| GLOBAL ECONOMIC INDICATORS — MULTIVARIATE ANALYSIS & PREDICTIVE MODELING

Slide 2 — Research Objectives & Data Overview

- **Objectives:**
- Identify global macroeconomic patterns through PCA.
- Classify countries based on fiscal and economic characteristics.
- Predict government revenue (% GDP) using machine learning regression models.
- **Dataset:**
- 14 macro indicators (2010–2024).
- Source: World Bank, IMF, OECD.
- Variables: GDP, CPI, Unemployment, Tax, Interest Rate, Gov. Revenue, etc.
- **Visuals:**
- Insert data summary table / descriptive statistics (from your Python EDA output).
- Add a correlation heatmap (Seaborn sns.heatmap) to show variable relationships.

Slide 3 — Exploratory Data Analysis (EDA)

- **Key EDA Insights:**
- High correlation: Government revenue ↔ government expenditure ↔ tax revenue.
- Moderate correlation: GDP & GNI with fiscal variables.
- Negative correlation: Inflation, interest rates vs. fiscal strength.
- Distribution: Fiscal indicators right-skewed; GDP per capita varies widely.
- **Visuals:**
- Insert:
 - Correlation matrix heatmap.
- Pairplot (or scatterplot matrix).
- Histogram of key variables (GDP, CPI, Gov. Revenue).

Slide 4 — PCA Overview

- **Goal:** Simplify 14 correlated variables → uncorrelated principal components.
- **Steps:** Standardization → PCA → retain top 3 components explaining ~70% variance.

| PC1 | GDP, GNI, Gov. Expense, Tax, Revenue | Fiscal & Economic Scale | ~40% |

| PC2 | Inflation, GDP Growth, Current Account | Monetary & External Balance | ~18% |

| PC3 | Inflation (CPI & Deflator), Unemployment | Price-Pressure & Labor Slack | ~12% |

- **Visuals:**
- Insert: Scree plot (variance explained by PCs).
- Insert: PCA loading plot (PC1–PC2 or PC1–PC3 biplot).

Slide 5 — Interpretation of Principal Components

- **PC1:** Represents fiscal scale and macroeconomic size.
- **PC2:** Captures monetary & inflation dynamics vs. fiscal stability.
- **PC3:** Reflects inflation-employment interactions.

Key takeaway:

Fiscal capacity (PC1) is the dominant structural factor differentiating economies.

- **Visuals:**
- Insert: Loading vectors (bar chart of top features per PC).

Slide 6 — PCA & Government Revenue Correlation

- **Government Revenue (% GDP):**
- Correlated positively with PC1 (r = 0.53) and PC3 (r = 0.47).
- Correlated negatively with PC2 (r = -0.39).

Interpretation:

Large, stable economies with fiscal capacity → higher revenue. Inflationary volatility (high PC2) reduces fiscal sustainability.

- **Visuals:**
- Insert: Scatter plot of Gov. Revenue vs PC1, PC2, PC3.
- Optional: Bar chart of correlation coefficients.

Slide 7 — PCA Visualization (3D + Quantiles)

- **3D PCA scatter (colored by Gov. Revenue % GDP):**
- Low revenue (<20%) → negative PC1.
- Medium (20–30%) → moderate PC1.
- High (>30%) → positive PC1, slightly negative PC3.
- **Quantile bands:** Smooth fiscal gradient along PC1.
- **Visuals:**
- Insert: 3D PCA scatter plot (the 4 screenshots you attached earlier).
- Highlight: "Fiscal Capacity Axis (PC1)" label on the plot.

Slide 8 — K-Means Clustering

- **Goal:** Classify economies into structural groups.
- **Optimal k:** 2 (from elbow method).
- **Silhouette Score:** 0.365 → moderate separation.
- **Clusters:**
- **Cluster 0 (orange):** High-revenue, stable, developed (+PC1, -PC3).
- **Cluster 1 (blue):** Low-revenue, inflation-sensitive (-PC1, -PC2, -PC3).

- **Visuals:**
- Insert:
 - Elbow plot (inertia vs. k).
 - 3D cluster plot (color-coded).
- Use your cluster screenshots (e.g., "Digital Segments Map" style).

Slide 9 — Regression Models Comparison

Models Tested: Linear | Ridge | Lasso | ElasticNet | Decision Tree | Random Forest | Gradient Boosting | SVR | PCR

Metrics: MSE, MAE, R2, 5-Fold CV

| Rank | Model | Test MSE ↓ | Test R² ↑ | CV R² ↑ | Notes |

|-----|------|-----|

- | 1 | Random Forest | 4.43 | 0.93 | 0.92 | Best accuracy & stability |
- | 2 | Gradient Boosting | 6.53 | 0.90 | 0.89 | Slightly lower stability |
- | 3 | Decision Tree | 11.04 | 0.83 | 0.84 | Simple, interpretable |
- | 4 | SVR | 12.46 | 0.81 | 0.80 | Non-linear, moderate bias |
- | 5-9 | Linear, Ridge, Lasso, ElasticNet, PCR | 13-18 | 0.70-0.80 | Underfit linear data |
- **Visuals:**
- Insert: Horizontal bar chart comparing Test R² or MSE across models.

