

```

1  # =====
2  # Script: Comparison Source and Synthetic DSWB Training.py
3  # Purpose: Compare between the source and the synthetic dataset
4  # Author: Tathagata Bhattacharjee
5  # Dependencies: pandas, SQLAlchemy, matplotlib, seaborn, dox and io
6  # # Note: Ensure all required Python packages are installed before running.
7  # =====
8
9  # Import required libraries
10 import pandas as pd # For data manipulation and analysis
11 from sqlalchemy import create_engine, text # For connecting to SQL databases
12 import matplotlib.pyplot as plt # For creating plots
13 import seaborn as sns # For enhanced statistical plots
14 from docx import Document # For creating Word (.docx) documents
15 from docx.shared import Inches # For resizing images in the Word document
16 from io import BytesIO # For in-memory byte buffer to handle images
17
18 #####
19 # STEP 1: Establish Database Connection
20 #####
21 print("\n#####")
22 print("\nSTEP 1: Establish Database Connection\n")
23
24 # Create a connection to a PostgreSQL database using SQLAlchemy
25 # Replace credentials with your actual username, password, host, port, and database name
26 engine = create_engine('postgresql://postgres:password1234@localhost:5432/PhD')
27
28 #####
29 # STEP 2: Load Data from Source and Synthetic Tables
30 #####
31 print("\n#####")
32 print("\nSTEP 2: Load Data from Source and Synthetic Tables\n")
33
34 # SQL queries to select relevant fields from the source and synthetic datasets
35 source_query = """
36     SELECT id, first_name, last_name, sex, dob, village_name, occupation,
37            covid_19_first_vacc_date, age_on_first_vacc, vacc_manufacturer
38     FROM synthetic_dswb_training.source_dataset_fake;
39 """
40 synthetic_query = """
41     SELECT id, first_name, last_name, sex, dob, village_name, occupation,
42            covid_19_first_vacc_date, age_on_first_vacc, vacc_manufacturer
43     FROM synthetic_dswb_training.synthetic_dataset_dswb_training;
44 """
45
46 # Try to load data into pandas DataFrames using SQL queries
47 # Print basic shape info or error if loading fails
48 try:
49     source_data = pd.read_sql(source_query, engine)
50     synthetic_data = pd.read_sql(synthetic_query, engine)
51     print("Data loaded successfully from both tables.")
52     print(f"source Data Shape: {source_data.shape}")
53     print(f"Synthetic Data Shape: {synthetic_data.shape}")
54 except Exception as e:
55     print(f"Error loading data: {e}")
56     exit() # Exit if there's an error
57
58 #####
59 # STEP 3: Data Preprocessing for Comparison (if needed)
60 #####
61 print("\n#####")
62 print("\nSTEP 3: Data Preprocessing for Comparison (if needed)\n")
63
64 # Convert 'dob' and 'covid_19_first_vacc_date' to datetime format for both datasets
65 source_data['dob'] = pd.to_datetime(source_data['dob'], errors='coerce')
66 source_data['covid_19_first_vacc_date'] =
pd.to_datetime(source_data['covid_19_first_vacc_date'], errors='coerce')
67 synthetic_data['dob'] = pd.to_datetime(synthetic_data['dob'], errors='coerce')
68 synthetic_data['covid_19_first_vacc_date'] =

