# Synthetic Data Generation for Soil Data

## Phase 2: Project Execution and Demonstration

## 1. Project Title:

Synthetic Data Generation for Soil Data

## 2. Objective Recap:

The objective of this project is to **generate synthetic data for soil properties** using the **CTGAN (Conditional Tabular Generative Adversarial Networks)** model. This aims to address privacy concerns associated with sharing real soil data while providing a realistic dataset for research and development purposes.

## 3. Technologies Used:

- Python
- SDV (Synthetic Data Vault) library
- Pandas
- Plotly
- Scikit-learn
- Google Colab

## 4. Full Code Implementation:

### Step 1: Install Necessary Libraries

```
pip install sdv
```

### Step 2: Data Loading and Preprocessing

```
import pandas as pd
df = pd.read_csv('soil.csv')

from sklearn.preprocessing import LabelEncoder
lr = LabelEncoder()
df['District_id'] = lr.fit_transform(df['District '])
df = df.drop('District ',axis =1)
```

### Step 3: Metadata Definition and Validation

```
from sdv.metadata import SingleTableMetadata
metadata = SingleTableMetadata()
metadata.detect_from_dataframe(df)
metadata.validate_data(data=df)
```

### Step 4: Synthesizer Training

```
from sdv.single_table import CTGANSynthesizer
synthesizer = CTGANSynthesizer(metadata, enforce_rounding=False,
epochs=3000, verbose=True)
synthesizer.fit(df)
```

**College Name:** VIT Bhopal
**Student Name:** Garvit

**Step 5: Synthetic Data Generation**

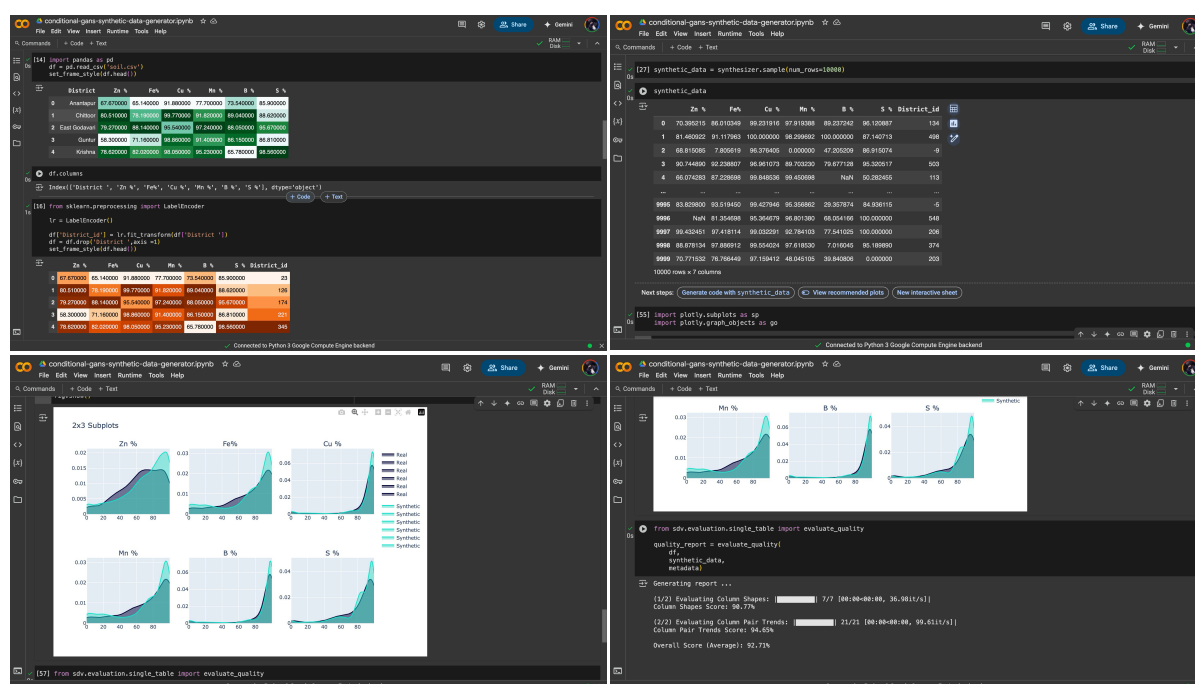```
synthetic_data = synthesizer.sample(num_rows=10000)
```

**Step 6: Evaluation and Visualization**

```
from sdv.evaluation.single_table import get_column_plot, evaluate_quality

# Generate plots for comparing real and synthetic data distributions
# ... (Code for creating Plotly subplots) ...

# Evaluate the quality of the synthetic data
quality_report = evaluate_quality(df, synthetic_data, metadata)
```

# 5. Output Screenshots:



# 6. Conclusion:

This project successfully generated **synthetic soil data** using the **CTGAN** model. The evaluation results and visualizations demonstrate the quality and similarity of the synthetic data to the original dataset. This synthetic data can be used for research and development while protecting the privacy of sensitive information in the original soil data.

# 7. References:

- SDV (Synthetic Data Vault) documentation
- CTGAN research paper
- Pandas documentation
- Plotly documentation
- Scikit-learn documentation