Slide 10 — Regression Interpretation

- **Linear & Ridge:** Full interpretability, stable but underfit.
- **Lasso & ElasticNet:** Shrink less-important features → key fiscal drivers remain (Tax, Gov. Expense, Current Account).
- **Random Forest & GBM:** Capture non-linear fiscal-monetary relationships.
- **PCR Equation:**

\$\$

 $hat{Y} = 24.66 + 2.52 PC1 - 2.19 PC2 + 2.90 PC3, \quad R^2 = 0.66$

- \rightarrow 66% of fiscal performance explained by first three macro factors.
- **Visuals:**
- Insert: Bar plot of linear model coefficients (optional).
- Add: Predicted vs Actual scatter for Random Forest or PCR.

Slide 11 — Integrated Discussion

- **PCA:** Fiscal scale (PC1) drives most cross-country variance.
- **Regression:** Non-linear models outperform linear → confirms complex interactions.
- **Clustering:** Distinct fiscal regimes → developed vs. developing economies.
- **PCR Insight:**

```
$$
```

 $hat{Y} = 24.66 + 2.52 PC1 - 2.19 PC2 + 2.90 PC3, \quad R^2 = 0.66$

- → +PC1 → higher revenue (fiscal strength).
- -PC2 → less volatility.
- +PC3 → moderate inflation aids nominal capacity.
- **Visuals:**
- Infographic: PC1-PC3 axes with arrows showing direction of positive effects.

Slide 12 — Policy & Strategic Recommendations

- 1 ** Strengthen Fiscal Capacity**
 - Broaden tax bases, digitalize collection, align expenditure with revenue.
- 2 ** Promote Sustainable Growth**
 - Focus on productivity & private sector diversification.
- 3 ** Maintain Monetary Stability **
 - Coordinate fiscal & monetary policies; control inflation.
- 4 ** Manage Employment & Inflation **
 - Counter-cyclical fiscal tools; labor reforms; price stability.
- 5 **Cluster-Specific Strategy**
 - **Cluster 0:** Efficiency, innovation, debt sustainability.
 - **Cluster 1:** Tax reform, governance, diversification.
- 6 **Overall Message:**

Economies balancing fiscal strength (PC1), monetary stability (low PC2), and controlled inflation (PC3) achieve the most sustainable revenue performance.

- **Visuals:**
- Use icons (is Fiscal, is Growth, Stability, gentlement).
- Add map or cluster chart for context.

Slide 13 — Conclusion

- **PCA, clustering, and regression jointly reveal:**
- Fiscal capacity is the strongest structural differentiator.
- Inflation control enhances fiscal sustainability.
- Data-driven modeling can predict fiscal performance (R² ≈ 0.93 with RF).

Final takeaway:

Strong, diversified, inflation-stable economies achieve higher and more sustainable government revenue ratios.

📚 METHODOLOGY

1. Data Collection & Preprocessing

- **Data Sources:** World Bank, IMF, OECD.
- **Time Frame:** 2010–2024.
- **Variables:** 14 macroeconomic indicators (GDP, CPI, Unemployment, Tax, Interest Rate, Gov. Revenue, etc.).
- **Cleaning:** Handle missing values, outliers, and ensure consistent units.
- **Standardization:** Scale features to zero mean and unit variance for PCA and regression.

2. Exploratory Data Analysis (EDA)

- **Descriptive Statistics:** Compute mean, median, standard deviation, min, max for each variable.
- **Correlation Analysis:** Generate a correlation matrix and heatmap to identify relationships.
- **Distribution Analysis:** Plot histograms and pairplots to understand data distribution and relationships.

3. Principal Component Analysis (PCA)

- **Standardization:** Ensure all variables are on the same scale.
- **PCA Execution:** Apply PCA to reduce dimensionality.
- **Component Selection:** Retain top 3 principal components explaining ~70% of variance.
- **Interpretation:** Analyze loadings to understand the contribution of original variables to each component.

4. Clustering Analysis

- **K-Means Clustering:** Determine the optimal number of clusters using the elbow method.
- **Silhouette Score:** Evaluate cluster quality.
- **Cluster Interpretation:** Analyze cluster characteristics based on principal components and original variables.

5. Regression Modeling

- **Model Selection:** Test multiple regression models including Linear, Ridge, Lasso, ElasticNet, Decision Tree, Random Forest, Gradient Boosting, SVR, and PCR.
- **Model Training:** Use 5-fold cross-validation to assess model performance.
- **Model Evaluation:** Compare models based on MSE, MAE, and R² metrics.
- **Feature Importance:** Analyze feature importance for models like Random Forest and GBM.

6. Model Interpretation & Insights

- **PCA Insights:** Understand the impact of principal components on fiscal capacity and government revenue.
- **Clustering Insights:** Identify distinct fiscal regimes and their characteristics.
- **Regression Insights:** Interpret coefficients and feature importance to understand the drivers of government revenue.

7. Integrated Analysis & Policy Recommendations

- **Synthesis:** Combine insights from PCA, clustering, and regression to form a comprehensive understanding of macroeconomic patterns.
- **Policy Recommendations:** Develop actionable strategies based on the integrated analysis to enhance fiscal capacity, promote sustainable growth, maintain monetary stability, and manage employment and inflation.

