

# Methodology

This section describes the systematic process by which the AI-driven budget optimization model is developed and evaluated. We proceed in five stages: (1) Data Collection & Preprocessing; (2) Regression Analysis for Lever Identification; (3) Markov Decision Process (MDP) Formulation; (4) Policy Computation; and (5) Monte Carlo Stress-Testing.

## 1. Data Collection & Preprocessing

- **Scope & Sources:** Annual budget records over the past 5–10 years, including Approved Budget  $B_t^{app}$ , Released Budget  $B_t^{rel}$ , and Actual Expenditure  $B_t^{act}$  for each sector.
- **Macroeconomic indicators:** GDP ( $G_t$ ), inflation rate ( $\pi_t$ ), and sector-specific performance measures (e.g., enrolment, health outcomes).
- **Cleaning & Imputation:** Identify missing entries in  $B_t^{rel}$  or  $B_t^{act}$ . Impute using a time-series-aware approach (e.g., linear interpolation or Kalman filter). Flag imputed values and record an indicator variable  $\delta_t^{imp} \in \{0,1\}$ .
- **Feature Engineering:** Compute Derived Metrics and normalize variables:

```
E_t = B_t^act / B_t^app
E^GDP_t = B_t^act / G_t
ΔG_t = (G_t - G_{t-1}) / G_{t-1}
Normalize variables to zero mean and unit variance for regression stability.
```

## 2. Regression Analysis for Lever Identification

**Model Specification:** We model the key outcome  $Y_t = E_t^{GDP}$  as a function of candidate drivers  $X_t = [\Delta G_t, E_{t-1}, \delta_t^{imp}, \dots]$ .

$$Y_t = \beta_0 + \sum_i \beta_i X_{t,i} + \epsilon_t, \quad \epsilon_t \sim N(0, \sigma^2)$$

**Estimation & Selection:** Fit an Elastic Net regression (L1 + L2 penalties):

```
min_β (1/T) Σ_{t=1}^T (Y_t - β^T X_t)^2 + λ_1 ||β||_1 + λ_2 ||β||_2^2
Select significant coefficients {β_i ≠ 0} as decision levers a_i.
```

## 3. Markov Decision Process Formulation

Define the tuple  $\langle S, A, T, R, \gamma \rangle$ :

- **States (S):** Discrete fiscal states {Under█Spend, On█Track, Over█Commit}.
- **Actions (A):** Budget adjustments from regression levers (e.g., Increase/Decrease by 5%).

- **Transition (T):** Use regression to estimate  $\mathbf{E}_{t+1}$  and map to next state,  $T(s'|s,a)$ .
- **Reward (R):**  $R(s,a) = w_1 E_{t+1} - w_2 \mathbf{1}\{E_{t+1} < \alpha\}$ .
- **Discount ( $\gamma$ ):** Discount factor (e.g.,  $\gamma=0.95$ ).

## 4. Policy Computation

Solve for the optimal policy  $\pi^*: S \rightarrow A$  via Value Iteration:

```
Initialize  $V_0(s) = 0$  for all  $s \in S$ 
Repeat until convergence:
     $V_{k+1}(s) = \max_a [ R(s,a) + \gamma \sum_{s'} T(s'|s,a) V_k(s') ]$ 
 $\pi^*(s) = \operatorname{argmax}_a [ R(s,a) + \gamma \sum_{s'} T(s'|s,a) V^*(s') ]$ 
Convergence:  $||V_{k+1} - V_k||_{\infty} < \epsilon$ .
```

## 5. Monte Carlo Stress-Testing

**Uncertainty Characterization:** Fit distributions to regression residuals  $\varepsilon_t$ .

**Simulation Protocol:** For  $N$  runs (e.g., 10,000):

- Sample random shocks  $\varepsilon_t$  each period.
- Begin from state  $s_0$  and apply  $\pi^*(s)$ .
- Propagate via  $T(s'|s,a)$  with sampled shocks; record  $\{E_t^{\text{GDP}}\}$ .

**Outcome Analysis:** Compute mean execution, Value-at-Risk (VaR), expected shortfall; identify policy failure scenarios ( $E_t < \alpha$ ).