

Cairo University
Faculty of Graduate Studies for
Statistical Research



Panel Data Methodologies

With

Application for Macroeconometrics

(Inflation Forecasting)

By: Mohamed Sayed Abd AL-mgyd

Statistician and Data Scientist

Supervised By

Dr. Amal Mohamed Abd-Elfatah

Assistant Professor of Statistics

Department of Applied Statistics and Econometrics

Faculty of Graduate Studies for Statistical Research (FSSR)

Table of Contents

Acknowledgement:	3
Abstract	4
1. Introduction	5
1.1. Study Objectives	6
1.2. Structure of the Study	6
2. Literature Review	7
2.1. Overview	7
2.2. Applications	10
3. Estimation Methods for Panel Data Models	12
3.1. Pooled Ordinary Least Squares (POLS)	12
3.2. Fixed Effects (FE)	14
3.3. Random Effects (RE)	18
3.4. Generalized Method of Moments (GMM)	25
4. Applied	30
4.1. Dataset Overview	30
4.2. Dataset Preprocessing	34
4.3. Models	45
Pooled Ordinary Least Squares (Pooled OLS) Model:	45
Fixed Effects (FE) Model:	46
Two-Step Difference Generalized Method of Moments (Difference GMM, Arellano–Bond):	53
4.4. Conclusion: Model Comparison and Final Selection	57
5. Recommendations and Future Work	58
Appendix (Codes)	59
References	60

List of Tables

Table 1) Classification and Description of Variables:-	31
Table 2) Number of Missing Values by Variable.....	34
Table 3) Descriptive Statistics	39
Table 4) Variance Inflation Factor (VIF).....	42
Table 5) Outlier Detection by Country	43
Table 6) (Pooled OLS) Summary:.....	45
Table 7) (Pooled OLS) Coefficients:.....	45
Table 8) (FE) Summary:	46
Table 9 (FE) Coefficients:	46
Table 10) (RE) Summary	47
Table 11 (RE) Coefficients:	47
Table 12) Corrected Coefficient Significance – RE Model with (SCC) Standard Errors: ..	51
Table 13) (GMM) Coefficients.....	54

List of Figures

Figure 1) Literature Evolution Tree.....	9
Figure 2) Number of Missing Values by Country	35
Figure 3) Apply MICE-BR to missing values for Egypt.....	38
Figure 4) Correlation	41

Acknowledgement:

Thanks, and appreciation to my parents and my supervisor Prof. Amal Mohamed Abdelfatah.

Abstract

This study investigates the effectiveness of panel data methodologies in Macroeconometrics modeling, with a specific focus on inflation forecasting. Utilizing a balanced panel dataset of 77 countries from 1980 to 2024, the analysis evaluates the empirical performance of several estimators: Pooled Ordinary Least Squares (Pooled OLS), Fixed Effects (FE), Random Effects (RE), and the Two-Step Difference Generalized Method of Moments (Difference GMM) developed by Arellano–Bond.

The dependent variable is average consumer price inflation (PCPIPCH), modeled against a set of macroeconomic indicators including fiscal balances, public debt ratios, trade volumes, investment shares, and PPP-adjusted income—sourced primarily from the IMF’s World Economic Outlook and the World Bank.

The methodology includes diagnostic testing for heteroskedasticity, serial correlation, cross-sectional dependence, and unit roots, followed by model selection using the Wald and Hausman tests. To ensure robust inference, standard errors are corrected using the Driscoll–Kraay estimator, and the Arellano–Bond GMM framework is applied to address endogeneity and inflation persistence using internal instruments for the lagged dependent variable.

Empirical results indicate that standard static models (e.g., Pooled OLS, FE) are sensitive to violations of econometric assumptions and deliver less reliable estimates under macro-panel conditions. In contrast, the Two-Step Difference GMM model consistently outperforms these alternatives, passing all diagnostic checks—including the Sargan test for instrument validity and the AR(2) test for residual autocorrelation—while producing efficient and theoretically consistent estimates.

This study highlights the importance of dynamic panel specifications in capturing structural inflation behavior and offers methodological guidance for researchers and policy insights for institutions involved in inflation targeting. By identifying the most suitable model, it contributes to the applied macroeconometrics literature and supports evidence-based policymaking in diverse economic contexts.

Keywords: Panel Data, Macroeconometrics, Inflation Forecasting, Fixed Effects, Random Effects, Pooled OLS, Endogeneity, Dynamic GMM, Arellano–Bond.