8. Conclusion

- **Summary:** Recap the key findings and their implications for macroeconomic policy and economic stability.
- **Future Work:** Suggest areas for further research and analysis.

| GLOBAL ECONOMIC INDICATORS — MULTIVARIATE ANALYSIS & PREDICTIVE MODELING

Slide 1 — Title Slide

- **Title:** Global Economic Indicators: PCA, Clustering, and Regression Analysis (2010–2024)
- **Subtitle:** A Multivariate Study of Fiscal Capacity and Macroeconomic Stability **Visuals:**
- Background: World map or macroeconomic visualization.
- Add your name, course, institution, and date.

Slide 2 — Research Objectives & Data Overview

Objectives:

- Identify global macroeconomic patterns through PCA.
- Classify countries based on fiscal and economic characteristics.
- Predict government revenue (% GDP) using machine learning regression models.
- **Dataset:**
- 14 macro indicators (2010-2024).
- Source: World Bank, IMF, OECD.
- Variables: GDP, CPI, Unemployment, Tax, Interest Rate, Gov. Revenue, etc.
- **Visuals:**
- Insert data summary table / descriptive statistics (from your Python EDA output).
- Add a correlation heatmap (Seaborn sns.heatmap) to show variable relationships.

Slide 3 — Exploratory Data Analysis (EDA)

- **Key EDA Insights:**
- High correlation: Government revenue ↔ government expenditure ↔ tax revenue.
- Moderate correlation: GDP & GNI with fiscal variables.
- Negative correlation: Inflation, interest rates vs. fiscal strength.
- Distribution: Fiscal indicators right-skewed; GDP per capita varies widely.
- **Visuals:**
- Insert:
 - Correlation matrix heatmap.
- Pairplot (or scatterplot matrix).
- Histogram of key variables (GDP, CPI, Gov. Revenue).
- ### Slide 4 PCA Overview
- **Goal:** Simplify 14 correlated variables → uncorrelated principal components.
- **Steps:** Standardization → PCA → retain top 3 components explaining ~70% variance.

```
| Component | Dominant Variables | Theme | Variance Explained |
```

| PC1 | GDP, GNI, Gov. Expense, Tax, Revenue | Fiscal & Economic Scale | ~40% |

| PC2 | Inflation, GDP Growth, Current Account | Monetary & External Balance | ~18% |

| PC3 | Inflation (CPI & Deflator), Unemployment | Price-Pressure & Labor Slack | ~12% |

- **Visuals:**
- Insert: Scree plot (variance explained by PCs).
- Insert: PCA loading plot (PC1–PC2 or PC1–PC3 biplot).
- ### Slide 5 Interpretation of Principal Components
- **PC1:** Represents fiscal scale and macroeconomic size.
- **PC2:** Captures monetary & inflation dynamics vs. fiscal stability.
- **PC3:** Reflects inflation-employment interactions.
- **Key takeaway:**

Fiscal capacity (PC1) is the dominant structural factor differentiating economies.

- **Visuals:**
- Insert: Loading vectors (bar chart of top features per PC).
- ### Slide 6 PCA & Government Revenue Correlation
- **Government Revenue (% GDP):**
- Correlated positively with PC1 (r = 0.53) and PC3 (r = 0.47).
- Correlated negatively with PC2 (r = -0.39).
- **Interpretation:**

Large, stable economies with fiscal capacity → higher revenue. Inflationary volatility (high PC2) reduces fiscal sustainability.

```
**Visuals:**
- Insert: Scatter plot of Gov. Revenue vs PC1, PC2, PC3.
- Optional: Bar chart of correlation coefficients.
### Slide 7 — PCA Visualization (3D + Quantiles)
**3D PCA scatter (colored by Gov. Revenue % GDP):**
- Low revenue (<20%) → negative PC1.
- Medium (20–30%) → moderate PC1.
- High (>30%) → positive PC1, slightly negative PC3.
**Quantile bands:** Smooth fiscal gradient along PC1.
**Visuals:**
- Insert: 3D PCA scatter plot (the 4 screenshots you attached earlier).
- Highlight: "Fiscal Capacity Axis (PC1)" label on the plot.
### Slide 8 — K-Means Clustering
**Goal:** Classify economies into structural groups.
**Optimal k:** 2 (from elbow method).
**Silhouette Score:** 0.365 → moderate separation.
**Clusters:**
- **Cluster 0 (orange):** High-revenue, stable, developed (+PC1, -PC3).
- **Cluster 1 (blue):** Low-revenue, inflation-sensitive (-PC1, -PC2, -PC3).
**Visuals:**
- Insert:
 - Elbow plot (inertia vs. k).
 - 3D cluster plot (color-coded).
- Use your cluster screenshots (e.g., "Digital Segments Map" style).
### Slide 9 — Regression Models Comparison
**Models Tested:** Linear | Ridge | Lasso | ElasticNet | Decision Tree | Random Forest |
Gradient Boosting | SVR | PCR
**Metrics:** MSE, MAE, R2, 5-Fold CV
| Rank | Model | Test MSE ↓ | Test R<sup>2</sup> ↑ | CV R<sup>2</sup> ↑ | Notes |
|-----|-----|-----|-----|
| 1 | Random Forest | 4.43 | 0.93 | 0.92 | Best accuracy & stability |
2 | Gradient Boosting | 6.53 | 0.90 | 0.89 | Slightly lower stability |
| 3 | Decision Tree | 11.04 | 0.83 | 0.84 | Simple, interpretable |
| 4 | SVR | 12.46 | 0.81 | 0.80 | Non-linear, moderate bias |
| 5–9 | Linear, Ridge, Lasso, ElasticNet, PCR | 13–18 | 0.70–0.80 | Underfit linear data |
```

- **Visuals:**
- Insert: Horizontal bar chart comparing Test R2 or MSE across models.

