

# **Synthetic Data Factory**

Automated, Statistically-Driven Synthetic Data Generation  
and Validation Platform

## **Project Overview:**

The Synthetic Data Factory (SDF) is an intelligent web application that enables users to generate realistic, privacy-safe, and statistically validated synthetic datasets with ease. It allows developers, data scientists, testers, and researchers to create custom datasets for analytics, testing, and machine learning without relying on sensitive real-world data.

This system bridges the gap between data accessibility and privacy by learning the structure, distribution, and relationships within uploaded datasets or defined schemas, and then reproducing that same statistical behaviour in newly generated synthetic data.

By integrating AI-based data pattern detection, Faker-driven content generation, and statistical validation techniques, SDF ensures every dataset is both realistic and unique, making it ideal for simulation, system testing, research, and training environments.

## **Scenario 1: Data Science and Machine Learning Development**

Data scientists and machine learning engineers often struggle to find large, high-quality datasets that are both diverse and privacy-compliant. The Synthetic Data Factory (SDF) addresses this challenge by generating realistic synthetic data that mirrors the statistical patterns of real-world datasets without exposing sensitive information. For instance, a financial technology startup developing a loan prediction model can use SDF to create synthetic data representing customer profiles, income ranges, credit scores, and loan histories. Instead of relying on confidential client data, they can upload a small sample dataset or define the required schema, and SDF will generate thousands of synthetic records maintaining the same relationships and distributions.

With SDF's built-in visualization and validation tools, teams can ensure that the generated data accurately reflects the original structure. This allows safe, large-scale model training and testing while remaining compliant with data privacy standards. Through its integration of Streamlit, Faker, and statistical validation frameworks, the Synthetic Data Factory not only accelerates data preparation but also empowers organizations to innovate responsibly. Ultimately, SDF transforms data scarcity into data abundance driving progress in AI, analytics, and decision intelligence securely and efficiently.

## **Scenario 2: Application Development and Database Testing**

In software development and testing, teams often require large, diverse, and realistic datasets to validate performance, database operations, and system integration. However, using production data introduces serious privacy and compliance challenges. The Synthetic Data Factory (SDF)

overcomes this by generating synthetic datasets that replicate the structure and statistical behaviour of real data without revealing any sensitive information.

For example, a banking enterprise can use SDF to generate synthetic customer profiles, account details, and transaction histories to test APIs, data pipelines, and application performance under production-like conditions. Through its intelligent combination of Streamlit, Faker, and validation frameworks, SDF enables secure, scalable, and automated data population for development environments.

Its integrated MySQL connectivity module supports direct schema detection and synthetic data insertion, simplifying test data management. By ensuring realistic simulation, privacy compliance, and data consistency, SDF helps teams accelerate release cycles, enhance reliability, and minimize time and cost in test data generation.

### Architecture Overview:

The Synthetic Data Factory is a modular, intelligent data generation platform designed to create realistic, privacy-safe, and statistically balanced synthetic datasets for AI, analytics, and application testing. Built using Streamlit as the front-end framework, SDF integrates a seamless, interactive interface with a robust backend powered by Python, Faker, NumPy, Pandas, and PyMySQL for advanced data handling and simulation.

The system follows a multi-layered architecture consisting of three core backend modules: *UniversalDataGenerator*, *DatabaseHandler*, and *DataValidator*, each serving a distinct role. The *UniversalDataGenerator* module generates synthetic datasets either from uploaded CSV files or user-defined column schemas. The *DatabaseHandler* enables direct integration with MySQL databases, allowing users to generate and insert synthetic data directly based on table structures. The *DataValidator* module ensures statistical fidelity by comparing real and synthetic datasets using distribution metrics and visual validations.

The Streamlit-based front-end ties these components together through an intuitive tabbed interface featuring CSV extension, custom data creation, and database connectivity sections. With integrated validation, progress tracking, and responsive visual design, SDF provides users with a complete synthetic data ecosystem enabling them to generate, validate, and visualize data securely and efficiently for diverse enterprise and research applications.

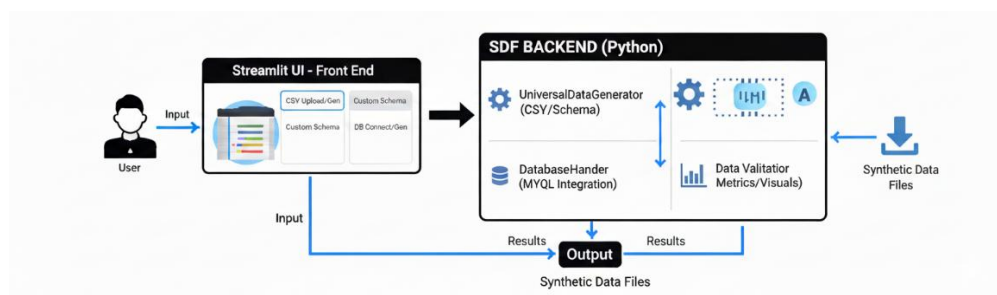


Figure 1: System architecture diagram

## Core Technologies:

- **Streamlit (Frontend Framework):** Provides an intuitive and interactive web interface for dataset generation, schema customization, and visualization. It handles real-time user inputs, displays synthetic data previews, and integrates dynamic components like tabs and progress indicators.
- **Faker (Synthetic Data Engine):** Powers the data generation process by creating realistic yet artificial values for various data types such as names, addresses, transactions, and demographic details. It supports both predefined templates and custom schema definitions.
- **NumPy & Pandas (Data Processing Core):** Enable efficient data manipulation, statistical computation, and tabular structure handling. These libraries ensure scalable and optimized performance during large-scale synthetic dataset creation and validation.
- **PyMySQL (Database Integration):** Allows direct connection with MySQL databases for real-time schema extraction and automatic insertion of generated synthetic data, enabling seamless backend testing and database population.
- **Matplotlib & Seaborn (Visualization Framework):** Used to compare and validate real vs. synthetic datasets through visual metrics such as histograms, distributions, and correlation heatmaps, ensuring statistical integrity and data balance.
- **Data Validation & Comparison Module:** Ensures that synthetic datasets maintain consistency, diversity, and similarity to real datasets by applying distribution checks, statistical summaries, and visualization-based validation.
- **Python Backend (Core Logic Layer):** Manages communication between all modules data generation, validation, and database operations that's ensuring smooth execution, modular scalability, and high-performance processing.

## Component-Wise Architecture:

Component	Description
<b>User Interface (Streamlit)</b>	Interactive dashboard for CSV upload, schema creation, and database connection with real-time data display.
<b>Universal Data Generator</b>	Creates realistic synthetic datasets using Faker, NumPy, and Pandas based on sample or custom schema.
<b>Database Handler (MySQL)</b>	Connects to MySQL, detects table schemas, and inserts generated data directly for testing and simulation.
<b>Data Validator &amp; Visualizer</b>	Compares real vs. synthetic data through statistical checks and visualizations using Matplotlib and Seaborn.
<b>Custom Schema Builder</b>	Let users define columns, types, and record limits for domain-specific synthetic data generation.

<b>Visualization &amp; Reporting</b>	Generates charts and summaries for quick validation of data quality and distribution patterns.
<b>Backend Logic Controller</b>	Manages coordination between data generation, validation, and visualization modules.
<b>Help &amp; Documentation</b>	Provides usage guidance and workflow support within the Streamlit interface.

### Pre-requisites:

- 1. Python Environment Setup:** The Synthetic Data Factory (SDF) is developed using Python 3.9+, leveraging its extensive data handling and visualization libraries. Create and activate a dedicated virtual environment (e.g., `synthetic_data_factory`) to manage dependencies and maintain a clean setup.

Official Documentation: <https://www.python.org/downloads/>

- 2. Streamlit Installation and Configuration:** Streamlit provides the front-end interface for data generation, visualization, and validation. Install Streamlit to build and run the interactive web dashboard.

Docs: <https://docs.streamlit.io/library/get-started/installation>

Tutorial: <https://www.youtube.com/watch?v=JwSS70SZdyM>

- 3. Library Installation and Core Dependencies:** SDF relies on several Python libraries for data generation, analysis, and visualization. Install all dependencies listed in `requirements.txt`.

Key Libraries:

streamlit – Interactive UI and data display

faker – Synthetic data generation

pandas – Data manipulation and structure handling

numpy – Numerical and statistical computation

matplotlib & seaborn – Data visualization and validation plots

pymysql – MySQL database connectivity and data insertion

- 4. Database Setup:** For projects requiring live database testing, install and configure MySQL Server. SDF connects via PyMySQL, enabling table schema reading and direct data insertion.

MySQL Installation: <https://dev.mysql.com/downloads/>

PyMySQL Docs: <https://pypi.org/project/PyMySQL/>

**5. Development Environment:**

Recommended IDEs for efficient development and debugging

Visual Studio Code: <https://code.visualstudio.com/>

PyCharm (Community Edition): <https://www.jetbrains.com/pycharm/download/>

**6. Optional Learning Resources:**

Enhance understanding of the tools and concepts used

Streamlit Components Gallery: <https://streamlit.io/components>

Faker Library Guide: <https://faker.readthedocs.io/>

## Project Flow:

### 1. Environment Setup and Dependency Configuration

- **Activity 1.1:** Create and activate a virtual environment and install required dependencies.
- **Activity 1.2:** Organize project structure and add UI styling with styles.css.
- **Activity 1.3:** Initialize Streamlit app and validate integration between modules.

