

Comparing Traditional vs. AI-Based Synthetic Data Generation

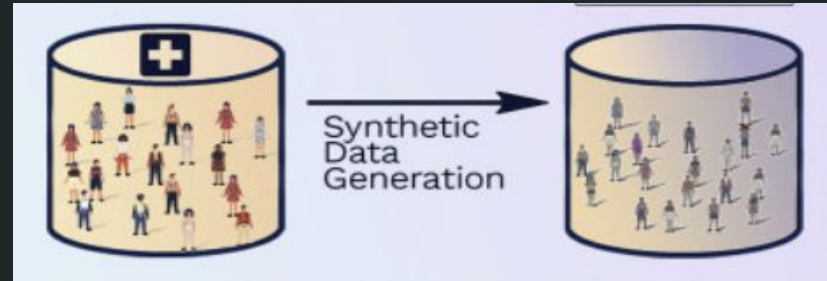
Business Insights and Code Walkthrough

Synthetic data

Refers to information that is artificially generated rather than collected from real-world events.

used in various applications

- Training machine learning models
- Testing Data Pipelines
- Testing software and systems
- Data privacy and security



Traditional Synthetic Data Generation

- **Tools Used:** Faker, NumPy
- **Key Points:**
 - Generates random but realistic names, job titles, dates (using Faker).
 - Salaries within predefined ranges (NumPy).
 - **Limitation:** Manual updates needed for new fields or changes in structure.
 - Lacks flexibility for complex or dynamic datasets.

AI-Based Synthetic Data Generation

- **Key Points:**
 - OpenAI interprets instructions and creates the dataset automatically.
 - No need for manual logic. Automatically handles complex coding and structures.
 - **Natural Language Interface:** Describe datasets in plain English, and the AI generates it.
 - **Result:** JSON output, converted into Pandas DataFrame, saved as CSV.

When to Use Each Method

Key Observations for Professionals

- **When to Use Traditional Methods:**

- **Performance-critical Applications:** Large datasets (>1M rows).
- **Controlled Environments:** Scenarios with strict rules.
- **Cost-sensitive Projects:** No external API costs.

- **When to Use AI-Based Methods:**

- **Prototyping and Exploration:** Quickly test ideas or train ML models.
- **Data with Realism and Context:** For industry-aware data patterns.
- **Flexibility and Speed:** Minimal time to generate data with prompt adjustments.

Challenges with Traditional vs. AI-Based Methods

- **Traditional Method:**
 - Requires detailed programming knowledge.
 - Harder to modify or scale for complex datasets.
 - Limited by Faker's realism.
- **AI-Based Method:**
 - Requires OpenAI API key and incurs costs.
 - Output might need validation or manual review.
 - Large datasets may require iterative generation due to token limits.
 - **Visual:** Table or comparison graphic of challenges for each.

Conclusion

Recommendations

- **For Simplicity and Flexibility**

Use OpenAI for quick, realistic datasets, especially for exploring new fields.

- **For Performance and Large Datasets**

Traditional methods are more cost-effective and scalable for large datasets.

Considerations for Businesses

- **Hybrid Use Case**

Businesses can adopt a **hybrid approach**

- Use AI for quick prototyping and testing ideas, because AI is quick and efficient.
- Switch to traditional methods for scaling and stable implementation."

- **Data Sensitivity**

Traditional methods ensure control and security. AI needs external tools, so validate data use and compliance carefully.

- **Realism Needs**

AI provides unmatched realism for applications like customer behavior analysis.

- **Cost vs. Time Trade-Off:**

Traditional for cost-saving, AI for quicker development.

Walkthrough & Code Example

Traditional Code Example: Faker, NumPy.

AI Method Code Example: OpenAI API, generating dataset with prompt.

Feature	Traditional Method	AI-Based Method
Ease of Setup	Requires coding for each field and manual logic.	Simple, relies on natural language prompts.
Flexibility	Adding fields requires extra coding effort.	Easily modify fields via prompt.
Data Realism	Limited by the randomness and library capabilities.	Context-aware and realistic data.
Scalability	Efficient for very large datasets.	Limited by API token restrictions.
Custom Constraints	Needs coding for complex constraints or patterns.	AI can interpret constraints directly.
Skill Requirement	Requires programming expertise.	Minimal; anyone can write prompts.
Cost	Free (except compute resources).	Paid (depends on OpenAI API usage).
Speed	Fast for simple datasets, slower for complex logic.	Faster for both simple and complex datasets.

