

```

1  # =====
2  # Script: Synthetic Data DSWB Training.py
3  # Purpose: Generate synthetic data from a PostgreSQL source using CTGAN
4  # Author: Tathagata Bhattacharjee
5  # Dependencies: CTGAN (via GitHub), scikit-learn, pandas, SQLAlchemy, matplotlib, seaborn
6  # Install ctgan from Git: pip install git+https://github.com/sdv-dev/CTGAN.git
7  # Note: Ensure all required Python packages are installed before running.
8  # =====
9
10 # ----- #
11 # Import required dependencies #
12 # ----- #
13 import pandas as pd # For data handling and manipulation
14 import numpy as np # For numerical operations and sampling
15 from sqlalchemy import create_engine, text # To connect and query a PostgreSQL database
16 from sklearn.preprocessing import LabelEncoder # To encode categorical variables for
    modeling
17 from ctgan import CTGAN # Conditional Tabular GAN for generating synthetic tabular data
18 import matplotlib.pyplot as plt # (Optional) For visualization
19 import seaborn as sns # (Optional) For enhanced plots
20 from datetime import timedelta # For date calculations
21
22 # ----- #
23 # STEP 1: Load Source Data from PostgreSQL #
24 # ----- #
25 print("\n#####")
26 print("\nSTEP 2: STEP 1: Load Source Data\n")
27
28 # Create connection to PostgreSQL database
29 engine = create_engine('postgresql://postgres:password1234@localhost:5432/dswb_training')
30
31 # SQL query to extract relevant fields from the table
32 query = """
33     SELECT id, first_name, last_name, sex, dob, village_name, occupation,
34            covid_19_first_vacc_date, age_on_first_vacc, vacc_manufacturer
35     FROM synthetic_dswb_training.source_dataset_fake
36     ORDER BY id;
37 """
38
39 # Load query result into a pandas DataFrame
40 source_data = pd.read_sql(query, engine)
41
42 # Print number of rows and columns for verification
43 print(f'Source Data Shape: {source_data.shape[0]} rows, {source_data.shape[1]} columns')
44
45 # ----- #
46 # STEP 2: Preprocess the dataset #
47 # ----- #
48 print("\n#####")
49 print("\nSTEP 2: Preprocess Data\n")
50
51 # Convert 'dob' and 'covid_19_first_vacc_date' to datetime format; drop rows with
    invalid dates
52 source_data['dob'] = pd.to_datetime(source_data['dob'], errors='coerce')
53 source_data.dropna(subset=['dob'], inplace=True)
54
55 source_data['covid_19_first_vacc_date'] =
    pd.to_datetime(source_data['covid_19_first_vacc_date'], errors='coerce')
56 source_data.dropna(subset=['covid_19_first_vacc_date'], inplace=True)
57
58 print("\n Date conversion for DoB and Vaccination Date and NaN removal done")
59
60 # Extract date components into separate columns (useful for CTGAN modeling)
61 source_data['year_of_birth'] = source_data['dob'].dt.year
62 source_data['month_of_birth'] = source_data['dob'].dt.month
63 source_data['day_of_birth'] = source_data['dob'].dt.day
64 source_data['year_of_vacc'] = source_data['covid_19_first_vacc_date'].dt.year
65 source_data['month_of_vacc'] = source_data['covid_19_first_vacc_date'].dt.month
66 source_data['day_of_vacc'] = source_data['covid_19_first_vacc_date'].dt.day

```

```

66
67 # Drop source datetime columns after extraction
68 source_data.drop(columns=['dob', 'covid_19_first_vacc_date'], inplace=True)
69
70 # List of categorical columns to encode
71 categorical_cols = ['id', 'first_name', 'last_name', 'sex', 'village_name',
72 'occupation', 'vacc_manufacturer']
73 source_value_maps = {} # Dictionary to store LabelEncoders for decoding later
74 print("\n Handling Categorical columns done")
75
76 # Apply Label Encoding on each categorical column