### 2. Core Logic Development (Synthetic Data Generation Engine)

- **Activity 2.1:** Implement UniversalDataGenerator for CSV and schema-based synthetic data.
- **Activity 2.2:** Develop validation module to analyze and compare real vs synthetic data.
- **Activity 2.3:** Integrate MySQL database for schema retrieval and synthetic data insertion.

### 3. Streamlit UI Implementation and User Interaction

- **Activity 3.1:** Design multi-tab layout with navigation and professional styling.
- **Activity 3.2:** Configure input system for columns, rows, and database details with session state.
- **Activity 3.3:** Display synthetic data and visualizations with computed metrics.

### 4. Testing, Optimization, and Deployment

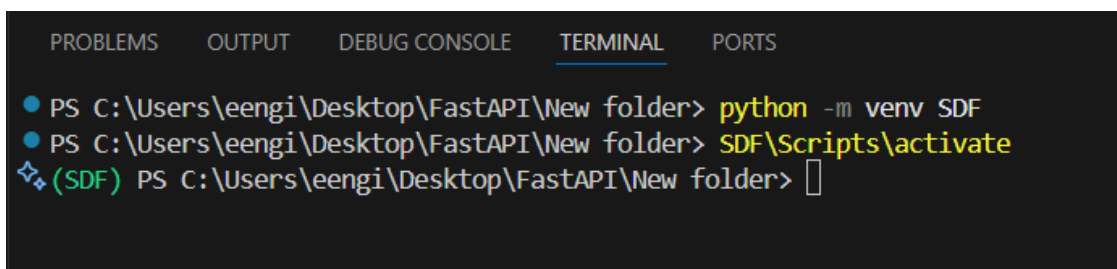
- **Activity 4.1:** Test UI components and verify accurate data generation.
- **Activity 4.2:** Validate database integration and large dataset performance.
- **Activity 4.3:** Prepare deployment and perform end-to-end validation for stability and performance.

## MILESTONE 1: Environment Setup and Dependency Configuration

This foundational milestone establishes the technical environment required for building and deploying the Synthetic Data Factory (SDF). It ensures that all dependencies, frameworks, and integrations are configured correctly for seamless execution of synthetic data generation, visualization, and validation workflows.

### Activity 1.1: Python Environment and Dependency Installation

- Create and activate a dedicated virtual environment for SDF to maintain dependency isolation.

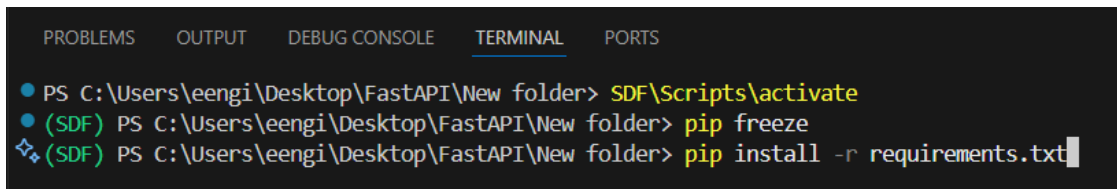


```
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL  PORTS

• PS C:\Users\eengi\Desktop\FastAPI\New folder> python -m venv SDF
• PS C:\Users\eengi\Desktop\FastAPI\New folder> SDF\Scripts\activate
❖ (SDF) PS C:\Users\eengi\Desktop\FastAPI\New folder> 
```

Figure 2: Creating & Activating Environment

- Install all required dependencies listed in requirements.txt using:



```
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL  PORTS

• PS C:\Users\eengi\Desktop\FastAPI\New folder> SDF\Scripts\activate
• (SDF) PS C:\Users\eengi\Desktop\FastAPI\New folder> pip freeze
❖ (SDF) PS C:\Users\eengi\Desktop\FastAPI\New folder> pip install -r requirements.txt
```

Figure 3: Installing requirements

- Verify compatibility across Python 3.9+ and confirm successful installation of each module.
- Test library imports and ensure data-related packages (NumPy, Pandas) and visualization tools (Matplotlib, Seaborn) work correctly.

## Activity 1.2: Project Structure Initialization

- Organize project modules (`streamlit_app.py`, `data_generator.py`, `database_handler.py`, `validation.py`) into a structured hierarchy for clarity and maintainability.

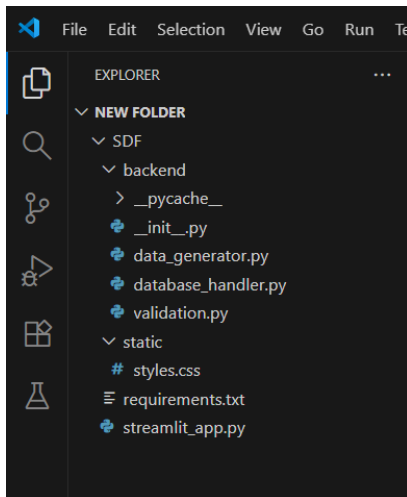


Figure 4: Folder Structure

- Add `styles.css` to define Streamlit UI design, theme consistency, and responsive layout behaviour.
- Verify file paths and imports to ensure modular independence and scalability.

## Activity 1.3: Streamlit Application Initialization

- Set up `streamlit_app.py` as the central entry point of the system using:

```
st.set_page_config(
    page_title="Synthetic Data Factory",
    layout="wide",
    initial_sidebar_state="expanded"
)

# Load custom CSS
with open(r"C:\Users\eengi\Desktop\SDG\New folder\SDF\static\styles.css") as f:
    st.markdown(f"<style>{f.read()}</style>", unsafe_allow_html=True)

# Initialize session state
if 'generator' not in st.session_state:
    st.session_state.generator = UniversalDataGenerator()
if 'db_handler' not in st.session_state:
    st.session_state.db_handler = DatabaseHandler(st.session_state.generator)
if 'validator' not in st.session_state:
    st.session_state.validator = DataValidator()
if 'synthetic_data' not in st.session_state:
    st.session_state.synthetic_data = None
```

Figure 5: Streamlit Configuration



- Ensure consistent theming, responsive layout, and component rendering across devices.
- Validate integration between all modules, confirming smooth data flow between UI and backend functions.

## MILESTONE 2: Core Logic Development (Synthetic Data Generation Engine)

This milestone builds the analytical and generation backbone of SDF. The goal is to automate realistic, statistically valid data generation using the UniversalDataGenerator and connect it with analysis, validation, and storage pipelines.

### Activity 2.1: Universal Data Generator Implementation (data\_generator.py)

- Develop the UniversalDataGenerator class to handle both CSV-based and custom schema-based synthetic data generation.

```

10 class UniversalDataGenerator:
11     def __init__(self):
12         self.fake = Faker()
13
14     def generate_from_csv(self, csv_file: str, num_rows: int = 200, validate: bool = True):
15         progress_bar = st.progress(0)
16         status_text = st.empty()
17
18         status_text.text("Uploading CSV file...")
19
20         try:
21             original_df = pd.read_csv(csv_file)
22             progress_bar.progress(25)
23
24             status_text.text("Analyzing data structure...")
25
26             with st.expander("Original Data Preview", expanded=True):
27                 st.write(f"Dataset Info: ** {len(original_df)} rows, {len(original_df.columns)} columns")
28                 st.dataframe(original_df.head(), use_container_width=True)
29
30             synthetic = {}
31             total_columns = len(original_df.columns)
32
33             for i, col in enumerate(original_df.columns):
34                 status_text.text(f"Generating column: {col} ({i+1}/{total_columns})")
35                 progress_bar.progress(25 + int(25 * i / total_columns))
36
37                 if pd.api.types.is_numeric_dtype(original_df[col]):
38                     synthetic[col] = self._generate_numeric_data(original_df[col], num_rows)
39                 elif pd.api.types.is_string_dtype(original_df[col]) or pd.api.types.is_object_dtype(original_df[col]):
40                     synthetic[col] = self._generate_text_data(col, original_df[col], num_rows)
41                 else:
42                     synthetic[col] = [f"Data_{i}" for i in range(num_rows)]
43
44             synthetic_df = pd.DataFrame(synthetic)
45
46             status_text.text("Processing relationships between columns...")
47             progress_bar.progress(75)
48
49             synthetic_df = self._handle_age_dob_relationship(synthetic_df)
50
51             progress_bar.progress(100)
52
53             status_text.text("Generation complete!")
54
55             with st.expander("Synthetic Data Preview", expanded=True):
56                 st.dataframe(synthetic_df.head(10), use_container_width=True)
57
58             return synthetic_df, original_df
59
60     except Exception as e:
61         st.error(f"Error reading CSV file: {e}")
62         return None, None