```

```

pd.to_datetime(synthetic_data['covid_19_first_vacc_date'], errors='coerce')
67
68 # Extract date components from date of birth for fine-grained comparison
69 source_data['year_of_birth'] = source_data['dob'].dt.year
70 source_data['month_of_birth'] = source_data['dob'].dt.month
71 source_data['day_of_birth'] = source_data['dob'].dt.day
72 synthetic_data['year_of_birth'] = synthetic_data['dob'].dt.year
73 synthetic_data['month_of_birth'] = synthetic_data['dob'].dt.month
74 synthetic_data['day_of_birth'] = synthetic_data['dob'].dt.day
75
76 # Extract date components from vaccination date for detailed analysis
77 source_data['vacc_year'] = source_data['covid_19_first_vacc_date'].dt.year
78 source_data['vacc_month'] = source_data['covid_19_first_vacc_date'].dt.month
79 synthetic_data['vacc_year'] = synthetic_data['covid_19_first_vacc_date'].dt.year
80 synthetic_data['vacc_month'] = synthetic_data['covid_19_first_vacc_date'].dt.month
81
82 #####
83 # STEP 4: Graphical and Tabular Comparisons and Saving to Word File
84 #####
85 print("\n#####")
86 print("\nSTEP 4: Graphical and Tabular Comparisons and Saving to Word File\n")
87
88 # Set a nice seaborn style for all plots
89 sns.set_style("whitegrid")
90
91 # Initialize Word document for saving plots and tables
92 document = Document()
93 document.add_heading('Comparison of Source and Synthetic Data', level=0)
94
95 # Define reusable function to compare a column visually and tabularly
96 def plot_comparison_and_save(source, synthetic, column, title, kind='bar'):
97     # Add heading for the section in the Word document
98     document.add_heading(title, level=1)
99
100     # Create two subplots: one for source and one for synthetic
101     fig, axes = plt.subplots(1, 2, figsize=(12, 5))
102
103     # Plot histogram or bar chart for source data
104     if kind == 'hist':
105         sns.histplot(source[column].dropna(), kde=False, ax=axes[0], bins=30)
106     else:
107         source[column].value_counts().nlargest(10).plot(kind='barh', ax=axes[0])
108         axes[0].invert_yaxis() # Most frequent at top
109
110     axes[0].set_title(f'Source Data - {title}')
111     axes[0].set_xlabel(column)
112
113     # Plot histogram or bar chart for synthetic data
114     if kind == 'hist':
115         sns.histplot(synthetic[column].dropna(), kde=False, ax=axes[1], bins=30)
116     else:
117         synthetic[column].value_counts().nlargest(10).plot(kind='barh', ax=axes[1])
118         axes[1].invert_yaxis() # Most frequent at top
119
120     axes[1].set_title(f'Synthetic Data - {title}')
121     axes[1].set_xlabel(column)
122
123     plt.tight_layout()
124
125     # Save the plot image in memory and insert into Word
126     buffer = BytesIO()
127     plt.savefig(buffer, format='png')
128     document.add_picture(buffer, width=Inches(6))
129     plt.close(fig)
130
131     # Add table comparing frequency and percentage of values
132     document.add_heading(f'Distribution of {column}', level=2)
133
134     # Compute value counts and percentages

```

```

135     source_counts = source[column].value_counts()
136     source_percentages = source[column].value_counts(normalize=True) * 100
137     synthetic_counts = synthetic[column].value_counts()
138     synthetic_percentages = synthetic[column].value_counts(normalize=True) * 100
139
140     # Merge into one DataFrame for side-by-side comparison
141     comparison_df = pd.DataFrame({
142         'Source Count': source_counts,
143         'Source %': source_percentages,
144         'Synthetic Count': synthetic_counts,
145         'Synthetic %': synthetic_percentages
146     }).fillna(0).sort_index()
147
148     # Create a table in Word with headers + rows
149     table = document.add_table(rows=comparison_df.shape[0] + 1,
150                                cols=comparison_df.shape[1] + 1)
151     table.style = 'Table Grid'
152
153     # Header row
154     hdr_cells = table.rows[0].cells
155     hdr_cells[0].text = column # Column name for index
156     for i, col in enumerate(comparison_df.columns):
157         hdr_cells[i + 1].text = col
158
159     # Fill in data rows
160     for i, index in enumerate(comparison_df.index):
161         row_cells = table.rows[i + 1].cells
162         row_cells[0].text = str(index)
163         for j, val in enumerate(comparison_df.loc[index]):
164             row_cells[j + 1].text = str(round(val, 2))
165
166     document.add_paragraph("\n") # Add space between plots
167
168     # Call the function above for all relevant columns
169     plot_comparison_and_save(source_data, synthetic_data, 'sex', 'Sex Distribution')
170     plot_comparison_and_save(source_data, synthetic_data, 'village_name', 'Village Name
171     Distribution')
172     plot_comparison_and_save(source_data, synthetic_data, 'occupation', 'Occupation
173     Distribution')
174     plot_comparison_and_save(source_data, synthetic_data, 'vacc_manufacturer', 'Vaccine
175     Manufacturer Distribution')
176     plot_comparison_and_save(source_data, synthetic_data, 'year_of_birth', 'Year of Birth
177     Distribution', kind='hist')
178     plot_comparison_and_save(source_data, synthetic_data, 'month_of_birth', 'Month of Birth
179     Distribution', kind='hist')
180     plot_comparison_and_save(source_data, synthetic_data, 'day_of_birth', 'Day of Birth
181     Distribution', kind='hist')
182     plot_comparison_and_save(source_data, synthetic_data, 'age_on_first_vacc', 'Age on First
183     Vaccination Distribution', kind='hist')
184     plot_comparison_and_save(source_data, synthetic_data, 'vacc_year', 'Vaccination Year
185     Distribution', kind='hist')
186     plot_comparison_and_save(source_data, synthetic_data, 'vacc_month', 'Vaccination Month
187     Distribution', kind='hist')
188
189     print("\nGraphical and tabular comparisons generated and saved to Word.")
190
191     #####
192     # STEP 5: Save the Word Document
193     #####
194     print("\n#####")
195     print("\nSTEP 5: Save the Word Document\n")
196
197     # Save the document as a .docx file
198     try:
199         document.save('comparison_report.docx')
200         print("Comparison report saved to comparison_report.docx")
201     except Exception as e:
202         print(f"Error saving Word document: {e}")
203
204

```

```
194 print("##### End of Code #####")
195
```