1. Introduction

Over the past several decades, the analysis of panel data has emerged as a fundamental approach in empirical economics, revolutionizing the way economists examine complex phenomena involving multiple entities observed over numerous time periods (Baltagi, 2008). Panel data, also known as longitudinal or cross-sectional time-series data, combines temporal depth with cross-sectional breadth, offering a powerful framework that enhances the precision and richness of econometric inference (Hsiao, 2003).

The evolution of panel methodologies began in the mid-20th century when researchers recognized the limitations of using only cross-sectional or time-series data. Early models treated observations independently, often overlooking unobserved factors and persistence in economic processes. Foundational work by Verbeek and Nijman (1992) and Moulton (1990) introduced techniques that model individual-specific effects and temporal dependencies. Arellano and Bond (1991) further advanced the field with dynamic panel estimators addressing endogeneity, solidifying panel analysis in modern econometrics (Arellano & Bond, 1991).

Panel data methods offer several crucial advantages over purely cross-sectional or time-series analyses. By including individual-specific parameters, either as fixed or random effects, panel models account for unobservable characteristics (e.g., institutional quality, cultural factors) that remain constant over time but vary across entities (Baltagi, 2008).

Panel methods offer several advantages. They incorporate individual-specific effects—fixed or random—that account for time-invariant unobservables like institutional quality (Baltagi, 2008). Pooling over time increases sample size and efficiency, reducing variance inflation (Wooldridge, 2010). Dynamic panels with lagged variables capture persistence in macro indicators, such as GDP growth trends (Arellano & Bond, 1991).

Depending on the data characteristics and research questions, economists can choose between pooled OLS, fixed effects, random effects, or more advanced techniques like Generalized Method of Moments (GMM). Panel models have shown great value in macroeconomic studies: Barro and Sala-i-Martin (2004) explored growth determinants across countries, Fischer (1993) showed how inflation inertia varies by regime, and Nickell (1997) applied fixed effects to study labor market rigidities across the OECD countries.

The policy implications are significant. By capturing both heterogeneity and dynamics, panel analysis helps design targeted policies—such as identifying effective fiscal stimuli or evaluating how interest rate changes affect different economies.

1.1. Study Objectives

This *project's primary objectives* are to:

1. Evaluate pooled OLS, fixed effects, random effects, and dynamic panel (Generalized Method of Moments) estimators in terms of consistency, efficiency, and applicability to macroeconomic data.
2. Utilize a panel dataset of macroeconomic indicators from multiple countries over two decades to demonstrate each model's performance, including diagnostic tests (Hausman test, Arellano-Bond test, Sargan test).
3. Develop recommendations for selecting appropriate panel methodologies based on data properties (e.g., N vs. T dimensions, presence of serial correlation, endogeneity risks).

By achieving these aims, the study will contribute both to econometric methodology and to evidence-based macroeconomic policymaking.

1.2. Structure of the Study

The thesis is organized as follows:

1. **Literature Review** – Synthesizes theoretical developments and empirical findings in panel data econometrics.
2. **Models** – Details the statistical foundations, estimation procedures, and diagnostic tests for each panel model.
3. **Empirical Analysis** – Applies the models to macroeconomic data, presents results, and interprets findings.
4. **Conclusion and Policy Implications** – Summarizes key insights, discusses limitations, and offers policy recommendations.

2. Literature Review

Panel data econometrics has undergone significant evolution over the past several decades. Early theoretical foundations emerged to address the limitations of cross-sectional and time-series models, paving the way for comprehensive methods that control unobserved heterogeneity, serial correlation, and endogeneity.

2.1. Overview

Hsiao (2003) provided the first systematic framework for panel analysis, introducing the within-transformation for fixed effects estimation and discussing the challenges of serial correlation and missing observations. Baltagi (2008) formalized the asymptotic properties of fixed effects (FE) and random effects (RE) estimators, deriving the random effects generalized least squares (GLS) formula and comparing bias-variance trade-offs. Wooldridge (2010) enriched these foundations by integrating diagnostic tools—such as the Hausman specification test and cluster-robust standard errors—and by addressing cross-sectional dependence using Pesaran’s CD test.

Mundlak (1978) demonstrated that including unit means of regressors captures correlation between individual effects and explanatory variables, underpinning the RE model intuition. Hausman (1978) introduced the Hausman test to choose between FE and RE by detecting inconsistent RE assumptions.

Arellano and Bond (1991) revolutionized dynamic panel analysis with the Difference GMM estimator, which first-differences to remove fixed effects and uses lagged levels as instruments to address endogeneity. Blundell and Bond (1998) extended this to System GMM, combining level and difference equations to improve efficiency with persistent data.

Pesaran (2004) introduced the CD test for cross-sectional dependence, guiding the use of factor-augmented regressions. Moon and Weidner (2015) developed interactive fixed effects models that estimate unobserved common factors varying over time, refined by Chudik and Pesaran (2018) to allow multiple latent factors via generalized least squares corrections.