77 for col in categorical_cols:
78     le = LabelEncoder()
79     source_data[col] = le.fit_transform(source_data[col].astype(str)) # Encode as string
80     source_value_maps[col] = le # Save encoder for later decoding
81 print("\n Label Encoding done")
82
83 # Define numerical columns for CTGAN (includes date parts and age)
84 numerical_cols = ['year_of_birth', 'month_of_birth', 'day_of_birth', 'year_of_vacc',
85 'month_of_vacc', 'day_of_vacc', 'age_on_first_vacc']
86
87 # Store male/female first names from source data for name assignment later
88 male_names_source = source_data[source_data['sex'] == 1]['first_name'].unique()
89 female_names_source = source_data[source_data['sex'] == 0]['first_name'].unique()
90
91 # Store source sex distribution (used to validate synthetic output)
92 sex_distribution_source = source_data['sex'].value_counts(normalize=True)
93
94 # ----- #
95 # STEP 3: Train the CTGAN model #
96 # ----- #
97 print("\n#####")
98 print("\nSTEP 3: Fit the CTGAN Model\n")
99
100 # Initialize CTGAN with specified training epochs
101 ctgan = CTGAN(epochs=1000, verbose=True)
102
103 # Train the CTGAN on preprocessed source data
104 ctgan.fit(source_data, discrete_columns=categorical_cols)
105 print("\n Step 3 to fit the CTGAN model completed")
106
107 # ----- #
108 # STEP 4: Generate Synthetic Data from CTGAN #
109 # ----- #
110 print("\n#####")
111 print("\nSTEP 4: Generate Synthetic Data\n")
112
113 # Generate synthetic data of the same size as source
114 synthetic_data = ctgan.sample(len(source_data))
115
116 # Display diagnostic information
117 print(f'Synthetic Data Shape: {synthetic_data.shape[0]} rows, {synthetic_data.shape[1]}
118 columns')
119 print("\nColumns in synthetic_data after sampling:")
120 print(synthetic_data.columns)
121 print("\nFirst few rows of synthetic_data after sampling with encoding:")
122 print(synthetic_data.head())
123 print("\nData types of synthetic_data after sampling:")
124 print(synthetic_data.dtypes)
125
126 # ----- #
127 # STEP 5: Reconstruct dates and filter for adults (age >= 18 at vaccination)
128 # ----- #
129 print("\n#####")
130 print("\nSTEP 5: Reconstructing dates and ensuring adult vaccination & sex-name
131 consistency\n")
132
133 # Round date components to ensure valid integers, then convert to Int64 (nullable)
134 for col in ['year_of_birth', 'month_of_birth', 'day_of_birth', 'year_of_vacc',

```

```

131     'month_of_vacc', 'day_of_vacc']):
132         synthetic_data[col] = synthetic_data[col].round().astype('Int64')
133
134 # Rebuild full date columns from components
135 synthetic_data['date_of_birth'] = pd.to_datetime({
136     'year': synthetic_data['year_of_birth'],
137     'month': synthetic_data['month_of_birth'],
138     'day': synthetic_data['day_of_birth']
139 }, errors='coerce')
140
141 synthetic_data['date_of_vacc'] = pd.to_datetime({
142     'year': synthetic_data['year_of_vacc'],
143     'month': synthetic_data['month_of_vacc'],
144     'day': synthetic_data['day_of_vacc']
145 }, errors='coerce')
146
147 # Compute age at vaccination in years
148 synthetic_data['age_at_vaccination'] = (synthetic_data['date_of_vacc'] -
149 synthetic_data['date_of_birth']).dt.days / 365.25
150
151 # Iteratively re-sample non-adult records (<18) to ensure adult population
152 adult_threshold = 18
153 while True:
154     non_adults = synthetic_data[synthetic_data['age_at_vaccination'] < adult_threshold]
155     if non_adults.empty:
156         break
157     print(f"\nGenerated {len(non_adults)} non-adult records. Re-sampling...")