```

Figure 6: Dataset Generator from CSV

```

64     def generate_from_columns(self, columns: List[str], num_rows: int = 200):
65         progress_bar = st.progress(0)
66         status_text = st.empty()
67
68         status_text.text(f"Generating data for {len(columns)} columns...")
69
70         synthetic = {}
71         total_columns = len(columns)
72
73         for i, col in enumerate(columns):
74             status_text.text(f"Generating: {col} ({i+1}/{total_columns})")
75             progress_bar.progress(50)
76             synthetic[col] = self._generate_from_column_name(col, num_rows)
77
78         synthetic_df = pd.DataFrame(synthetic)
79
80         status_text.text("Processing relationships between columns...")
81         progress_bar.progress(75)
82
83         synthetic_df = self._handle_age_dob_relationship(synthetic_df)
84
85         progress_bar.progress(100)
86
87         status_text.text("Generation complete!")
88
89         with st.expander("Generated Data Preview", expanded=True):
90             st.dataframe(synthetic_df.head(10), use_container_width=True)
91
92         return synthetic_df

```

Figure 7: Data Generator from Column Names

- Implement configurable parameters for number of rows, data types, and randomization logic.

```

133     def _generate_numeric_data(self, column_data: pd.Series, num_rows: int):
134         clean_data = column_data.dropna()
135
136         if len(clean_data) == 0:
137             if pd.api.types.is_integer_dtype(column_data):
138                 return np.random.randint(0, 100, num_rows)
139             else:
140                 return np.round(np.random.uniform(0, 100, num_rows), 2)
141
142         min_val = clean_data.min()
143         max_val = clean_data.max()
144         mean_val = clean_data.mean()
145         std_val = clean_data.std()
146
147         if std_val > 0:
148             generated = np.random.normal(mean_val, std_val, num_rows)
149         else:
150             generated = np.random.uniform(min_val, max_val, num_rows)
151
152         generated = np.clip(generated, min_val * 0.8, max_val * 1.2)
153
154         if pd.api.types.is_integer_dtype(column_data):
155             return generated.astype(int)
156         else:
157             return np.round(generated, 2)
158
159     def _generate_text_data(self, column_name: str, column_data: pd.Series, num_rows: int):
160         clean_data = column_data.dropna()
161
162         if len(clean_data) > 0 and clean_data.nunique() <= 15:
163             return np.random.choice(clean_data.unique(), num_rows)
164
165         return self._generate_from_column_name(column_name, num_rows)

```

Figure 8: Numeric data processing

- Integrate Faker to simulate text, numeric, categorical, and temporal data types.

```

167     def _generate_from_column_name(self, column_name: str, num_rows: int):
168         col_lower = column_name.lower()
169
170         faker_methods = {}
171
172         # Personal Information
173         'name': self.fake.name,
174         'first_name': self.fake.first_name,
175         'last_name': self.fake.last_name,
176         'full_name': self.fake.name,
177         'username': self.fake.user_name,
178         'password': self.fake.password,
179
180         # Contact Information
181         'email': self.fake.email,
182         'phone': lambda: f"{random.randint(6000000000, 9999999999)}",
183         'mobile': lambda: f"{random.randint(6000000000, 9999999999)}",
184
185         # Location Information
186         'address': lambda: self.fake.address().replace('\n', ' '),
187         'street': self.fake.street_address,
188         'city': self.fake.city,
189         'state': self.fake.state,
190         'country': self.fake.country,
191         'zip': self.fake.zipcode,
192         'postal': self.fake.postcode,
193         'location': self.fake.city,
194
195         # Company & Professional
196         'company': self.fake.company,
197         'job': self.fake.job,
198         'job_title': self.fake.job,
199         'industry': self.fake.bs,
200
201         # Financial
202         'credit_card': self.fake.credit_card_number,
203         'iban': self.fake.iban,
204         'currency': self.fake.currency_code,
205
206         # Internet & Tech
207         'url': self.fake.url,
208         'website': self.fake.url,
209         'domain': self.fake.domain_name,
210         'ip': self.fake.ipv4,

```

Figure 9: Faker Mapping

```

212         # Dates & Times
213         'date': self.fake.date,
214         'time': self.fake.time,
215         'year': lambda: self.fake.year(),
216         'month': lambda: self.fake.month_name(),
217
218         # Date of Birth variations
219         'dob': self.fake.date_of_birth,
220         'date_of_birth': self.fake.date_of_birth,
221         'birth_date': self.fake.date_of_birth,
222         'birthdate': self.fake.date_of_birth,
223         'birthday': self.fake.date_of_birth,
224
225         # Products & Commerce
226         'product': self.fake.word,
227         'brand': self.fake.company,
228         'color': self.fake.color_name,
229
230         # Education
231         'school': self.fake.company,
232         'university': self.fake.company,
233         'grade': lambda: random.choice(['A', 'B', 'C', 'D', 'F']),
234
235         # Medical
236         'hospital': self.fake.company,
237         'doctor': self.fake.name,
238         'disease': lambda: random.choice(['Flu', 'Cold', 'Headache', 'Fever', 'Allergy']),
239
240         # Vehicles
241         'car': lambda: f"{self.fake.company()} {self.fake.word()}",
242         'license': self.fake.license_plate,
243
244         # Identification
245         'id': lambda: f"ID_{self.fake.random_int(1000, 9999)}",
246         'ssn': self.fake.ssn,
247         'passport': self.fake.passport_number,
248
249         # Numeric types
250         'age': lambda: random.randint(18, 70),
251         'salary': lambda: random.randint(30000, 150000),
252         'price': lambda: round(random.uniform(10, 1000), 2),
253         'quantity': lambda: random.randint(1, 100),
254         'count': lambda: random.randint(0, 100)

```

Figure 10: Faker Mapping

```

249 # Numeric types
250 'age': lambda: random.randint(18, 70),
251 'salary': lambda: random.randint(30000, 150000),
252 'price': lambda: round(random.uniform(10, 1000), 2),
253 'quantity': lambda: random.randint(1, 100),
254 'score': lambda: random.randint(0, 100),
255 'rating': lambda: random.randint(1, 5),
256 'serial': lambda: f"SN{self.fake.random_number(digits=8)}",
257
258 # Boolean types
259 'is_': lambda: random.choice([True, False]),
260 'has_': lambda: random.choice([True, False]),
261 'active': lambda: random.choice([True, False]),
262 'status': lambda: random.choice(['Active', 'Inactive', 'Pending']),
263
264
265 for pattern, faker_method in faker_methods.items():
266     if pattern in col_lower:
267         try:
268             return [faker_method() for _ in range(num_rows)]
269         except:
270             continue
271
272 if any(word in col_lower for word in ['first', 'given']):
273     return [self.fake.first_name() for _ in range(num_rows)]
274 elif any(word in col_lower for word in ['last', 'surname', 'family']):
275     return [self.fake.last_name() for _ in range(num_rows)]
276 elif any(word in col_lower for word in ['street', 'road', 'avenue']):
277     return [self.fake.street_address() for _ in range(num_rows)]
278 elif any(word in col_lower for word in ['state', 'province', 'region']):
279     return [self.fake.state() for _ in range(num_rows)]
280 elif any(word in col_lower for word in ['gender', 'sex']):
281     return [random.choice(['Male', 'Female']) for _ in range(num_rows)]
282
283 elif any(word in col_lower for word in ['number', 'count', 'total', 'amount']):
284     return [random.randint(1, 1000) for _ in range(num_rows)]
285 elif any(word in col_lower for word in ['percent', 'percentage', 'rate']):
286     return [round(random.uniform(0, 100), 2) for _ in range(num_rows)]
287
288 else:
289     return [self.fake.word() for _ in range(num_rows)]

```

Figure 11: Faker Mapping

## Activity 2.2: Data Analysis and Validation Module (validation.py)

- Develop analytical functions to calculate core statistical metrics such as mean, median, mode, variance, and correlation.

```

8 class DataValidator:
9     def __init__(self):
10         pass
11
12     def validate_synthetic_data(self, original_df: pd.DataFrame, synthetic_df: pd.DataFrame, alpha: float = 0.05):
13         st.subheader("Validation and Quality Metrics")
14
15         validation_results = {
16             'columns': [],
17             'ks_statistic': [],
18             'ks_pvalue': [],
19             'ks_significant': [],
20             'mean_original': [],
21             'mean_synthetic': [],
22             'mean_diff': [],
23             'mean_diff_percent': []
24         }
25
26         common_columns = set(original_df.columns) & set(synthetic_df.columns)
27
28         for col in common_columns:
29             if not pd.api.types.is_numeric_dtype(original_df[col]) or not pd.api.types.is_numeric_dtype(synthetic_df[col]):
30                 continue
31
32             orig_data = original_df[col].dropna()
33             synth_data = synthetic_df[col].dropna()
34
35             if len(orig_data) == 0 or len(synth_data) == 0:
36                 continue
37
38             ks_stat, ks_pvalue = ks_2samp(orig_data, synth_data)
39
40             mean_orig = orig_data.mean()
41             mean_synth = synth_data.mean()
42             mean_diff = abs(mean_orig - mean_synth)
43             mean_diff_percent = (mean_diff / abs(mean_orig)) * 100 if mean_orig != 0 else 0
44
45             validation_results['columns'].append(col)
46             validation_results['ks_statistic'].append(ks_stat)
47             validation_results['ks_pvalue'].append(ks_pvalue)
48             validation_results['ks_significant'].append(ks_pvalue < alpha)
49             validation_results['mean_original'].append(mean_orig)
50             validation_results['mean_synthetic'].append(mean_synth)
51             validation_results['mean_diff'].append(mean_diff)
52             validation_results['mean_diff_percent'].append(mean_diff_percent)
53
54         validation_df = pd.DataFrame(validation_results)