Ahn, Lee, and Schmidt (2019) proposed jackknife bias reduction for GMM in panels with highly persistent dynamics. Bai, Liao, and Shi (2020) integrated factor estimation into system GMM to jointly address endogeneity and cross-dependence. Aghion et al. (2021) introduced bias correction for network spillovers and measurement errors using higher-order instruments, while Huang and

Pesaran (2022) incorporated spatial weight matrices into interactive effects. Sun and Kim (2023) applied LASSO regularization within GMM to select optimal instruments, and Zhang and Lee (2024) leveraged machine learning (random forests) to generate non-linear instruments for panels with structural breaks.

Panel data methodologies rely on rigorous diagnostic checks to ensure estimator validity and robustness. The Hausman Test (Hausman, 1978) serves as a pivotal specification test, comparing fixed effects (FE) and random effects (RE) estimates. By testing the null hypothesis that individual effects are uncorrelated with regressors, a significant Hausman statistic indicates that RE assumptions fail, favoring the FE model for consistent coefficient estimates in the presence of endogeneity.

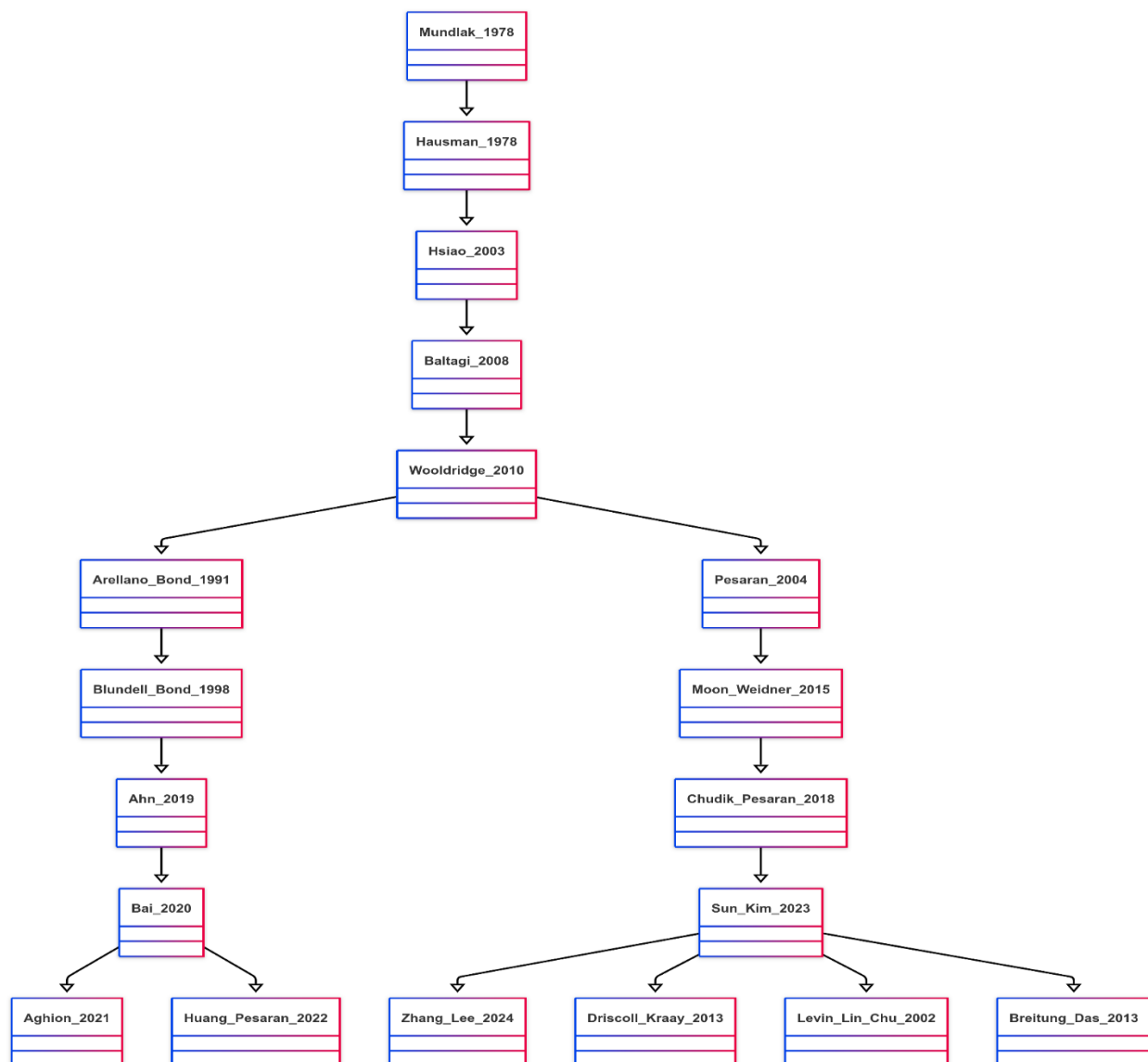
Serial Correlation Tests are crucial in dynamic panel contexts to validate instrument use. Arellano and Bond (1991) introduced tests for first order and second-order autocorrelation in differenced residuals. Since first-difference mechanically induces negative first-order autocorrelation, researchers focus on the absence of second-order autocorrelation (AR(2)). A rejection of the null of no AR(2) suggests instrument invalidity, undermining GMM estimates.

