsoft Fabric, PySpark, and SDV (CTGAN). This optimized approach addresses high-cardinality issues and improves CTGAN performance while ensuring privacy compliance (GDPR, PCI-DSS).

Step 1: Import necessary libraries and Load Raw Data

Source: Files/transactions.csv stored in the Microsoft Fabric Lakehouse. (1 million rows)

Schema Example: CustomerName | Email | CardNumber | Age | Amount | Country | TransactionDate
Loads data into a Spark DataFrame with automatic schema inference.

```
1 1 from pyspark.sql.functions import col, sha2, concat_ws
2 2 import pandas as pd
3 3 from sdv.single_table import CTGANSynthesizer
4 4 from sdv.metadata import SingleTableMetadata
5 5
6 6 # -----
7 7 # LOAD RAW DATA
8 8 # -----
9 9 raw_df = spark.read.csv("Files/transactions.csv", header=True, inferSchema=True)
```

Step 2: Mask Direct Identifiers (Hashing)

Purpose: Remove PII by hashing CustomerName and Email with SHA-256 and dropping sensitive columns.

```
14 14 # -----
15 15 # MASK DIRECT IDENTIFIERS (HASHING)
16 16 # -----
17 17 masked_df = (
18 18     raw_df
19 19     .withColumn("CustomerHash", sha2(concat_ws(":", col("CustomerName"), col("Email")), 256))
20 20     .drop("CustomerName", "Email", "CardNumber")  # Remove PII
21 21 )
```

Step 3: Sample Data to Pandas

Limit to 200,000 rows to improve CTGAN performance. Required since CTGAN works on pandas DataFrames.

```
23 23 # -----
24 24 # SAMPLE DATA TO PANDAS
25 25 # -----
26 26 sample_data = masked_df.limit(200_000).toPandas()
```

Step 4: Drop Unique Identifiers

CustomerHash is dropped to prevent extremely high cardinality which would cause CTGAN to create hundreds of thousands of one-hot encoded columns.

```
28 # -----
29 # DROP UNIQUE IDENTIFIERS
30 #
31 if "CustomerHash" in sample_data.columns:
32     sample_data = sample_data.drop(columns=["CustomerHash"])
```

Step 5: Handle High-Cardinality Columns

- Group rare countries (<50 occurrences) into 'Other'.
- Convert TransactionDate to numeric features (DayOfWeek, Month) instead of treating every date as a unique category.

```
34 # -----
35 # HANDLE HIGH-CARDINALITY COLUMNS
36 #
37 # Reduce rare countries to 'Other'
38 if 'Country' in sample_data.columns:
39     value_counts = sample_data['Country'].value_counts()
40     rare_countries = value_counts[value_counts < 50].index
41     sample_data['Country'] = sample_data['Country'].replace(rare_countries, 'Other')
42
43 # Extract date features
44 if 'TransactionDate' in sample_data.columns:
45     sample_data['TransactionDate'] = pd.to_datetime(sample_data['TransactionDate'], errors='coerce')
46     sample_data['DayOfWeek'] = sample_data['TransactionDate'].dt.dayofweek
47     sample_data['Month'] = sample_data['TransactionDate'].dt.month
48     sample_data = sample_data.drop(columns=['TransactionDate'])
```

Step 6: Detect and Update Metadata

Auto-detects data types, then explicitly sets:

- Age, Amount: numerical
- Country: categorical

```
50 # -----
51 # DETECT AND UPDATE METADATA
52 #
53 metadata = SingleTableMetadata()
54 metadata.detect_from_dataframe(data=sample_data)
55
56 # Ensure correct data types
57 if 'Age' in sample_data.columns:
58     metadata.update_column('Age', sdtype='numerical')
59 if 'Amount' in sample_data.columns:
60     metadata.update_column('Amount', sdtype='numerical')
61 if 'Country' in sample_data.columns:
62     metadata.update_column('Country', sdtype='categorical')
63
64 print("✓ Metadata processed successfully.")
```

Step 7: Train CTGAN Synthesizer

Trains the CTGAN model for 80 epochs to learn statistical patterns of the dataset.

```
66 # -----
67 # TRAIN CTGAN SYNTHESIZER
68 #
69 synthesizer = CTGANSynthesizer(metadata, epochs=80, verbose=True)
70 synthesizer.fit(sample_data)
```

Step 8: Generate Synthetic Data

Creates 1,000,000 synthetic records and converts them back to Spark DataFrame for large-scale processing.

```
72 # -----
73 # GENERATE SYNTHETIC DATA
74 #
75 synthetic_data = synthesizer.sample(1_000_000)
76 synthetic_df = spark.createDataFrame(synthetic_data)
77
78 print(f"✓ Generated synthetic dataset with {synthetic_df.count():,} rows.")
```

Step 9: Apply Differential Privacy

Adds noise to Age and Amount columns to strengthen privacy guarantees using epsilon=0.5.

```
80 # -----
81 # APPLY DIFFERENTIAL PRIVACY (NOISE INJECTION)
82 #
83 epsilon = 0.5
84 dp_df = (
85     synthetic_df
86     .withColumn("Age", (col("Age") + (rand() * 2 - 1) / epsilon).cast("int"))
87     .withColumn("Amount", spark_round(col("Amount") + (rand() * 10 - 5) / epsilon, 2))
88 )
89
90 print("✓ Differential privacy noise applied.")
```

Step 10: Save Dataset to Fabric Lakehouse

Final synthetic dataset is stored as a Delta table for analytics and ML pipelines.

```
92 # -----
93 # SAVE PRIVACY-COMPLIANT SYNTHETIC DATASET TO FABRIC LAKEHOUSE
94 #
95 dp_df.write.format("delta").mode("overwrite").saveAsTable('synthetic_transactions_delta')
```

Flow Validation

- PII is securely hashed and dropped before model training.
- Unique identifiers are removed to prevent dimensionality explosion.
- High-cardinality and date columns are optimized for CTGAN.
- Metadata is explicitly updated for better model quality.
- GAN-based synthetic data generation ensures realistic but non-identifiable data.
- Differential privacy noise adds extra protection to numeric fields.
- Data is stored in Delta Lakehouse, ready for analytics and ML pipelines.

Output



Recommended Enhancements

1. Automate rare category threshold selection based on dataset size.
2. Store metadata JSON for reproducibility.
3. Parameterize epsilon for configurable privacy levels.
4. Add quality evaluation metrics for synthetic data fidelity.
5. Integrate the entire pipeline as an automated Fabric notebook task.
6. The above output is **not ideal** because the synthetic data should be closer to the real distribution. However, it is common CTGAN to produce this pattern without tuning. Further research needs to be done, for example, increasing the **training epochs** (e.g., 150–200) for better convergence. Or **Data Normalization/Scaling** for Amount and Age before training (GANs work better with normalized inputs).