158
159     resampled_non_adults = ctgan.sample(len(non_adults))
160     for col in ['year_of_birth', 'month_of_birth', 'day_of_birth', 'year_of_vacc',
161 'month_of_vacc', 'day_of_vacc']:
162         resampled_non_adults[col] = resampled_non_adults[col].round().astype('Int64')
163
164     resampled_non_adults['date_of_birth'] = pd.to_datetime({
165         'year': resampled_non_adults['year_of_birth'],
166         'month': resampled_non_adults['month_of_birth'],
167         'day': resampled_non_adults['day_of_birth']
168     }, errors='coerce')
169     resampled_non_adults['date_of_vacc'] = pd.to_datetime({
170         'year': resampled_non_adults['year_of_vacc'],
171         'month': resampled_non_adults['month_of_vacc'],
172         'day': resampled_non_adults['day_of_vacc']
173     }, errors='coerce')
174
175     resampled_non_adults['age_at_vaccination'] = (resampled_non_adults['date_of_vacc'] -
176 resampled_non_adults['date_of_birth']).dt.days / 365.25
177
178     # Replace non-adult rows with resampled valid rows
179     synthetic_data.loc[non_adults.index, resampled_non_adults.columns] =
180 resampled_non_adults
181
182 # Finalize age column and remove temporary ones
183 synthetic_data['age_on_first_vacc'] =
184 synthetic_data['age_at_vaccination'].fillna(-1).round().astype(int)
185 synthetic_data.loc[synthetic_data['age_on_first_vacc'] == -1, 'age_on_first_vacc'] =
186 np.nan
187 synthetic_data.dropna(subset=['age_on_first_vacc'], inplace=True)
188 synthetic_data['age_on_first_vacc'] = synthetic_data['age_on_first_vacc'].astype(int)
189 synthetic_data.drop(columns=['age_at_vaccination'], inplace=True)
190 synthetic_data.drop(columns=['year_of_birth', 'month_of_birth', 'day_of_birth',
191 'year_of_vacc', 'month_of_vacc', 'day_of_vacc'], inplace=True)
192 synthetic_data.dropna(subset=['date_of_birth', 'date_of_vacc'], inplace=True)
193
194 # ----- #
195 # STEP 6: Ensure first names are consistent with sex assignment #
196 # ----- #
197 print("\n#####")
198 print("\n STEP 6: Ensure sex matches gender-specific names and maintain statistical
199 properties\n")

```

```

191
192 # Map numeric sex values to readable labels temporarily
193 synthetic_data['sex_source'] = synthetic_data['sex'].map({0: 'Female', 1: 'Male'})
194
195 # Replace first names with sex-consistent samples from source data
196 def assign_gender_consistent_name(row):
197     if row['sex_source'] == 'Male':
198         return np.random.choice(male_names_source)
199     elif row['sex_source'] == 'Female':
200         return np.random.choice(female_names_source)
201     else:
202         return np.nan
203
204 synthetic_data['first_name_consistent'] =
205 synthetic_data.apply(assign_gender_consistent_name, axis=1)
206 synthetic_data['first_name'] = synthetic_data['first_name_consistent']
207 synthetic_data.drop(columns=['first_name_consistent', 'sex_source'], inplace=True)
208
209 # Show sex distribution comparison for validation
210 synthetic_sex_distribution = synthetic_data['sex'].value_counts(normalize=True)
211 print("\nSource Sex Distribution:")
212 print(synthetic_sex_distribution)
213 print("\nInitial Synthetic Sex Distribution:")
214 print(synthetic_sex_distribution)
215
216 # ----- #
217 # STEP 7: Handle missing values in selected string fields #
218 # ----- #
219 print("\n#####")
220 print("\n STEP 7: Fill missing or blank last names, villages and occupation with random
221 values from the source dataset\n")
222
223 # Fill blank/missing values with random samples from source
224 for col in ['last_name', 'village_name', 'occupation', 'vacc_manufacturer']:
225     random_values = source_data[col].dropna().unique()
226     blank_or_missing = (synthetic_data[col].isna()) | (synthetic_data[col] == '')
227     blank_count = blank_or_missing.sum()