Cross-Sectional Dependence undermines standard error calculations if unaddressed. Pesaran's CD test (Pesaran, 2004) computes the average pairwise correlation of residuals across panel units under the null of cross-sectional independence. Significant CD statistics alert to latent common factors, prompting the use of factor-augmented regressions or panel estimators that incorporate common correlated effects. Driscoll and Kraay (2013) further developed robust covariance matrix estimators that remain consistent under general forms of spatial and serial dependence, offering an alternative when factor structure is difficult to specify.

Unit Root and Cointegration Tests guide model specifications by diagnosing non-stationarity. Levin, Lin, and Chu (2002) proposed panel unit root tests under a common autoregressive parameter, controlling for entity-specific deterministic trends and serial correlation. When series exhibit unit roots, differencing or incorporating error-correction terms becomes necessary. Building on this, Breitung and Das (2013) introduced panel cointegration methods that test long-run equilibrium relationships among non-stationary variables by extending Pedroni's group mean statistics to account for heterogeneity and cross-dependence.

Finally, Overidentification Tests such as the Sargan and Hansen J-tests evaluate the joint validity of instruments in GMM estimation. Under the null that instruments are orthogonal to the error term, a high p-value confirms instrument exogeneity. However, overfitting with too many instruments can weaken test power, requiring careful instrument selection and potential use of instrument reduction techniques such as the collapsed instrument matrix or LASSO-based selection (Sun & Kim, 2023).

Figure 1) Literature Evolution Tree



Source: <https://www.mermaidchart.com/app/projects/1c581e1e-0fe6-4ce3-b041-a72174f9ca0b/diagrams/dc07b120-1935-46e0-8043-60440a4b7f35/>

2.2. Applications

Barro and Sala-i-Martin (2004) conducted one of the earliest panel studies on GDP convergence by applying both fixed effects and dynamic GMM estimators to a large cross-country panel. They demonstrated that poorer countries' growth rates converge more slowly toward richer countries once country-specific unobservable—captured via fixed effects—are controlled for. Their dynamic GMM implementation, using lagged GDP levels as instruments, provided robust evidence against simple pooled OLS conclusions, highlighting the persistence of growth dynamics over time.

Fischer (1993) applied difference GMM to panel inflation data from a sample of industrial economies, uncovering significant inflation inertia and the differential impact of monetary regimes. By first-difference and using lagged inflation rates as instruments, he isolated genuine serial correlation in the inflation process. Building on this, Breitung and Das (2013) extended the analysis to emerging markets by employing panel cointegration techniques. They showed that while short-run inflation-output trade-offs vary across countries, long-run relationships adhere to a stable Phillips curve, validated through Pedroni-style group mean statistics adapted for cross-dependence.

Nickell (1997) utilized fixed effects models to examine unemployment dynamics within OECD countries, focusing on the role of labor market regulations. His within-transformation approach removed time-invariant country effects, revealing that stricter employment protection and higher unemployment benefits substantially increase equilibrium unemployment. Later, Ciccone (2015) used system GMM to assess how the 2008 financial crisis affected unemployment persistence in advanced economies. Ciccone's study leveraged internal instruments to control for endogeneity of policy responses and demonstrated that crisis-induced policy shifts had long-lasting labor market effects.

Becker, Fetzer, and Novy (2010) introduced interactive fixed effects to panel studies of fiscal policy, capturing unobserved global shocks while estimating heterogeneous fiscal multipliers across countries. Their approach combined factor-augmented regressions with time-varying loadings, revealing that fiscal stimulus efficacy depends critically on country-specific characteristics and global business cycle phases. Imbs and Wacziarg (2018) applied factor-augmented dynamic panels to study globalization's impact on productivity, showing that latent

common factors—representing global integration forces—significantly drive productivity convergence among countries.