Slide 10 — Regression Interpretation

- **Linear & Ridge:** Full interpretability, stable but underfit.
- **Lasso & ElasticNet:** Shrink less-important features → key fiscal drivers remain (Tax, Gov. Expense, Current Account).
- **Random Forest & GBM:** Capture non-linear fiscal-monetary relationships.
- **PCR Equation:**

\$\$

$$hat{Y} = 24.66 + 2.52 PC1 - 2.19 PC2 + 2.90 PC3, \quad R^2 = 0.66$$

- \rightarrow 66% of fiscal performance explained by first three macro factors.
- **Visuals:**
- Insert: Bar plot of linear model coefficients (optional).
- Add: Predicted vs Actual scatter for Random Forest or PCR.

Slide 11 — Integrated Discussion

- **PCA:** Fiscal scale (PC1) drives most cross-country variance.
- **Regression:** Non-linear models outperform linear → confirms complex interactions.
- **Clustering:** Distinct fiscal regimes → developed vs. developing economies.
- **PCR Insight:**

\$\$

$$hat{Y} = 24.66 + 2.52 PC1 - 2.19 PC2 + 2.90 PC3, \quad R^2 = 0.66$$

- → +PC1 → higher revenue (fiscal strength).
- -PC2 → less volatility.
- +PC3 → moderate inflation aids nominal capacity.
- **Visuals:**
- Infographic: PC1–PC3 axes with arrows showing direction of positive effects.

Slide 12 — Policy & Strategic Recommendations

- 1 **Strengthen Fiscal Capacity**
 - Broaden tax bases, digitalize collection, align expenditure with revenue.
- 2 ** Promote Sustainable Growth **
 - Focus on productivity & private sector diversification.
- 3 ** Maintain Monetary Stability**

- Coordinate fiscal & monetary policies; control inflation.
- 4 ** Manage Employment & Inflation **
 - Counter-cyclical fiscal tools; labor reforms; price stability.
- 5 **Cluster-Specific Strategy**
 - **Cluster 0:** Efficiency, innovation, debt sustainability.
 - **Cluster 1:** Tax reform, governance, diversification.
- 6 **Overall Message:**

Economies balancing fiscal strength (PC1), monetary stability (low PC2), and controlled inflation (PC3) achieve the most sustainable revenue performance.

- **Visuals:**
- Use icons (Kriscal, Kriscal, Kriscal) Fiscal, Kriscal, Kriscal,
- Add map or cluster chart for context.

Slide 13 — Conclusion

- **PCA, clustering, and regression jointly reveal:**
- Fiscal capacity is the strongest structural differentiator.
- Inflation control enhances fiscal sustainability.
- Data-driven modeling can predict fiscal performance (R² ≈ 0.93 with RF).
- **Final takeaway:**

Strong, diversified, inflation-stable economies achieve higher and more sustainable government revenue ratios.

Below are a handful of **GitHub repository concepts** that blend a **semi-formal tone** with a **data-science / machine-learning focus**.

Each idea includes a short description, suggested folder structure, and a few "nice-to-have" features that make the repo look polished and professional.

1. | **Economic-Indicators-ML-Toolkit**

Goal: Provide a reusable, well-documented codebase for analyzing macro-economic time-series (e.g., GDP, inflation, unemployment) and building predictive models for fiscal outcomes.

```
| Folder | Contents | |------| | 'data/` | Raw CSV/Parquet files, cleaned Parquet snapshots, `README` describing sources (World Bank, IMF, OECD). | | 'src/` | Python modules: `data_ingest.py`, `eda.py`, `pca.py`, `clustering.py`, `regression.py`. | | `notebooks/` | Jupyter notebooks for step-by-step tutorials (`01_EDA.ipynb`, `02_PCA.ipynb`, `03_Clustering.ipynb`, `04_Regression.ipynb`). |
```

```
| `models/` | Pickled scikit-learn models, ONNX exports, and a `model_registry.yaml` listing
version, metrics, and training date. |
| `docs/` | Sphinx or MkDocs site with API reference, usage guide, and a "Getting Started"
tutorial.
'tests/' | Unit tests ('pytest') for each module, plus a 'requirements-dev.txt'. |
| `scripts/` | CLI entry points (`train_regression.py`, `run_clustering.py`). |
| `requirements.txt` | Production dependencies. |
| `environment.yml` | Conda environment definition. |
| `README.md` | Project overview, motivation, data provenance, installation, quick-start,
contribution guide. |
| `CONTRIBUTING.md` | Guidelines for pull requests, code style (PEP 8), testing, and
documentation.
| `LICENSE` | MIT or Apache-2.0. |
**Why it works:**
- **Semi-formal tone** – formal documentation + clear, reproducible notebooks.
- **ML-centric** – PCA, clustering, regression, model versioning.
- **Data provenance** – transparent source attribution, reproducible pipelines.
## 2. **Time-Series-Forecast-Competition**
**Goal:** Host a mini-competition where participants forecast a macro-economic variable
(e.g., next-quarter GDP growth) using a public dataset. The repo contains the data, baseline
models, evaluation scripts, and a leaderboard.
| Folder | Contents |
|-----|
| `data/` | Historical quarterly data (`gdp_quarterly.csv`), train/test split (`train.csv`, `test.csv`).
| `baseline/` | Simple baselines: `naive.py` (last-value), `arima.py`, `prophet_baseline.ipynb`.
| 'submissions/' | Sample CSV ('sample submission.csv') and a 'README' on submission
format. |
| `evaluation/` | Python script `evaluate.py` (RMSE, MAE, MAPE) and a `leaderboard.csv`. |
| `notebooks/` | Example notebooks: `01_Exploratory.ipynb`, `02_ARIMA_Tuning.ipynb`,
`03_Prophet_Tuning.ipynb`, `04_DeepAR.ipynb`. |
| `models/` | Serialized models from participants (optional). |
| 'docs/' | Competition rules, scoring rubric, and a "How to Contribute" guide. |
requirements.txt | Packages needed for the competition (pandas, statsmodels, fbprophet,
torch, etc.).
| `README.md` | Competition overview, dataset description, evaluation metric, timeline, and
contact info. I
| `CONTRIBUTING.md` | How to submit a model, code style, and leaderboard update
process. |
| `LICENSE` | MIT. |
**Why it works:**
```