```

Figure 12: Validation of dataset

- Implement comparison logic for real vs synthetic data distributions to validate fidelity.

```

56     if len(validation_df) > 0:
57         col1, col2, col3 = st.columns(3)
58
59         significant_cols = validation_df['ks_significant'].sum()
60         total_cols = len(validation_df)
61         quality_score = self.calculate_quality_score(validation_df)
62
63         with col1:
64             st.metric("Quality Score", f"{quality_score:.1f}/100")
65
66         with col2:
67             st.metric("Distribution Match", f"({total_cols - significant_cols})/{total_cols}")
68
69         with col3:
70             st.metric("Avg Mean Difference", f"{validation_df['mean_diff_percent'].mean():.2f}%")
71
72         if quality_score >= 80:
73             st.success("EXCELLENT: Synthetic data closely matches original distribution")
74         elif quality_score >= 60:
75             st.warning("GOOD: Synthetic data reasonably matches original distribution")
76         elif quality_score >= 40:
77             st.info("FAIR: Some differences detected in synthetic data")
78         else:
79             st.error("POOR: Significant differences in synthetic data")
80
81         st.subheader("Detailed Column Analysis")
82         display_df = validation_df.copy()
83         display_df['ks_significant'] = display_df['ks_significant'].map({True: '✗', False: '✓'})
84         display_df = display_df.round(4)
85         st.dataframe(display_df, use_container_width=True)
86
87         self.plot_distributions(original_df, synthetic_df, common_columns)
88
89     else:
90         st.warning("No numeric columns available for validation.")
91
92     return validation_df

```

Figure 13: Comparing Logic

- Structure results using **Pandas DataFrames** for compatibility with the visualization layer.

```

94     def _plot_distributions(self, original_df: pd.DataFrame, synthetic_df: pd.DataFrame, common_columns: set):
95         st.subheader("Distribution Comparison")
96         for col in common_columns:
97             if pd.api.types.is_numeric_dtype(original_df[col]):
98                 fig = go.Figure()
99                 fig.add_trace(go.Histogram(
100                     x=original_df[col].dropna(),
101                     name='Original',
102                     opacity=0.7,
103                     nbinsx=20
104                 ))
105                 fig.add_trace(go.Histogram(
106                     x=synthetic_df[col].dropna(),
107                     name='Synthetic',
108                     opacity=0.7,
109                     nbinsx=20
110                 ))
111                 fig.update_layout(
112                     title=f'Distribution of {col}',
113                     xaxis_title=col,
114                     yaxis_title='Frequency',
115                     barmode='overlay'
116                 )
117                 st.plotly_chart(fig, use_container_width=True)

```

Figure 14: Plotting Function

- Integrate automated summary reports highlighting statistical differences and similarities.

```

119     def calculate_quality_score(self, validation_df: pd.DataFrame) -> float:
120         if len(validation_df) == 0:
121             return 0
122
123         ks_penalty = validation_df['ks_significant'].mean() * 40
124         mean_diff_penalty = min(30, validation_df['mean_diff_percent'].mean() / 2)
125         base_score = 100
126         quality_score = base_score - ks_penalty - mean_diff_penalty
127
128         return max(0, quality_score)

```

Figure 15: Score Calculation

### Activity 2.3: Database Integration (database\_handler.py)

- Implement secure and optimized **MySQL** database integration using **PyMySQL**.
- Allow users to:
  - Retrieve existing table schemas automatically.
  - Generate data conforming to those schemas.

```

9 class DatabaseHandler:
10     def __init__(self, data_generator):
11         self.generator = data_generator
12
13     def connect_to_mysql(self, host: str, user: str, password: str, database: str):
14         try:
15             connection = pymysql.connect(
16                 host=host,
17                 user=user,
18                 password=password,
19                 database=database,
20                 charset='utf8mb4',
21                 cursorclass=pymysql.cursors.DictCursor
22             )
23             st.success("Successfully connected to MySQL database!")
24             return connection
25         except MySQLError as err:
26             st.error(f"Failed to connect to database: {err}")
27             return None
28
29     def get_mysql_tables(self, connection) -> List[str]:
30         try:
31             with connection.cursor() as cursor:
32                 cursor.execute("SHOW TABLES")
33                 tables = [list(table.values())[0] for table in cursor.fetchall()]
34             return tables
35         except MySQLError as err:
36             st.error(f"Error fetching tables: {err}")
37             return []
38
39     def get_table_schema(self, connection, table_name: str) -> Dict[str, Any]:
40         try:
41             with connection.cursor() as cursor:
42                 cursor.execute(f"DESCRIBE {table_name}")
43                 schema = cursor.fetchall()
44
45                 schema_info = {}
46                 for column in schema:
47                     schema_info[column['Field']] = {
48                         'field': column['Field'],
49                         'type': column['Type'],
50                         'null': column['Null'],
51                         'key': column['Key'],
52                         'default': column['Default'],
53                         'extra': column['Extra']
54                     }
55                 return schema_info
56         except MySQLError as err:
57             st.error(f"Error fetching schema for {table_name}: {err}")
58             return {}
59

```

Figure 16: SQL operation function

- Insert synthetic records directly into connected databases.

```

60 def generate_from_mysql_table(self, connection, table_name: str, num_rows: int = 200):
61     progress_bar = st.progress(0)
62     status_text = st.empty()
63
64     status_text.text(f"Generating synthetic data for table: {table_name}")
65
66     schema = self.get_table_schema(connection, table_name)
67     if not schema:
68         return None, None
69
70     status_text.text(f"Analyzing table schema: {len(schema)} columns")
71     progress_bar.progress(25)
72
73     original_data = self.get_table_data(connection, table_name, limit=1000)
74     original_df = pd.DataFrame(original_data) if original_data else None
75
76     synthetic_data = {}
77     total_columns = len(schema)
78
79     for i, (column_name, column_info) in enumerate(schema.items()):
80         status_text.text(f"Generating: {column_name} ({i+1}/{total_columns})")
81         progress_bar.progress(25 + int(50 * (i / total_columns)))
82         synthetic_data[column_name] = self._generate_from_mysql_column(column_name, column_info, num_rows)
83
84     synthetic_df = pd.DataFrame(synthetic_data)
85     progress_bar.progress(100)
86
87     status_text.text("Generation complete!")
88
89     with st.expander("Synthetic Data Preview", expanded=True):
90         st.dataframe(synthetic_df.head(10), use_container_width=True)
91
92     return synthetic_df, original_df
93
94 def insert_to_mysql_table(self, connection, table_name: str, dataframe: pd.DataFrame) -> bool:
95     try:
96         with connection.cursor() as cursor:
97             columns = ', '.join(dataframe.columns)
98             placeholders = ', '.join(['%s'] * len(dataframe.columns))
99             insert_query = f"INSERT IGNORE INTO {table_name} ({columns}) VALUES ({placeholders})"
100             data_tuples = [tuple(row) for row in dataframe.values]
101             cursor.executemany(insert_query, data_tuples)
102             connection.commit()
103
104             st.success(f"Successfully inserted {cursor.rowcount} rows into {table_name}")
105             return True
106
107     except MySQLError as err:
108         st.error(f"Error inserting data into {table_name}: {err}")
109         connection.rollback()
110         return False
111

```

Figure 17: Creating and Inserting the dummy data to database

- Include robust error handling for connection failures, authentication issues, and SQL exceptions.

```

try:
    with connection.cursor() as cursor:
        columns = ', '.join(dataframe.columns)
        placeholders = ', '.join(['%s'] * len(dataframe.columns))
        insert_query = f"INSERT IGNORE INTO {table_name} ({columns}) VALUES ({placeholders})"
        data_tuples = [tuple(row) for row in dataframe.values]
        cursor.executemany(insert_query, data_tuples)
        connection.commit()

        st.success(f"Successfully inserted {cursor.rowcount} rows into {table_name}")
        return True
except MySQLError as err:
    st.error(f"Error inserting data into {table_name}: {err}")
    connection.rollback()
    return False

```

Figure 18: Exception Handling



## MILESTONE 3: Streamlit UI Implementation and User Interaction

This milestone focuses on building an **interactive, intuitive user interface** that connects backend intelligence with real-time visualization. The UI allows users to configure schemas, generate data, validate results, and interact with synthetic datasets dynamically.