Levine, Loayza, and Beck (2000) pioneered the use of system GMM to study the relationship between financial development and economic growth, instrumentalizing financial depth indicators with their own lags. They provided early evidence that deeper banking systems promote growth after addressing simultaneity and omitted variables. More recently, Aghion et al. (2021) expanded this line of research by incorporating network spillovers and measurement error corrections in R&D panels. Their bias-corrected GMM framework, using higher-order spatial and temporal lags as instruments, offered more precise estimates of R&D's productivity spillovers across OECD countries.

3. Estimation Methods for Panel Data Models

This chapter presents the classical panel data models along with their respective assumptions and estimation methods, focusing on Pooled Ordinary Least Squares (POLS), Fixed Effects (FE), and Random Effects (RE) models. These models are essential for handling panel data and provide different approaches to managing individual-specific effects and time-related variations. In addition, the chapter introduces Dynamic Panel Data models, which incorporate lagged dependent variables to capture temporal dynamics, and discusses estimation techniques such as the Generalized Method of Moments (GMM), commonly used to address endogeneity in dynamic settings.

3.1. Pooled Ordinary Least Squares (POLS)

Panel data, also called longitudinal data, includes observations on multiple entities (e.g., individuals, firms) over multiple time periods. This structure allows researchers to examine variations both across entities and over time. In this context, Pooled Ordinary Least Squares (POLS) are a fundamental method for analyzing panel data by combining observations across both dimensions into a single, extensive cross-sectional dataset (Wooldridge, 2010).

The panel data model for POLS is specified as follows (Baltagi, 2008):

$$Y_{it} = \beta_0 + \beta_1 X_{it} + \epsilon_{it}$$

Where Y_{it} represents the dependent variable for entity i at time t , X_{it} denotes the independent variable for entity i at time t , β_0 is the intercept, β_1 is the slope coefficient, and ϵ_{it} is the error term.

The key assumptions of the POLS model include linearity, independence, homoscedasticity, and no autocorrelation of the error terms. Specifically, it is assumed that the relationship between Y_{it} and X_{it} is linear, the error terms are independently distributed and homoscedastic, and there is no autocorrelation across time periods.

To validate these assumptions, several diagnostic tests can be employed. The F-test is used to assess the joint significance of the regression coefficients. The Breusch-Pagan test helps detect heteroscedasticity, while the Durbin-Watson test evaluates autocorrelation in the residuals.

We minimize the residual sum of squares (RSS) to estimate the POLS model. The RSS is calculated as (Baltagi, 2008):

$$RSS = \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \beta_0 - \beta_1 X_{it})^2$$

The POLS estimator for β_0 and β_1 is derived by setting the partial derivatives of the RSS to zero with respect to these parameters. This yields the following equations:

$$\begin{aligned} \frac{\partial RSS}{\partial \beta_0} &= -2 \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \beta_0 - \beta_1 X_{it}) = 0 \\ \frac{\partial RSS}{\partial \beta_1} &= -2 \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \beta_0 - \beta_1 X_{it}) X_{it} = 0 \end{aligned}$$

Solving these equations provides the POLS estimates:

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^N \sum_{t=1}^T (X_{it} - \bar{X})(Y_{it} - \bar{Y})}{\sum_{i=1}^N \sum_{t=1}^T (X_{it} - \bar{X})^2} \\ \hat{\beta}_0 &= \bar{Y} - \hat{\beta}_1 \bar{X} \end{aligned}$$

Where \bar{X} and \bar{Y} represent the overall means of X_{it} and Y_{it} , respectively.

While POLS is a straightforward method and easy to implement, it may not always be optimal for panel data analysis. POLS assumes homoscedastic errors and ignores individual-specific effects, leading to inefficient or biased estimates if these assumptions are unmet. Alternative methods, such as Fixed Effects (FE) and Random Effects (RE) models, are designed to address these issues by accounting for unobserved heterogeneity and correlation within entities. Pooled Ordinary Least Squares remains a foundational technique in panel data analysis due to its simplicity and ease of application. However, it is crucial to verify its underlying assumptions and

consider more advanced models when necessary to ensure the robustness and reliability of the results (Greene, 2012).

3.2. Fixed Effects (FE)

Fixed Effects (FE) models are essential in panel data analysis for controlling unobserved individual-specific heterogeneity that remains constant over time. These models are particularly advantageous when analyzing datasets with repeated observations on the same entities—individuals, firms, or countries—where the individual-specific effects correlate with the explanatory variables (Wooldridge, 2010). If not accounted for, this correlation can lead to biased estimates in simpler models such as Pooled OLS (Ordinary Least Squares). The FE model addresses this issue by introducing entity-specific intercepts that capture these unobserved characteristics, thus controlling for any time-invariant attributes (Baltagi, 2008).

