# Indian Government Scheme Tracker - Implementation Report

## Overview

This project, **Indian Government Scheme Tracker**, introduces beginner Data Analysts to the world of **GovTech**, data transparency, and citizen-oriented analytics. It leverages **Python, SQL, and visualization tools** to ingest, clean, and analyze scheme-related datasets. The focus is on understanding **Digital India schemes**, designing dashboards, and presenting insights for better citizen engagement and policy transparency. Minimal ML techniques are included only as an optional extension.

## Goals & Learning Objectives

- Learn how to ingest and clean real-world government datasets.

- Understand the structure of Indian government schemes (Digital India, PMAY, Ayushman Bharat, etc.).

- Use SQL to query and organize citizen–scheme data.

- Perform **EDA (Exploratory Data Analysis)** using Python and Seaborn.

- Create **interactive dashboards** with Streamlit and Power BI.

- Communicate policy insights with simple, clear visualizations.

- Build confidence in data storytelling for public policy applications.

## Applications & Use Cases

- **Citizen transparency**: Present clear dashboards showing scheme beneficiaries, fund allocations, state/district-level performance.

- **Policy monitoring**: Help policymakers track scheme adoption and success rates.

- **Resource allocation**: Identify gaps in fund utilization or scheme coverage.

- **Awareness tools**: Provide citizens easy access to check scheme eligibility, reach, and usage.

- **GovTech research**: Serve as a student project for open data analysis in governance.

## Tech Stack (Modules & Tools)

**Python**

- **Core Libraries**: pandas, numpy
- **Visualization**: seaborn, matplotlib, plotly
- **Dashboard**: streamlit
- **Data Handling**: openpyxl (Excel), csv module

**SQL**

- MySQL or PostgreSQL (students can use SQLite for simplicity)
- SQLAlchemy (optional connector)

**Visualization & Reporting**

- **Power BI**: for rich visuals and drill-down dashboards
- **Streamlit**: for Python-based interactive dashboards

**Dev Tools**

- Jupyter Notebook / VS Code for analysis
- GitHub for version control and project sharing

## How the Project Works (Conceptual Flow)

1. **Data Sources**: Use provided government scheme datasets (CSV/Excel). These may include citizen enrollment, state-wise allocation, beneficiary counts, disbursement status, etc.
2. **Data Ingestion**: Load datasets into Python (pandas) and/or a SQL database.
3. **Data Cleaning**: Handle nulls, duplicates, inconsistent formats (e.g., state names, dates, currency).
4. **EDA & Insights**:
    - Summaries (total beneficiaries, funds allocated, funds utilized)
    - Group analysis (by scheme, state, district, gender, or socio-economic categories)

o Time-series trends (year-on-year adoption, seasonal spikes)

5. **SQL Queries**: Perform joins, aggregations, and subqueries for deeper insights (e.g., average fund utilization by scheme).

6. **Visualization**: Use seaborn/plotly for initial plots (bar, line, heatmap, pie). Then, move to Power BI and Streamlit for professional dashboards.

7. **Dashboards**: Build dashboards that show key metrics (fund allocation vs utilization, beneficiary demographics, regional disparities).

8. **Optional Minimal ML**: Use simple trend projection (linear regression) or outlier detection (boxplot-based) to highlight unusual allocation/beneficiary patterns.

# Step-by-step Implementation Instructions (Student-friendly)

**Phase 0 — Setup**

1. Create a project folder and initialize version control (git).

2. Set up a virtual environment and install dependencies (pandas, numpy, seaborn, matplotlib, streamlit, sqlalchemy).

3. Prepare datasets (CSV/Excel) and place them in a data/ folder.

**Phase 1 — Data Ingestion**

4. Open Jupyter Notebook and load datasets using pandas.

5. Inspect the first few rows and note data types.

6. Create a SQL database (SQLite/MySQL/PostgreSQL). Load cleaned data into SQL tables.

**Phase 2 — Data Cleaning**

7. Standardize column names (snake_case for consistency).

8. Handle missing values: fill with 0 (for numeric) or 'Unknown' (for categorical) where appropriate.

9. Remove duplicates by checking unique keys (e.g., scheme_id + citizen_id).

10. Convert dates into datetime format and extract useful parts (year, month, quarter).

11. Standardize categorical fields (e.g., state names with a reference list).

**Phase 3 — Exploratory Data Analysis (EDA)**

12. Perform descriptive statistics: mean, median, min, max, counts.

13. Group by scheme and state to check fund allocation and utilization.

14. Create seaborn plots:

   o Bar plots of top 10 states by allocation/utilization

   o Line plot of scheme adoption over time

   o Heatmap of state vs scheme utilization

15. Document findings: e.g., "Scheme X has high allocation but low utilization in State Y."

**Phase 4 — SQL Queries**

16. Write queries for:

   o Top 5 schemes by beneficiaries

   o State-wise utilization percentage

   o Schemes with maximum year-on-year growth

   o Citizens covered by multiple schemes (using JOIN)

17. Export query results into CSV for visualization.

**Phase 5 — Visualization (Python)**

18. Create seaborn/plotly visualizations for scheme-wise comparisons.

19. Save figures as PNGs to use in reports.

**Phase 6 — Power BI Dashboard**

20. Import cleaned data into Power BI.

21. Build visuals:

   o KPI cards (total allocation, utilization, beneficiaries)

   o Bar chart: state-wise allocation/utilization

   o Line chart: adoption trend

   o Map: scheme coverage by state/district

22. Add slicers/filters for scheme, year, and region.

**Phase 7 — Streamlit Dashboard**

23. Create a simple Streamlit app.

24. Add interactive filters (scheme, state, year).

25. Display summary statistics and plots dynamically.

26. Allow download of filtered data as CSV.

**Phase 8 — Optional Minimal ML (Beginner-friendly)**

27. Fit a simple linear regression to predict next year's allocation based on past data.

28. Highlight outlier states/schemes using boxplot/IQR methods.

29. Visualize projected trends alongside actuals.

**Phase 9 — Reporting & Documentation**

30. Write a README documenting the project, dataset sources, and insights.

31. Save EDA notebook, SQL query file, and Streamlit app in repo.

32. Record a short video demo of Power BI and Streamlit dashboards.

## Deliverables

1. Cleaned dataset in CSV and SQL format.

2. Jupyter Notebook with EDA.

3. SQL query file with analysis queries.

4. Power BI dashboard (PBIX file).

5. Streamlit dashboard (Python script).

6. README file with documentation.

7. Screenshots of key findings.

## Common Pitfalls & Tips

- **Dataset issues**: Government datasets often have missing/inconsistent values—document your cleaning decisions.

- **Visualization clutter**: Focus on clarity—avoid too many categories in one chart.

- **SQL joins**: Double-check keys to avoid duplicate counts.

- **Dashboard usability**: Ensure filters are intuitive and visuals load quickly.

## Extensions & Advanced Ideas

- Build a citizen-facing app where individuals can check scheme eligibility (prototype with Streamlit).

- Add NLP to parse scheme descriptions and cluster similar schemes.

- Automate updates by connecting directly to government open data portals (using APIs if available).

- Compare performance across multiple schemes using multi-index analysis in pandas.