### Activity 3.1: Layout Design and Navigation

- Develop a multi-tabbed interface for:

CSV Extension

Dummy Data Creation

SQL Operations

```
27 # Header
28 st.markdown('<h1 class="main-header">Synthetic Data Factory</h1>', unsafe_allow_html=True)
29 st.markdown('<h2 class="main-header">Generate intelligent, realistic, and statistically balanced synthetic data - with precision and style.</h2>', unsafe_allow_html=True)
30 st.markdown(" ")
31
32 # Main Tabs Navigation
33 tab1, tab2, tab3, tab4, tab5 = st.tabs([
34     "Home",
35     "CSV Extension",
36     "Dummy Data Creation",
37     "MySQL Operations",
38     "About"
39 ])
40
41 with tab1:
42     home_page()
43 with tab2:
44     csv_extension_page()
45 with tab3:
46     dummy_data_creation_page()
47 with tab4:
48     mysql_operations_page()
49 with tab5:
50     about_page()
51
52
```

Figure 19: Navigation buttons

- Used columns, custom Markdown, and HTML formatting for adding custom style and structure.

```

54 def home_page():
55     st.markdown("""
56     <div style='margin-bottom: 2rem;'>
57         <div style='
58             background: linear-gradient(135deg, #2082AA 0%, #8A2BE2 100%);
59             padding: 2.5rem;
60             border-radius: 15px;
61             color: white;
62             width: 100%;
63             box-shadow: 0 8px 25px rgba(0,0,0,0.15);
64             text-align: center;
65         '>
66             <h3 style='margin: 0 0 1rem 0; font-size: 1.5rem;'>Generate Smart Synthetic Data</h3>
67             <p style='margin: 0; font-size: 1.1rem; opacity: 0.95;'>
68                 Create realistic, statistically validated datasets for testing, development, and machine learning applications
69             </p>
70         </div>
71     """, unsafe_allow_html=True)
72
73     # Key Statistics Section
74     st.markdown("""
75     <div style='text-align: left; margin-bottom: 3rem;'>
76         <div style='color: #2c3e50; margin-bottom: 2rem; font-size: 2rem; font-weight: 600;'>Why Choose Our Data Generator?</div>
77     </div>
78     """, unsafe_allow_html=True)
79
80     col1, col2, col3 = st.columns(3)
81
82     with col1:
83         st.markdown("""
84         <div style='
85             background: white;
86             padding: 2.5rem 1.5rem;
87             border-radius: 15px;
88             text-align: center;
89             box-shadow: 0 8px 25px rgba(0,0,0,0.08);
90             border: 1px solid #f0f0f0;
91             height: 100%;
92         '>
93             <div style='font-size: 3rem; font-weight: 800; color: #6677aa; margin-bottom: 1rem;'>100+</div>
94             <div style='color: #5d6d7e; font-weight: 600; font-size: 1.1rem;'>Data Types Supported</div>
95         </div>
96         """, unsafe_allow_html=True)
97
98     with col2:
99         st.markdown("""
100         <div style='
101             background: white;
102             padding: 2.5rem 1.5rem;
103             border-radius: 15px;
104             text-align: center;
105             box-shadow: 0 8px 25px rgba(0,0,0,0.08);
106             border: 1px solid #f0f0f0;
107             height: 100%;
108         '>
109             <div style='font-size: 3rem; font-weight: 800; color: #f0932b; margin-bottom: 1rem;'>50+</div>
110             <div style='color: #5d6d7e; font-weight: 600; font-size: 1.1rem;'>Faster Generation</div>
111         </div>
112         """, unsafe_allow_html=True)
113
114     with col3:
115         st.markdown("""
116         <div style='
117             background: white;
118             padding: 2.5rem 1.5rem;

```

Figure 20: Column separation output

```

104 # Features Section
105 st.markdown("""
106 <div style='text-align: left; margin-bottom: 3rem;'>
107     <div style='color: #2c3e50; margin-bottom: 2rem; font-size: 2rem; font-weight: 600;'>Core Features</div>
108 </div>
109 """, unsafe_allow_html=True)
110
111 col1, col2, col3 = st.columns(3)
112
113 with col1:
114     st.markdown("""
115     <div style='
116         background: white;
117         padding: 2.5rem;
118         border-radius: 15px;
119         box-shadow: 0 10px 30px rgba(0,0,0,0.1);
120         border: 1px solid #e0e0e0;
121         height: 100%;
122         transition: transform 0.3s ease;
123     '>
124         <div style='font-size: 2.5rem; color: #6677aa; margin-bottom: 1.5rem; font-weight: 600;'>100+</div>
125         <div style='color: #2c3e50; margin-bottom: 1rem; font-size: 1.3rem; font-weight: 600;'>CSV Data Extension</div>
126         <p style='color: #5d6d7e; line-height: 1.6; margin: 0; font-size: 0.95rem;'>
127             Upload CSV files and generate extended synthetic data that maintains original statistical properties and distributions.
128         </p>
129     </div>
130     """, unsafe_allow_html=True)
131
132 with col2:
133     st.markdown("""
134     <div style='
135         background: white;
136         padding: 2.5rem;
137         border-radius: 15px;
138         box-shadow: 0 10px 30px rgba(0,0,0,0.1);
139         border: 1px solid #e0e0e0;
140         height: 100%;
141         transition: transform 0.3s ease;
142     '>
143         <div style='font-size: 2.5rem; color: #f0932b; margin-bottom: 1.5rem; font-weight: 600;'>50+</div>
144         <div style='color: #2c3e50; margin-bottom: 1rem; font-size: 1.3rem; font-weight: 600;'>Custom Data Creation</div>
145         <p style='color: #5d6d7e; line-height: 1.6; margin: 0; font-size: 0.95rem;'>
146             Create high-quality synthetic datasets from scratch using AI-powered column detection and intelligent data generation.
147         </p>
148     </div>
149     """, unsafe_allow_html=True)
150
151 with col3:
152     st.markdown("""
153     <div style='
154         background: white;
155         padding: 2.5rem;
156         border-radius: 15px;
157         box-shadow: 0 10px 30px rgba(0,0,0,0.1);
158         border: 1px solid #e0e0e0;
159         height: 100%;
160         transition: transform 0.3s ease;
161     '>
162         <div style='font-size: 2.5rem; color: #6677aa; margin-bottom: 1.5rem; font-weight: 600;'>100+</div>
163         <div style='color: #2c3e50; margin-bottom: 1rem; font-size: 1.3rem; font-weight: 600;'>MySQL Database Integration</div>
164         <p style='color: #5d6d7e; line-height: 1.6; margin: 0; font-size: 0.95rem;'>
165             Connect to MySQL databases, generate synthetic data based on table schemas, and insert directly into databases.
166         </p>
167     </div>
168     """, unsafe_allow_html=True)

```

Figure 21: Markdowns

- Integrate visual separators and typography for a professional, data-centric dashboard experience.

### Activity 3.2: Input System and Schema Configuration

Added text area for user input including:

- Column names and types format for CSV extension page

```

384 def csv_extension_page():
385     st.markdown("---")
386
387     uploaded_file = st.file_uploader("Upload CSV File", type='csv')
388     st.markdown(" ")
389
390     if uploaded_file is not None:
391         with open("temp_upload.csv", "wb") as f:
392             f.write(uploaded_file.getbuffer())
393
394         original_df = pd.read_csv("temp_upload.csv")
395
396         col1, col2 = st.columns(2)
397
398         with col1:
399             st.metric("Original Rows", len(original_df))
400         with col2:
401             st.metric("Columns", len(original_df.columns))
402
403         with st.expander("Original Data Preview", expanded=True):
404             st.dataframe(original_df.head(), use_container_width=True)
405
406         st.markdown("### Generation Settings")
407         col1, col2 = st.columns(2)
408
409         with col1:
410             num_rows = st.number_input("Number of Synthetic Rows", min_value=1, value=200, step=100)
411         with col2:
412             output_file = st.text_input("Output Filename", value="synthetic_data.csv")
413
414         validate = st.checkbox("Run Validation", value=True)
415
416         if st.button("Generate Synthetic Data", type="primary", use_container_width=True):
417             with st.spinner("Generating synthetic data..."):
418                 synthetic_df, original_df_full = st.session_state.generator.generate_from_csv(
419                     "temp_upload.csv",
420                     num_rows,
421                     validate=validate
422                 )
423
424             if synthetic_df is not None:
425                 st.session_state.synthetic_data = synthetic_df
426
427             if validate and original_df_full is not None:
428                 st.session_state.validator.validate_synthetic_data(original_df_full, synthetic_df)
429
430         if st.session_state.synthetic_data is not None:
431             st.markdown("### Download Generated Data")
432             csv = st.session_state.synthetic_data.to_csv(index=False)
433             st.download_button(
434                 "Download CSV",
435                 data=csv,
436                 file_name=output_file,
437                 mime="text/csv",
438                 use_container_width=True
439             )
440         else:
441             st.error("Please enter at least one column name.")
442     else:
443         st.error("Please enter column names.")

```

Figure 22: Input to CSV extension

- Row column name for generation for Dummy data creation

```

379 def dummy_data_creation_page():
380     st.markdown("---")
381
382     columns_input = st.text_area(
383         "Enter Column Names (one per line or comma-separated)",
384         height=80,
385         placeholder="name, email, age, salary, city, phone_number"
386     )
387
388     col1, col2 = st.columns(2)
389
390     with col1:
391         num_rows = st.number_input("Number of Rows", min_value=1, value=200, step=100)
392
393     with col2:
394         output_file = st.text_input("Output Filename", value="dummy_data.csv")
395
396     if st.button("Generate dummy Data", type="primary", use_container_width=True):
397         if columns_input:
398             columns = []
399             if ',' in columns_input:
400                 columns = [col.strip() for col in columns_input.split(',') if col.strip()]
401             else:
402                 columns = [col.strip() for col in columns_input.split('\n') if col.strip()]
403
404             if columns:
405                 with st.spinner(f"Generating {num_rows} rows with {len(columns)} columns..."):
406                     synthetic_df = st.session_state.generator.generate_from_columns(columns, num_rows)
407                     st.session_state.synthetic_data = synthetic_df
408
409             if st.session_state.synthetic_data is not None:
410                 st.markdown("### Download Generated Data")
411                 csv = st.session_state.synthetic_data.to_csv(index=False)
412                 st.download_button(
413                     "Download CSV",
414                     data=csv,
415                     file_name=output_file,
416                     mime="text/csv",
417                     use_container_width=True
418                 )
419             else:
420                 st.error("Please enter at least one column name.")
421         else:
422             st.error("Please enter column names.")