The Fixed Effects model can be formally represented as:

$$Y_{it} = \alpha_i + \beta_1 X_{it} + \epsilon_{it}$$

Where Y_{it} denotes the dependent variable for entity i at time t , and X_{it} represents the independent variable for the same entity and time period. The term α_i signifies the individual-specific fixed effect, which captures the unobserved characteristics that are constant over time but vary across entities. β_1 is the coefficient of the independent variable X_{it} , and ϵ_{it} represents the idiosyncratic error term. The inclusion of α_i in the model allows for the control of all time-invariant characteristics of the entities, which might otherwise confound the relationship between the dependent and independent variables (Greene, 2012).

The Fixed Effects model operates under several key assumptions. First, it assumes that the individual-specific effects α_i are correlated with the regressors X_{it} . This correlation is critical because it justifies using Fixed Effects to control for potential bias in parameter estimates. Second, it assumes that α_i are constant over time for each entity, thus capturing all time-invariant factors. Lastly, the model assumes that the error term ϵ_{it} is independently and identically distributed with zero mean and constant variance and is uncorrelated with the regressors X_{it} within each entity.

Several diagnostic tests are employed to ensure the Fixed Effects model's robustness. The F-test for Fixed Effects assesses whether the fixed effects α_i are jointly significant. If the F-test indicates that the fixed effects are significant, it supports using the Fixed Effects model over the Pooled OLS model. The Hausman test is also used to compare the Fixed Effects model with the Random Effects model. The null hypothesis of the Hausman test is that the Random Effects model is appropriate, which assumes that the entity-specific effects are uncorrelated with the regressors. If the Hausman test rejects this null hypothesis, it indicates that the Fixed Effects model is more suitable (Wooldridge, 2010).

Two prominent methods for estimating the Fixed Effects model are the Within Transformation (de-meaning) and the Least Squares Dummy Variable (LSDV) approach.

Within Transformation (De-meaning Approach)

The Within Transformation method involves removing the individual-specific effects by subtracting the entity-specific means from the observations. This transformation is applied as follows (Baltagi, 2021):

$$\begin{aligned}\tilde{Y}_{it} &= Y_{it} - \bar{Y}_i \\ \tilde{X}_{it} &= X_{it} - \bar{X}_i\end{aligned}$$

Here, \bar{Y}_i and \bar{X}_i are the averages of Y_{it} and X_{it} for entity i over time t . By substituting these transformed variables into the original model, we obtain:

$$\tilde{Y}_{it} = \beta_1 \tilde{X}_{it} + \tilde{\epsilon}_{it}$$

Where $\tilde{\epsilon}_{it}$ represents the transformed error term. This transformation effectively eliminates the individual-specific effect α_i , which was present in the original model as:

$$Y_{it} = \beta_1 X_{it} + \alpha_i + \epsilon_{it}$$

By removing α_i through the de-meaning process, we obtain a model that only contains the within-entity variation of the dependent and independent variables. This allows us to estimate β_1 while controlling for the entity-specific effects. The OLS estimator for β_1 is given by:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^N \sum_{t=1}^T \tilde{X}_{it} \tilde{Y}_{it}}{\sum_{i=1}^N \sum_{t=1}^T \tilde{X}_{it}^2}$$

This estimator is derived from minimizing the residual sum of squares in the transformed regression model. The consistency and unbiasedness of $\hat{\beta}_1$ are guaranteed under the assumption that the regressors are correlated with the error term ϵ_{it} , which means that the within-entity variation in X_{it} and Y_{it} provides valid information about the relationship between them despite entity-specific effects. The Within Transformation is particularly advantageous as it controls for unobserved heterogeneity that is constant over time within entities. This method removes the bias associated with omitted variable bias that could arise from these time-invariant individual-specific effects. However, it should be noted that this approach only uses within-entity variation and does not account for between-entity differences. As a result, any variation that is constant across entities or that does not change over time within entities is not used in the estimation of β_1 (Greene, 2012).