- **Engaging** participants can see their rank on a live leaderboard.
- **Educational** baseline notebooks teach classic and modern forecasting methods.
- **Community-driven** open to contributions, fostering collaboration.

```
## 3. **Fiscal-Policy-Impact-Dashboard**
```

Goal: Build an interactive dashboard (Streamlit / Panel) that lets users explore how fiscal variables (tax rates, government spending, debt-to-GDP) affect macro outcomes (inflation, growth) across countries.

```
| Folder | Contents |
|-----|
| `data/` | Cleaned macro dataset (`macro_panel.csv`), metadata (`metadata.yaml`). |
| `src/` | `dashboard.py` (Streamlit app), `utils.py` (data loaders, helper functions). |
| `models/` | Pre-trained regression models (e.g., Random Forest) used for on-the-fly
predictions. |
| `assets/` | CSS/JS for custom styling, logo, favicon. |
| `docs/` | MkDocs site with "Data Dictionary", "Methodology", and "Interpretation Guide". |
| `tests/` | Unit tests for data loading and model inference. |
| `requirements.txt` | Streamlit, pandas, scikit-learn, plotly, etc. |
| `environment.yml` | Conda environment. |
| `README.md` | Project purpose, demo GIF, installation, usage instructions, and a
"Contact" section. |
| `CONTRIBUTING.md` | Guidelines for adding new visualizations or models. |
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** - professional documentation + interactive UI.
- **ML-focused** – underlying predictive models, but the UI emphasizes exploration.
- **Reproducible** – all data transformations are version-controlled.
## 4. in **AutoML-Fiscal-Predictor**
**Goal:** Demonstrate an end-to-end AutoML pipeline that automatically selects features,
preprocesses data, and fits the best regression model for predicting government revenue
```