```

Figure 23: Input to the Dummy data creation

- Database details (host, user, password, table)

```

442 def mysql_operations_page():
443
444     with col1:
445         host = st.text_input("Host", value="localhost")
446         password = st.text_input("Password", type="password")
447     with col2:
448         user = st.text_input("Username", value="root")
449         database = st.text_input("Database", value="test")
450
451     if st.button("Connect to Database", use_container_width=True):
452         with st.spinner("Connecting to database..."):
453             connection = st.session_state.db_handler.connect_to_mysql(host, user, password, database)
454             if connection:
455                 st.session_state.db_connection = connection
456                 st.session_state.db_tables = st.session_state.db_handler.get_mysql_tables(connection)
457
458     if 'db_connection' in st.session_state and st.session_state.db_connection:
459         st.success("Connected to database!")
460
461     if st.session_state.db_tables:
462         st.subheader("Available Tables")
463         selected_table = st.selectbox("Select Table", st.session_state.db_tables)
464
465         if selected_table:
466             # Show table schema
467             schema = st.session_state.db_handler.get_table_schema(st.session_state.db_connection, selected_table)
468
469             st.subheader("Table Schema")
470             schema_df = pd.DataFrame.from_dict(schema, orient='index')
471             st.dataframe(schema_df, use_container_width=True)
472
473             num_rows = st.number_input("Number of Rows", min_value=1, value=100, key="mysql_rows")
474             validate = st.checkbox("Run Validation", value=True, key="mysql_validate")
475             insert_db = st.checkbox("Insert into Database", value=False)
476
477             if st.button("Generate from Table Schema", type="primary", use_container_width=True):
478                 with st.spinner(f"Generating data for table '{selected_table}'..."):
479                     synthetic_df, original_df = st.session_state.db_handler.generate_from_mysql_table(
480                         st.session_state.db_connection,
481                         selected_table,
482                         num_rows
483                     )
484                     st.session_state.synthetic_data = synthetic_df
485
486                     if synthetic_df is not None:
487                         if validate and original_df is not None and len(original_df) > 0:
488                             st.info(f"Validating against {len(original_df)} original rows...")
489                             st.session_state.validator.validate_synthetic_data(original_df, synthetic_df)
490
491                         if insert_db:
492                             success = st.session_state.db_handler.insert_to_mysql_table(
493                                 st.session_state.db_connection,
494                                 selected_table,
495                                 synthetic_df
496                             )

```

Figure 24: Database detail input section

- Integrate progress indicators and success messages using Streamlit's feedback components (st.spinner, st.success, etc.).
- Persist inputs via session state for consistent user experience.

```

# Initialize session state
if 'generator' not in st.session_state:
    st.session_state.generator = UniversalDataGenerator()
if 'db_handler' not in st.session_state:
    st.session_state.db_handler = DatabaseHandler(st.session_state.generator)
if 'validator' not in st.session_state:
    st.session_state.validator = DataValidator()
if 'synthetic_data' not in st.session_state:
    st.session_state.synthetic_data = None

```

Figure 25: Session State

### Activity 3.3: Data Visualization and Result Display

- Display generated synthetic data dynamically using **Pandas DataFrames** in Streamlit tables.

```

304 def csv_extension_page():
305     st.markdown("=====")
306
307     uploaded_file = st.file_uploader("Upload CSV File", type='csv')
308     st.markdown(" ")
309
310     if uploaded_file is not None:
311         with open("temp_upload.csv", "wb") as f:
312             f.write(uploaded_file.getbuffer())
313
314         original_df = pd.read_csv("temp_upload.csv")
315
316         col1, col2 = st.columns(2)
317
318         with col1:
319             st.metric("Original Rows", len(original_df))
320         with col2:
321             st.metric("Columns", len(original_df.columns))
322
323         with st.expander("Original Data Preview", expanded=True):
324             st.dataframe(original_df.head(), use_container_width=True)
325
326         st.markdown("### Generation Settings")
327         col1, col2 = st.columns(2)
328
329         with col1:
330             num_rows = st.number_input("Number of Synthetic Rows", min_value=1, value=200, step=100)
331         with col2:
332             output_file = st.text_input("Output Filename", value="synthetic_data.csv")
333
334         validate = st.checkbox("Run Validation", value=True)
335
336         if st.button("Generate Synthetic Data", type="primary", use_container_width=True):
337             with st.spinner("Generating synthetic data..."):
338                 synthetic_df, original_df_full = st.session_state.generator.generate_from_csv(
339                     "temp_upload.csv",
340                     num_rows,
341                     validate=False
342                 )
343
344                 if synthetic_df is not None:
345                     st.session_state.synthetic_data = synthetic_df
346
347                     if validate and original_df_full is not None:
348                         st.session_state.validator.validate_synthetic_data(original_df_full, synthetic_df)
349
350                 if st.session_state.synthetic_data is not None:
351                     st.markdown("### Download Generated Data")
352                     csv = st.session_state.synthetic_data.to_csv(index=False)
353                     st.download_button(
354                         label="Download CSV",
355                         data=csv,
356                         file_name=output_file,
357                         mime="text/csv",
358                         use_container_width=True
359                     )

```

Figure 26: Streamlit Output section for CSV extension

- Integrate Matplotlib and Plotly for generating statistical charts histograms with correlation of the inputted dataset.

```

94 def plot_distributions(self, original_df: pd.DataFrame, synthetic_df: pd.DataFrame, common_columns: set):
95     st.subheader("Distribution Comparison")
96     for col in common_columns:
97         if pd.api.types.is_numeric_dtype(original_df[col]):
98             fig = go.Figure()
99             fig.add_trace(go.Histogram(
100                 x=original_df[col].dropna(),
101                 name='Original',
102                 opacity=0.7,
103                 nbinsx=20
104             ))
105             fig.add_trace(go.Histogram(
106                 x=synthetic_df[col].dropna(),
107                 name='Synthetic',
108                 opacity=0.7,
109                 nbinsx=20
110             ))
111             fig.update_layout(
112                 title=f'Distribution of {col}',
113                 xaxis_title=col,
114                 yaxis_title='Frequency',
115                 barmode='overlay'
116             )
117             st.plotly_chart(fig, use_container_width=True)

```

Figure 27: Plotting the hist graph

- Show computed metrics like mean, standard deviation, and column count in summary cards.

```

327     if uploaded_file is not None:
328         with open("temp_upload.csv", "wb") as f:
329             f.write(uploaded_file.getbuffer())
330
331         original_df = pd.read_csv("temp_upload.csv")
332
333         col1, col2 = st.columns(2)
334
335         with col1:
336             st.metric("Original Rows", len(original_df))
337         with col2:
338             st.metric("Columns", len(original_df.columns))
339
340         with st.expander("Original Data Preview", expanded=True):
341             st.dataframe(original_df.head(), use_container_width=True)
342
343         st.markdown("### Generation Settings")
344         col1, col2 = st.columns(2)
345
346         with col1:
347             num_rows = st.number_input("Number of Synthetic Rows", min_value=1, value=200, step=100)
348         with col2:
349             output_file = st.text_input("Output Filename", value="synthetic_data.csv")
350
351         validate = st.checkbox("Run Validation", value=True)
352
353         if st.button("Generate Synthetic Data", type="primary", use_container_width=True):
354             with st.spinner("Generating synthetic data..."):
355                 synthetic_df, original_df_full = st.session_state.generator.generate_from_csv(
356                     "temp_upload.csv",
357                     num_rows,
358                     validate=False
359                 )
360
361                 if synthetic_df is not None:
362                     st.session_state.synthetic_data = synthetic_df
363
364                 if validate and original_df_full is not None:
365                     st.session_state.validator.validate_synthetic_data(original_df_full, synthetic_df)

```

Figure 28: Evaluation Matrix

## MILESTONE 4: Testing, Optimization, and Final Deployment

This final milestone ensures reliability, speed, and stability of the Synthetic Data Factory through rigorous front-end testing, optimization, and deployment preparation.

### Activity 4.1: Functional Testing

- Test all UI components including buttons, input forms, chart rendering, and data preview tables.