Least Squares Dummy Variable (LSDV) Approach

The Least Squares Dummy Variable (LSDV) approach is widely used for estimating Fixed Effects models in panel data analysis. This method explicitly includes dummy variables for each entity, capturing the individual-specific effects directly within the regression model. The LSDV approach can be advantageous when dealing with panel data where individual-specific heterogeneity is a significant concern. The model is specified as (Greene, 2012):

$$Y_{it} = \beta_1 X_{it} + \sum_{i=1}^N \delta_i D_i + \epsilon_{it}$$

Where D_i is a dummy variable for entity i , and δ_i denotes the fixed effect associated with each entity. The term $\sum_{i=1}^N \delta_i D_i$ represents the individual-specific effects α_i in the model. This method estimates β_1 while directly accounting for the α_i by including a separate dummy variable for each entity. The LSDV approach involves estimating a separate parameter for each entity's fixed effect, which can lead to computational challenges, mainly when the number of entities N is large. Each dummy variable D_i introduces an additional parameter, which increases the complexity of the estimation. Consequently, while the LSDV method provides a straightforward way to estimate

Fixed Effects, it may be less efficient regarding computational resources than other methods, especially in large datasets. To estimate β_1 using the LSDV approach, the minimization problem is formulated as follows:

$$\min_{\beta_1} \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \beta_1 X_{it} - \delta_i)^2$$

Here, the objective is to minimize the residual sum of squares (RSS), which is the sum of the squared differences between the observed and predicted values of Y_{it} . The residuals include the fixed effects captured by the dummy variables. Taking the derivative of the RSS with respect to β_1 . Moreover, by setting it to zero, we obtain the first-order condition for minimization:

$$\frac{\partial}{\partial \beta_1} \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \beta_1 X_{it} - \delta_i)^2 = -2 \sum_{i=1}^N \sum_{t=1}^T X_{it} (Y_{it} - \beta_1 X_{it} - \delta_i) = 0$$

Solving for β_1 , we derive the estimator:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^N \sum_{t=1}^T X_{it} (Y_{it} - \hat{\delta}_i)}{\sum_{i=1}^N \sum_{t=1}^T X_{it}^2}$$

where $\hat{\delta}_i$ represents the estimated fixed effects for each entity. This estimator adjusts for the individual-specific effects by focusing on the within-entity variations. The LSDV approach provides a straightforward and interpretable method for estimating Fixed Effects models. Including dummy variables for each entity directly models the individual-specific effects, which helps control unobserved heterogeneity. However, the computational burden associated with many dummy variables can be significant, and the method might be less efficient than other techniques, such as the Within Transformation (demeaning approach) (Wooldridge, 2010; Greene, 2012).

The LSDV approach is particularly valuable when the goal is to account for fixed effects explicitly and when the number of entities is manageable. Alternative methods may be considered for large datasets or when computational efficiency is a concern.

3.3. Random Effects (RE)

The Random Effects (RE) model is a prevalent choice due to its ability to account for individual-specific variability that is assumed to be uncorrelated with the explanatory variables. Here, comprehensively explores the Random Effects model, including its formulation, estimation methods, and critical mathematical derivations. The Random Effects model assumes that individual-specific effects are randomly distributed and uncorrelated with the explanatory variables included in the model. This contrasts with the Fixed Effects model, where the individual-specific effects are treated as parameters to be estimated. Mathematically, the Random Effects model can be represented as:

$$y_{it} = \beta_0 + \beta_1 x_{it} + u_i + \epsilon_{it}$$

Where y_{it} is the dependent variable for individual i at time t , β_0 is the intercept, β_1 is the coefficient of the explanatory variable x_{it} , u_i represents the individual-specific effect, and ϵ_{it} is the idiosyncratic error term.

The individual-specific effect u_i is assumed to be randomly distributed with $u_i \sim N(0, \sigma_u^2)$, and the idiosyncratic error ϵ_{it} is assumed to be independently and identically distributed with $\epsilon_{it} \sim N(0, \sigma_\epsilon^2)$. Importantly, u_i and ϵ_{it} are assumed to be uncorrelated. The Random Effects model is beneficial when the variation across individuals is assumed to be random and unrelated to the explanatory variables. However, one must verify the appropriateness of the Random Effects model using statistical tests such as the Hausman test, which assesses whether the Random Effects assumption of no correlation between individual-specific effects and regressors is valid. A significant result from the Hausman test suggests that a Fixed Effects model may be more appropriate. The Random Effects model assumes that the random effects are uncorrelated with the regressors, which might not always be valid. When this assumption is violated, the estimates from the Random Effects model may be biased and inconsistent, making it crucial to assess the validity of this model assumption in empirical analyses.