(% GDP) from macro indicators.

```
| Folder | Contents |
|-----|
| `data/` | Raw and cleaned datasets (`raw.csv`, `clean.csv`). |
| `src/` | `pipeline.py` (scikit-learn `Pipeline` with `ColumnTransformer`), `automl.py` (calls
`tpot`, `auto-sklearn`, or `h2o`). |
| `configs/` | YAML files defining search spaces, cross-validation folds, and evaluation
metrics.
| `notebooks/` | `01_Data_Exploration.ipynb`, `02_AutoML_Comparison.ipynb`. |
| 'models/' | Best model artifacts ('best model.pkl', 'best model.onnx'). |
```

```
\'docs/\' | Sphinx docs describing the AutoML workflow, hyper-parameter choices, and
reproducibility. |
| `tests/` | Tests for pipeline integrity (`pytest`). |
| `scripts/` | CLI to run the pipeline (`run_automl.sh`). |
| 'requirements.txt' | 'pandas', 'scikit-learn', 'tpot', 'auto-sklearn', 'h2o', etc. |
| `environment.yml` | Conda environment. |
| `README.md` | Project overview, motivation, step-by-step guide, and performance
summary. |
| `CONTRIBUTING.md` | How to add new estimators or datasets. |
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** – clear documentation of the AutoML process, reproducible scripts.
- **ML-centric** - showcases multiple AutoML frameworks, model comparison, and
versioning.
- **Educational** – readers can see the trade-offs between different AutoML tools.
## 5. **Macro-ML-Case-Studies**
**Goal:** A collection of self-contained case studies (one per notebook) that walk through a
full data-science workflow on a specific macro-economic question (e.g., "Can we predict
sovereign default using macro-variables?").
| Folder | Contents |
|-----|
case studies/ One sub-folder per case ('sovereign default/', 'inflation forecast/',
`revenue prediction/`). |
| `data/` | Raw and cleaned data per case (`sovereign_default/raw.csv`,
`sovereign default/clean.csv`). |
| `src/` | Shared utility functions (`utils.py`). |
| `notebooks/` | One notebook per case (`sovereign_default.ipynb`,
`inflation forecast.ipynb`). |
| `models/` | Saved models for each case. |
| `docs/` | MkDocs site with a table of contents, methodology overview, and a "Lessons
Learned" page. |
| `requirements.txt` | Common dependencies plus case-specific ones. |
| `environment.yml` | Conda environment. |
| `README.md` | List of case studies, motivation, and how to run each notebook. |
| `CONTRIBUTING.md` | Guidelines for adding new case studies (data source, analysis
steps, write-up).
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** – each case study reads like a short research report with code.
- **ML-focused** – each notebook follows the full pipeline (EDA → feature engineering →
modeling \rightarrow evaluation).
- **Modular** – new case studies can be added without touching existing code.
```

```
## 6. **Fiscal-Policy-Literature-Review-with-Code**
**Goal:** Combine a literature review of fiscal policy and macro-economic modeling with
reproducible code snippets that implement the key methods discussed in the papers.
| Folder | Contents |
|-----|
| `literature/` | PDFs or links to seminal papers (`Kydland & Prescott 1977.pdf`,
`Romer_1990.pdf`). |
'src/' | Python modules implementing key estimators (e.g., 'cagan_model.py',
`new_keynesian_model.py`). |
| `notebooks/` | `literature_review.ipynb` – narrative with embedded code cells, plus
'replication.ipynb' for a selected paper. |
| `docs/` | MkDocs site with a "Methodology" section, bibliography (BibTeX), and a "Code
Appendix". |
| `requirements.txt` | Packages for econometrics (`statsmodels`, `linearmodels`, `arch`). |
| `environment.yml` | Conda environment. |
| `README.md` | Project description, list of covered papers, and instructions. |
| `CONTRIBUTING.md` | How to add new papers or improve code snippets. |
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** – scholarly tone in the narrative, formal code implementations.
- **ML-focused** - concrete implementations of econometric models that often underpin ML
pipelines.
- **Research-oriented** – bridges theory and practice, useful for students and researchers.
## 7. **Fiscal-Policy-Simulation-Engine**
**Goal:** Provide a lightweight simulation engine (Monte-Carlo or agent-based) that lets
users experiment with fiscal shocks (tax cuts, spending increases) and observe
macro-economic outcomes.
| Folder | Contents |
|-----|
| `engine/` | Core simulation code (`simulation.py`), parameter definitions (`params.yaml`). |
| `data/` | Baseline parameter sets, calibration data. |
| `notebooks/` | `simulation_demo.ipynb` (interactive sliders via ipywidgets),
`sensitivity_analysis.ipynb`. |
| `results/` | Sample simulation outputs, plots, and summary statistics. |
| 'docs/' | Documentation of model equations, calibration procedure, and interpretation
guide. |
| `tests/` | Unit tests for the simulation core. |
| `requirements.txt` | `numpy`, `pandas`, `matplotlib`, `seaborn`, `ipywidgets`. |
| `environment.yml` | Conda environment. |
```

```
| `README.md` | Project overview, simulation equations, usage guide, and example plots. |
| `CONTRIBUTING.md` | How to extend the model (add new variables, stochastic shocks). |
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** – technical description of the simulation model plus clear usage docs.
- **ML-adjacent** – simulation outputs can be fed into ML models for forecasting or policy
impact analysis.
- **Interactive** - ipywidgets make the repo feel dynamic and engaging.
## 8. **Fiscal-Policy-Risk-Analytics**
**Goal:** Build a risk-analytics framework that quantifies fiscal risk (e.g., probability of
debt-service default) using macro-economic variables and machine-learning classification.
| Folder | Contents |
|-----|
| `data/` | Historical fiscal and macro data (`fiscal_panel.csv`), default events
('default events.csv'). |
| `src/` | `features.py` (feature engineering), `models.py` (logistic regression, XGBoost,
LightGBM), `risk_metrics.py`. |
| `notebooks/` | `risk_modeling.ipynb` (EDA \rightarrow feature selection \rightarrow model training \rightarrow
calibration), `calibration_demo.ipynb`. |
| `models/` | Pickled classifiers, calibrated probability models. |
| `docs/` | Sphinx docs with sections: "Risk Metrics", "Model Evaluation", "Interpretation of
SHAP values". |
| `tests/` | Tests for feature pipelines and model evaluation. |
| `scripts/` | CLI to compute risk scores for a new country (`compute_risk.py`). |
| `requirements.txt` | `pandas`, `scikit-learn`, `xgboost`, `lightgbm`, `shap`, `plotly`. |
| `environment.yml` | Conda environment. |
| `README.md` | Project purpose, risk metrics (PD, LGD, expected loss), demo, and
installation. |
CONTRIBUTING.md` | Adding new risk indicators or alternative classifiers. |
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** – formal risk-analytics terminology, rigorous evaluation.
- **ML-focused** - classification, calibration, SHAP explanations, model interpretability.
- **Practical** – directly applicable to sovereign-risk departments or financial institutions.
## 9. **Fiscal-Policy-Visualization-Library**
**Goal:** A small library of reusable visualizations (e.g., stacked area charts of fiscal
aggregates, Sankey diagrams of tax-to-spending flows) that can be imported into any
Jupyter notebook or Streamlit app.
```

```
| Folder | Contents |
|-----|
'viz/' | Python module ('fiscal viz.py') with functions: 'stacked area()', 'sankey()',
`bubble_chart()`. |
| `examples/` | Jupyter notebooks showing each visualization on real macro data. |
| 'docs/' | Sphinx docs with API reference, gallery of plots, and styling guide. |
| `tests/` | Unit tests for each visualization function. |
| `requirements.txt` | `pandas`, `matplotlib`, `seaborn`, `plotly`, `networkx`. |
| `environment.yml` | Conda environment. |
| `README.md` | Library overview, installation, quick-start examples, and contribution guide.
| `CONTRIBUTING.md` | Guidelines for adding new chart types or styling themes. |
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** - professional documentation, clear API.
- **ML-adjacent** – visualizations help interpret ML results (e.g., feature importance plots).
- **Reusable** - other projects can import `fiscal viz` without pulling the whole analysis
pipeline.
## 10. 📚 **Fiscal-Policy-Data-Science-Curriculum**
**Goal:** A structured learning path (a "course") that teaches data-science and ML
techniques through the lens of fiscal policy analysis. Each module is a Jupyter notebook with
accompanying slides and quizzes.
| Folder | Contents |
|-----|
| `modules/` | `01_Intro_to_Macro_Data.ipynb`, `02_Exploratory_Analysis.ipynb`,
`03_Feature_Engineering.ipynb`, `04_Modeling_Techniques.ipynb`,
`05_Risk_Analytics.ipynb`. |
| `slides/` | PowerPoint/Google Slides PDFs for each module. |
| `quizzes/` | Auto-graded Jupyter notebooks (`quiz_01.ipynb`, etc.) or a separate `quizzes/`
folder with JSON files for an LMS. |
| `data/` | Small sample datasets for each module. |
| `solutions/` | Solutions notebooks (hidden from students). |
| `docs/` | MkDocs site with a syllabus, learning objectives, and a "Resources" page. |
| `requirements.txt` | Core data-science stack. |
| `environment.yml` | Conda environment. |
| `README.md` | Course description, target audience, prerequisites, and how to run the
notebooks.
| `CONTRIBUTING.md` | How to add new modules or improve guizzes. |
| `LICENSE` | MIT. |
**Why it works:**
- **Semi-formal** – academic-style notebooks + slide decks.
- **ML-focused** – each module builds ML skills while staying grounded in fiscal policy.
```