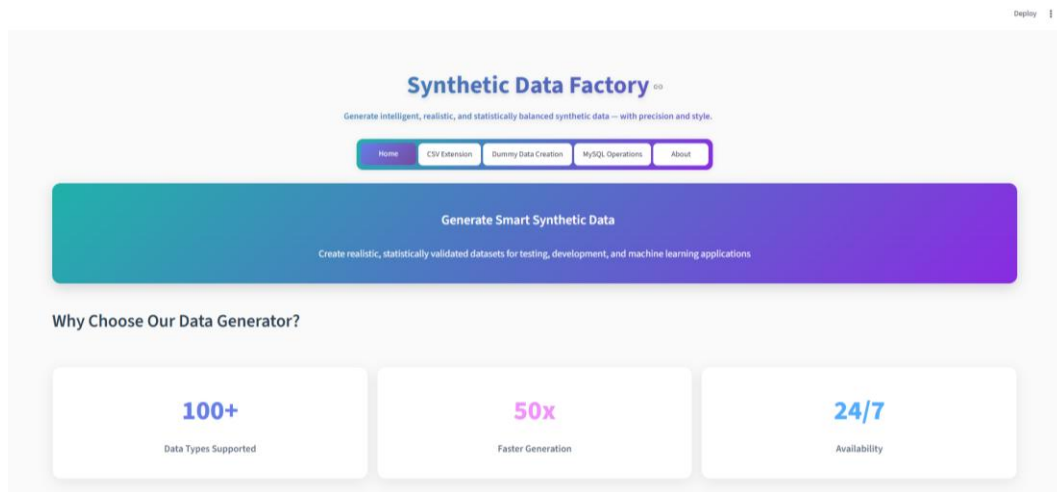


Figure 29: Home Page

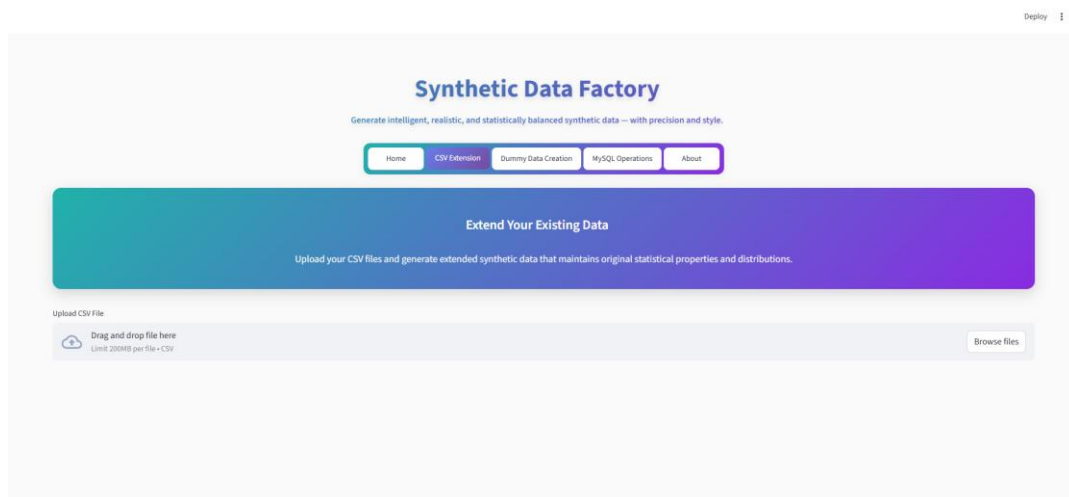


Figure 30: CSV Extension Page

## Synthetic Data Factory

Generate intelligent, realistic, and statistically balanced synthetic data — with precision and style.

Home CSV Extension **Dummy Data Creation** MySQL Operations About

### Create Custom Datasets

Enter column names and let our intelligent system automatically detect the best data type for each column.

Enter Column Names (one per line or comma-separated)

name, email, age, salary, city, phone\_number

Number of Rows

200

Output Filename

dummy\_data.csv

Generate dummy Data

Figure 31: Dummy Data Creation Page

## Synthetic Data Factory

Generate intelligent, realistic, and statistically balanced synthetic data — with precision and style.

Home CSV Extension Dummy Data Creation **MySQL Operations** About

### Database Integration

Connect to your MySQL database and generate synthetic data based on existing table schemas.

#### Database Connection

Host

localhost

Username

root

Password

Database

test

Connect to Database

Figure 32: SQL Operation Page



- Ensure data generation is accurate and consistent for both CSV and database modes.

Deploy |

Upload CSV File

Drag and drop file here  
 Limit: 200MB per file • CSV

Browse files

dummy\_data.csv 312.0B

Original Rows: 10      Columns: 3

Original Data Preview

	name	age	dob
0	Kimberly Rodriguez		64 1961-09-20
1	Ashley Williams		52 1973-08-09
2	Paul Stone		42 1983-11-28
3	Grace Pearson		20 2005-04-21
4	Larry Hoffman		55 1970-09-17

Generation Settings

Number of Synthetic Rows: 200      Output Filename: synthetic\_data.csv

☒ Run Validation

Generate Synthetic Data

Generation complete!

Original Data Preview

Dataset Info: 10 rows, 3 columns

Figure 33: Data Generation from CSV Extension

Deploy |

Generation complete!

Original Data Preview

Dataset Info: 10 rows, 3 columns

	name	age	dob
0	Kimberly Rodriguez		64 1961-09-20
1	Ashley Williams		52 1973-08-09
2	Paul Stone		42 1983-11-28
3	Grace Pearson		20 2005-04-21
4	Larry Hoffman		55 1970-09-17

Synthetic Data Preview

	name	age	dob
0	Tracie Williams		26 1999-10-01
1	Tracie Williams		30 1995-08-10
2	Collin Brown		34 1991-03-07
3	Christian Campbell		51 1974-05-13
4	Mary Morse		15 2010-08-16
5	Kimberly Rodriguez		36 1989-08-28
6	Grace Pearson		43 1982-02-19
7	Ashley Williams		34 1991-09-04
8	Larry Hoffman		15 2010-09-25
9	Vicki Rodriguez		59 1975-11-12

Validation and Quality Metrics

Quality Score: 99.4/100      Distribution Match: 1/1      Avg Mean Difference: 1.12%

Figure 34: Real VS Synthetic Data

Deploy 1

Enter Column Names (one per line or comma separated)  
 name, age, gender, address

Number of Rows  
200

Output Filename  
dummy\_data.csv

Generate dummy Data

Generation complete!

Generated Data Preview

	name	age	gender	address
0	Ryan Hall	79	Female	88394 Parker Cape Suite 725, North Robert, NC 28364
1	Dale Rangel	36	Female	1708 James Forge, Pktyville, SC 29512
2	Justin Graves	52	Female	155 Nelson, FPO AP 46866
3	Brian Krause	62	Female	2410 Vincent Valley Apt. 145, North Melissa, TN 38281
4	Brenda Edwards	41	Female	8055 Howell Prairie Apt. 353, North Angela, CO 30814
5	Nicholas Jones	46	Male	53653 Payne Mission, Lewismouth, PA 80412
6	Robert Wells	39	Female	2137 Anderson Mount Suite 430, East Patrickherst, SC 80500
7	Matthew Marshall	22	Male	1845 Simon Forge, Luke Richard, MD 80507
8	Denise Christian	58	Female	184 Townsend Isle Suite 280, Port Suzanne, HI 42783
9	Julie Hayes	43	Female	99542 Wallace Thoroughway, East Kaitlyn, CT 07811

Download Generated Data

Download CSV

Figure 35: Data Creation from Dummy Data creation Model

- Validate interactive chart responsiveness across different browsers and resolutions.

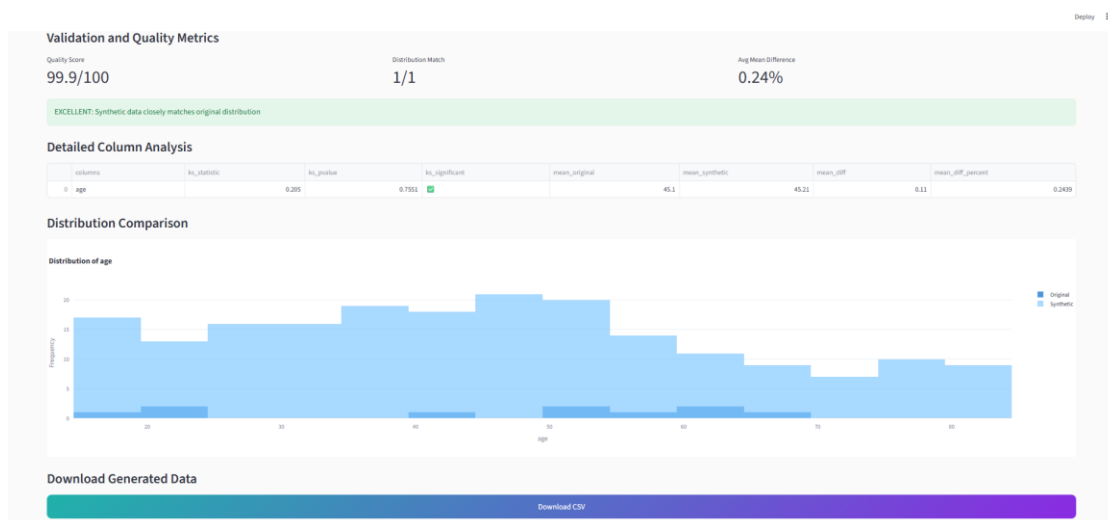
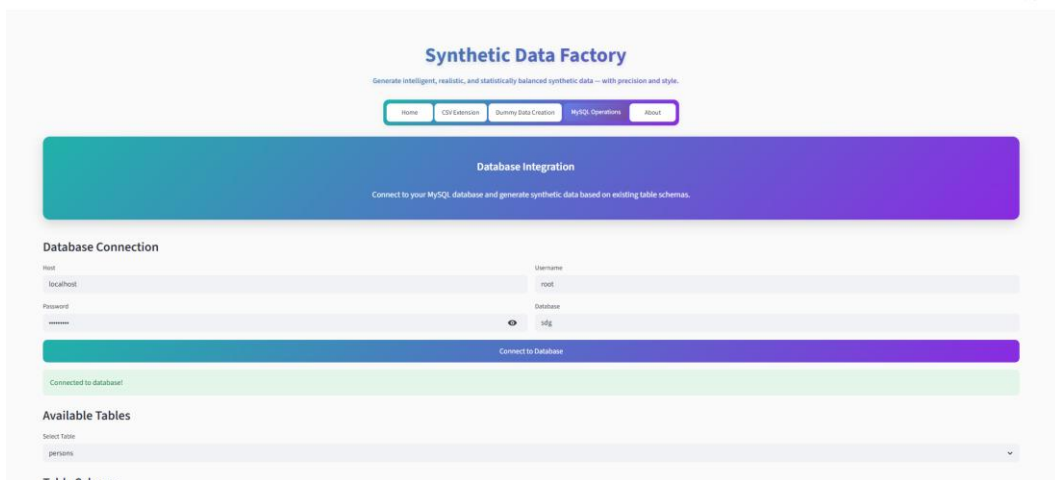