228     if blank_count > 0:
229         synthetic_data.loc[blank_or_missing, col] = np.random.choice(random_values,
230 size=blank_count)
231
232 # ----- #
233 # STEP 8: Decode encoded values (categoricals) to source strings #
234 # ----- #
235 print("\n#####")
236 print("\n STEP 8: Re-encode numeric columns back to source string values\n")
237
238 # Use saved LabelEncoders to reverse-encode fields back to human-readable values
239 for col in ['first_name', 'last_name', 'village_name', 'occupation', 'vacc_manufacturer']:
240     le = source_value_maps[col]
241     value_mapping = {i: le.inverse_transform([i])[0] for i in range(len(le.classes_))}
242     synthetic_data[col] = synthetic_data[col].map(value_mapping).fillna('')
243
244 # Final decoding for sex and ID formatting
245 synthetic_data['sex'] = synthetic_data['sex'].replace({1: 'Male', 0: 'Female'})
246 synthetic_data['id'] = np.abs(synthetic_data['id']) # Ensure positive values
247 synthetic_data['id'] = synthetic_data['id'].astype(int).astype(str) # Convert to string
248
249 print("\nSynthetic data sample after re-encoding")
250 print(synthetic_data.head())
251
252 # ----- #
253 # STEP 9: Save final synthetic dataset to a CSV #
254 # ----- #
255 print("\n#####")
256 print("\n STEP 9: Store the Synthetic Data to a CSV files\n")
257
258 # Map internal column names to output schema
259 columns_to_save = {

```

```

257     'id': 'id',
258     'first_name': 'first_name',
259     'last_name': 'last_name',
260     'sex': 'sex',
261     'date_of_birth': 'dob',
262     'village_name': 'village_name',
263     'occupation': 'occupation',
264     'date_of_vacc': 'covid_19_first_vacc_date',
265     'age_on_first_vacc': 'age_on_first_vacc',
266     'vacc_manufacturer': 'vacc_manufacturer'
267 }
268
269 # Select and rename columns for export
270 synthetic_data_for_csv =
synthetic_data[list(columns_to_save.keys())].rename(columns=columns_to_save)
271
272 # Preview final dataset and write to file
273 print("\nSynthetic data sample for CSV:")
274 print(synthetic_data_for_csv.head())
275 synthetic_data_for_csv.to_csv('synthetic_data_dswb_training.csv', index=False)
276 print("\n Synthetic data saved successfully to CSV file.")
277
278 #####
279 # STEP 10: Store the Synthetic Data to a PostgreSQL Table
280 #####
281 print("\n#####")
282 print("\n STEP 10: Store the Synthetic Data to a PostgreSQL Table \n")
283 # Step 10a: Create the table if not exists
284 create_table_query = """
285 CREATE TABLE IF NOT EXISTS synthetic_dswb_training.synthetic_dataset_dswb_training
286 (
287     id text,
288     first_name text,
289     last_name text,
290     sex text,
291     dob timestamp without time zone,
292     village_name text,
293     occupation text,
294     covid_19_first_vacc_date timestamp without time zone,
295     age_on_first_vacc integer,
296     vacc_manufacturer text
297 );
298 """
299
300 with engine.connect() as connection:
301     connection.execute(text(create_table_query)) # Create table if it doesn't exist
302
303 # Step 10b: Prepare data for insertion with correct column names
304 synthetic_data_for_db = synthetic_data.rename(columns={
305     'date_of_birth': 'dob',
306     'date_of_vacc': 'covid_19_first_vacc_date'
307 })
308
309 # Step 10c: Insert synthetic data into the table
310 synthetic_data_for_db[['id', 'first_name', 'last_name', 'sex', 'dob', 'village_name',
'occupation', 'covid_19_first_vacc_date', 'age_on_first_vacc',
'vacc_manufacturer']].to_sql(
311     'synthetic_dataset_dswb_training',
312     con=engine,
313     schema='synthetic_dswb_training',
314     if_exists='append', # Or 'replace' if you want to overwrite the table
315     index=False
316 )
317 print("\n Synthetic data saved successfully to PostgreSQL table.")
318
319 print("##### End of Code #####")
320
321
322

```