- **Educational** – ready-to-use for workshops, MOOCs, or self-study.

How to Choose the Best Fit

- 1. **Audience** If you target practitioners (analysts, policy makers), go with a **Dashboard** or **Risk-Analytics** repo.
- 2. **Community Interaction** A **Competition** or **AutoML** repo encourages contributions and leaderboard excitement.
- 3. **Educational Depth** **Case-Studies**, **Curriculum**, or **Literature-Review** are ideal for teaching or research groups.
- 4. **Reusability** **Toolkit**, **Visualization Library**, or **Simulation Engine** provide modular code that can be imported into other projects.

Quick "Starter" Checklist for Any Repo

- **README** clear purpose, one-sentence tagline, installation steps, demo GIF.
- **Code Style** enforce PEP 8, use `black`/`ruff` in CI.
- **Testing** `pytest` with >80 % coverage; CI badge (GitHub Actions).
- **Documentation** MkDocs/Sphinx site, auto

Data Cleaning

1. **Data Collection**

- Gathered macroeconomic indicators (2010-2024) from World Bank, IMF, and OECD.
- Variables include GDP, CPI, Unemployment, Tax Revenue, Government Expenditure, Interest Rate, Government Revenue (% GDP), etc.

2. **Data Import & Inspection**

- Loaded the raw CSV/Excel files into a Python `pandas` DataFrame.
- Checked for missing values, duplicate rows, and inconsistent data types.

3. **Handling Missing Data**

- Identified missing entries using `df.isnull().sum()`.
- For variables with <5 % missingness, imputed with the median of the column.
- For larger gaps, used linear interpolation or removed the series if the variable was not critical (e.g., short-term interest rate gaps).

4. **Outlier Detection & Treatment**

- Plotted box-plots and histograms for each variable.
- Applied the IQR rule (1.5 × IQR) to flag extreme values.
- Replaced outliers with the nearest non-outlier value (Winsorization) to avoid distorting the distribution.

5. **Standardization**

- Scaled numeric features to zero mean and unit variance (`StandardScaler` from `scikit-learn`).

- This step is essential for PCA and many regression algorithms to ensure variables contribute equally.
- 6. **Encoding Categorical Variables**
- Although most indicators were numeric, any categorical fields (e.g., region) were one-hot encoded.
- 7. **Final Clean Dataset**
 - Resulted in a tidy DataFrame with 14 core macro indicators, ready for analysis.

Exploratory Data Analysis (EDA)

- 1. **Descriptive Statistics**
 - Computed mean, median, standard deviation, min, max for each variable.
 - Summarized in a table (e.g., using `df.describe()`).
- 2. **Correlation Analysis**
 - Calculated Pearson correlation matrix ('df.corr()').
- Visualized with a heatmap (`sns.heatmap`) to spot strong relationships (e.g., Government Revenue ↔ Government Expenditure).
- 3. **Distribution Checks**
 - Plotted histograms and KDEs for key variables (GDP, Gov. Revenue, Inflation).
- Noted right-skewness in fiscal indicators, prompting log-transformation where appropriate.
- 4. **Pairwise Relationships**
 - Used `sns.pairplot` or a scatterplot matrix to explore bivariate relationships.
- Highlighted non-linear patterns (e.g., quadratic relationship between Inflation and Revenue).
- 5. **Time-Series Overview**
 - Plotted each indicator over time to detect trends, seasonality, or structural breaks.
 - Identified periods of economic shocks (e.g., 2008 crisis, COVID-19).
- 6. **Key Insights**
 - High positive correlation between fiscal variables (Revenue, Expenditure, Tax).
 - Moderate correlation between GDP/GNI and fiscal strength.
 - Negative correlation of Inflation and Interest Rates with fiscal stability.