Figure 36: Chart Visualization

## Activity 4.2: Integration and Performance Validation

- Test MySQL connection handling and data insertion with large datasets.



**Synthetic Data Factory**  
 Generate intelligent, realistic, and statistically balanced synthetic data — with precision and style.

Home | CSV Extension | Query Data Creation | **MySQL Operations** | About

**Database Integration**  
 Connect to your MySQL database and generate synthetic data based on existing table schemas.

**Database Connection**

Host: localhost Username: root  
 Password: Password Database: sfsg

Connect to Database

Connected to database!

**Available Tables**  
 Select Table: persons

**Table Schema**

Field	Type	Null	Key	Default	Extra
id	int	NO	PK	None	auto_increment
name	varchar(100)	YES		None	
age	int	YES		None	
sex	varchar(10)	YES		None	
date_of_birth	date	YES		None	

Number of Rows: 100

☒ Run Validation  
☐ Insert into Database

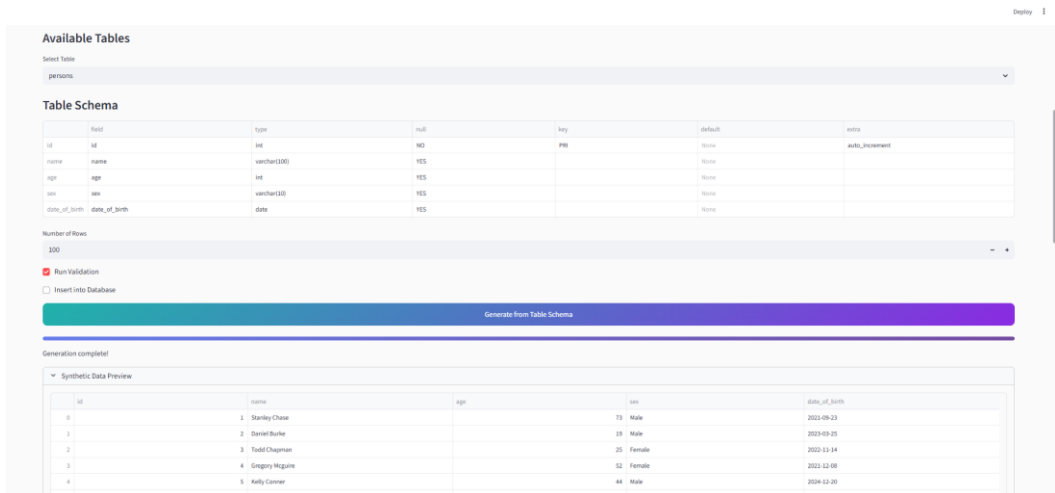
Generate from Table Schema

Generation completed!

**Synthetic Data Preview**

id	name	age	sex	date_of_birth
0	Stanley Chase	73	Male	2023-09-23
1	Daniel Burke	19	Male	2023-03-05
2	Todd Chapman	25	Female	2022-11-14
3	Gregory McGuire	52	Female	2022-12-08
4	Kelly Conner	44	Male	2024-12-20

Figure 37: Connection to the Database



**Available Tables**  
 Select Table: persons

**Table Schema**

Field	Type	Null	Key	Default	Extra
id	int	NO	PK	None	auto_increment
name	varchar(100)	YES		None	
age	int	YES		None	
sex	varchar(10)	YES		None	
date_of_birth	date	YES		None	

Number of Rows: 100

☒ Run Validation  
☐ Insert into Database

Generate from Table Schema

Generation completed!

**Synthetic Data Preview**

id	name	age	sex	date_of_birth
0	Stanley Chase	73	Male	2023-09-23
1	Daniel Burke	19	Male	2023-03-05
2	Todd Chapman	25	Female	2022-11-14
3	Gregory McGuire	52	Female	2022-12-08
4	Kelly Conner	44	Male	2024-12-20

Figure 38: Data Generation from SQL Model

- Configuring the system generation of 10K+ rows efficiently in the Database (MYSQL).

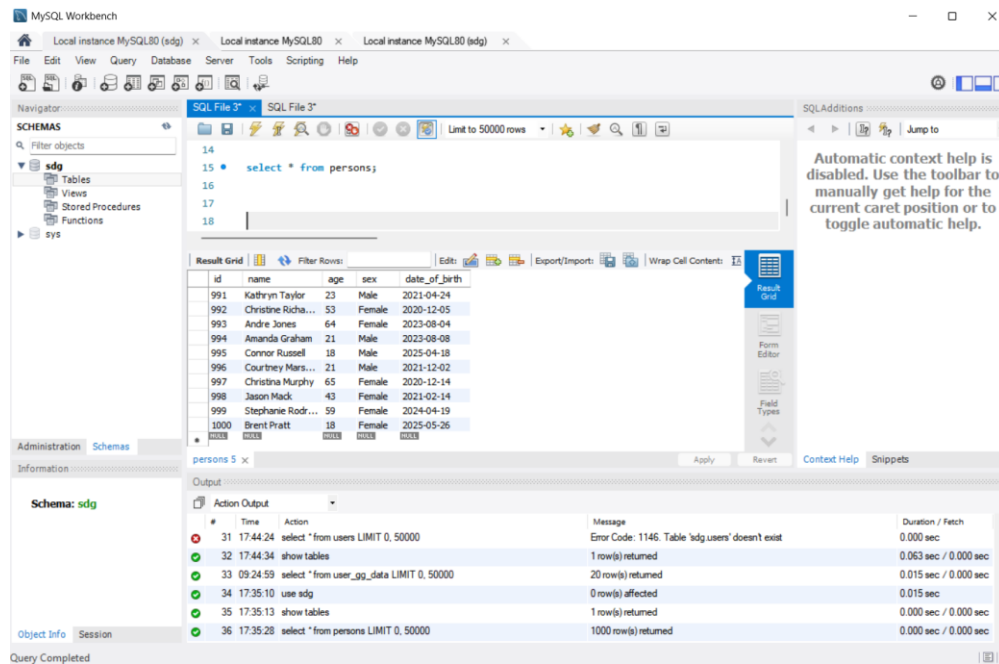


Figure 39: SQL data inserting verification

- Validate that Streamlit's refresh and rerun mechanisms preserve data state correctly with session states.

### Activity 4.3: Deployment Preparation and Final Validation

- The project was prepared for deployment with all dependencies documented in requirements.txt and environment setup instructions completed.
- The application was successfully deployed on Streamlit Cloud (or a local server) with proper configuration to ensure scalability and stability.
- Comprehensive end-to-end validation was performed across all modules, confirming UI consistency, functionality, and seamless integration.
- The deployed environment was verified to maintain data integrity, high performance, and responsiveness under real-world usage conditions.

## Conclusion

The Synthetic Data Factory (SDF) represents a significant advancement in intelligent data generation, transforming how organizations create, validate, and deploy synthetic datasets for testing, development, and analytics. By integrating AI-powered data generation techniques with robust statistical validation and database compatibility, the system ensures that every dataset is realistic, consistent, and tailored to the user's specific requirements. Whether it's extending existing CSV files, creating custom datasets from scratch, or populating MySQL databases, SDF empowers users to produce high-quality synthetic data efficiently and accurately.

Built on a Streamlit interface, the application delivers a seamless and interactive user experience. Users can intuitively navigate between multiple data generation modes, configure datasets, monitor validation metrics, and download or directly insert synthetic data into databases. The combination of Python's data science libraries, intelligent column detection, and database integration makes SDF both technically sophisticated and practically applicable for real-world workflows.

Looking ahead, Synthetic Data Factory has strong potential for further expansion. Features such as automated pattern detection, advanced anomaly generation, integration with additional database systems, and enhanced analytics dashboards could transform SDF into a comprehensive data simulation platform. Ultimately, this project demonstrates how AI-driven synthetic data generation can bridge gaps in data availability, accelerate development cycles, and empower organizations to make data-driven decisions with confidence and precision.