Principal Component Analysis (PCA)

1. **Standardization**

- Re-standardized the cleaned dataset (mean = 0, variance = 1) to ensure comparability across scales.
- 2. **Covariance Matrix & Eigen-Decomposition**
 - Computed the covariance matrix of the standardized data.
 - Calculated eigenvalues and eigenvectors to identify principal components (PCs).
- 3. **Variance Explained**
 - Sorted eigenvalues in descending order.
- Determined the number of components needed to capture \sim 70 % of total variance (typically the first 3 PCs).
- 4. **Component Loadings**
 - Examined the eigenvectors to see which original variables load heavily on each PC.
 - Example loadings:
- **PC1** GDP, GNI, Government Expenditure, Tax, Revenue (captures overall fiscal/economic size).
 - **PC2** Inflation, GDP Growth, Current Account (monetary & external balance).
 - **PC3** CPI, Unemployment (price-pressure & labor slack).

5. **Scree Plot**

- Plotted eigenvalues to visualize the "elbow" where additional components add little explanatory power.

6. **Biplot / Loading Plot**

- Created a biplot showing both observations (countries) and variable loadings, illustrating how each variable contributes to the PCs.

7. **Interpretation**

- **PC1** (≈40 % variance) → Fiscal scale & macroeconomic size.
- **PC2** (≈18 % variance) → Monetary dynamics vs. fiscal stability.
- **PC3** (≈12 % variance) → Inflation-employment interactions.

K-Means Clustering

- 1. **Choosing the Number of Clusters (k)**
 - Applied the elbow method: plotted within-cluster sum of squares (inertia) vs. k.
 - Observed a clear bend at k = 2, indicating two distinct groups.

2. **Silhouette Analysis**

- Computed silhouette scores for k = 2-6.
- Silhouette ≈ 0.365 for k = 2, suggesting moderate separation but still meaningful clusters.

3. **Running K-Means**

- Initialized centroids using the *k-means++* algorithm for better convergence.

- Ran the algorithm (default 300 iterations, tolerance = 1e-4) on the standardized PCA scores (PC1-PC3).

4. **Cluster Interpretation**

- **Cluster 0 (Orange)**: High-revenue, stable, developed economies (positive PC1, negative PC3).
- **Cluster 1 (Blue)**: Low-revenue, inflation-sensitive, developing economies (negative PC1, negative PC2, negative PC3).

5. **Visualization**

- 3-D scatter plot of PC1, PC2, PC3 colored by cluster.
- Elbow plot and silhouette score chart to justify the choice of k.

Regression Modeling

1. **Target Variable**

- Government Revenue as a percentage of GDP (continuous outcome).

2. **Feature Selection**

- Used the three principal components (PC1, PC2, PC3) as predictors, capturing the majority of variance.
 - Also retained original variables (Tax, Expenditure, Inflation, etc.) for interpretability.

3. **Model Candidates**

- **Linear Regression** baseline, fully interpretable.
- **Ridge & Lasso** regularized linear models to handle multicollinearity.
- **ElasticNet** combination of L1/L2 penalties.
- **Tree-Based Models** Decision Tree, Random Forest, Gradient Boosting.
- **Support Vector Regression (SVR)** non-linear kernel option.
- **Principal Component Regression (PCR)** linear regression on PCs.

4. **Model Training & Evaluation**

- Split data into training (80 %) and testing (20 %) sets, preserving temporal order to avoid leakage.
 - Employed 5-fold cross-validation for hyper-parameter tuning (e.g., `GridSearchCV`).
 - Evaluated each model using:
 - **Mean Squared Error (MSE)**
 - **Mean Absolute Error (MAE)**
 - **R2 (coefficient of determination)**

5. **Best-Performing Model**

- **Random Forest** achieved the lowest test MSE (≈ 4.43) and highest R² (≈ 0.93).
- Gradient Boosting was a close second, offering slightly better interpretability via feature importance.

6. **Model Interpretation**

- **Random Forest Feature Importances**: Tax, Government Expenditure, and PC1 (fiscal scale) were top drivers.
 - **Linear Model Coefficients** (on PCs):
 - `Revenuê = 24.66 + 2.52·PC1 2.19·PC2 + 2.90·PC3`
 - Explains \sim 66 % of variance (R² = 0.66).
- **Non-linear Models**: Capture complex interactions (e.g., diminishing returns of fiscal size when inflation is high).

7. **Residual Analysis**

- Plotted residuals vs. fitted values to check homoscedasticity.
- Verified normality of residuals with Q-Q plots.

Summary of Steps

These steps collectively enable a comprehensive understanding of how macroeconomic factors interact, how countries can be grouped based on fiscal and monetary characteristics, and how to predict government revenue with high accuracy.