The Random Effects model can be estimated using several techniques. The most common methods include:

Generalized Least Squares (GLS)

Generalized Least Squares (GLS) are used to handle heteroscedasticity and correlation in the error terms of panel data models, particularly in the context of Random Effects (RE) models. The primary advantage of GLS is its ability to provide efficient estimates when the assumptions of classical Ordinary Least Squares (OLS) are violated, specifically when errors are not independently and identically distributed (i.i.d.) or exhibit correlation across entities. In a Random Effects model, the observed data can be represented as (Greene, 2012):

$$Y_{it} = X_{it}\beta + u_i + \epsilon_{it}$$

Where Y_{it} is the dependent variable for entity i at time t , X_{it} is a vector of regressors, β is the vector of coefficients to be estimated, u_i represents the random effect specific to entity i , and ϵ_{it} is the idiosyncratic error term. The random effects u_i are assumed to be independently distributed with variance σ_u^2 , and the idiosyncratic errors ϵ_{it} are assumed to be independently distributed with variance σ_ϵ^2 . The covariance matrix of the error terms, Ω , reflects the combined variance from both sources:

$$\Omega = \sigma_u^2 J + \sigma_\epsilon^2 I$$

Where J is an $N \times N$ matrix with all entries equal to 1, representing the covariance due to random effects, and I is the identity matrix of dimension $T \times T$, representing the idiosyncratic variance.

To apply GLS, we first need to transform the model to remove the correlation introduced by the random effects. This is done by applying a matrix transformation that standardizes the error terms. Specifically, we pre-multiply the model by $\Omega^{-1/2}$, where $\Omega^{-1/2}$ is the matrix square root of the inverse of Ω . The matrix $\Omega^{-1/2}$ can be decomposed as follows (Baltagi, 2008):

$$\Omega^{-1/2} = \frac{1}{\sqrt{\sigma_\epsilon^2}} \left(I - \frac{\sigma_u^2}{\sigma_\epsilon^2 + \sigma_u^2} J \right)$$

The transformed model then becomes:

$$\Omega^{-1/2} Y_{it} = \Omega^{-1/2} X_{it} \beta + \Omega^{-1/2} \epsilon_{it}$$

where $\Omega^{-1/2}\epsilon_{it}$ is now homoscedastic and uncorrelated, simplifying the error structure to meet the assumptions of OLS. After transformation, the GLS estimator can be derived by applying OLS to the transformed model. The transformed regression equation is:

$$\tilde{Y}_{it} = \tilde{X}_{it}\beta + \tilde{\epsilon}_{it}$$

Where:

- $\tilde{Y}_{it} = \Omega^{-1/2}Y_{it}$
- $\tilde{X}_{it} = \Omega^{-1/2}X_{it}$,
- $\tilde{\epsilon}_{it} = \Omega^{-1/2}\epsilon_{it}$.

The GLS estimator for β is given by (Baltagi, 2008):

$$\hat{\beta}_{GLS} = (X'\Omega^{-1}X)^{-1}X'\Omega^{-1}y$$

Where X is the matrix of regressors in the original model, y is the vector of dependent variables in the original model, and Ω^{-1} is the inverse of the covariance matrix Ω .

The covariance matrix Ω can be decomposed into:

$$\Omega = \sigma_u^2 J + \sigma_\epsilon^2 I$$

The inverse of Ω is calculated as (Wooldridge, 2010):

$$\Omega^{-1} = \frac{1}{\sigma_\epsilon^2} \left(I - \frac{\sigma_u^2}{\sigma_\epsilon^2 + \sigma_u^2} J \right)$$

This decomposition allows for the efficient computation of Ω^{-1} and consequently, the GLS estimator. The GLS estimator adjusts for the correlation structure of the random effects, leading to more efficient and consistent parameter estimates than OLS when random effects are significant. In practice, estimating σ_u^2 and σ_ϵ^2 is crucial for applying for GLS. These parameters can be estimated using Maximum Likelihood Estimation (MLE) or Generalized Method of Moments (GMM). The MLE approach maximizes the likelihood function derived from the assumed

distribution of the errors, while GMM uses sample moments to estimate the variance components. GLS is a powerful method for handling the complexities of random effects in panel data models. By transforming the model to remove the random effects and applying OLS to the transformed data, GLS provides more efficient and reliable estimates. However, accurate specification of the covariance matrix is essential for the robustness of the GLS estimator (Greene, 2012; Wooldridge, 2010).

Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation (MLE) is a powerful statistical method used to estimate parameters in various models, including Random Effects (RE) models in panel data settings. MLE involves specifying a likelihood function based on the assumed distribution of the errors and then finding the parameter values that maximize this likelihood function. In the context of random effects models, MLE allows for estimating both the regression coefficients and the variance components associated with the random effects and the idiosyncratic errors. In the Random Effects model, the observed data are given by (Greene, 2012):

$$Y_{it} = X_{it}\beta + u_i + \epsilon_{it}$$

Where Y_{it} is the dependent variable for entity i at time t , X_{it} is a vector of regressors, β is the vector of coefficients to be estimated, u_i represents the random effect specific to entity i , and ϵ_{it} is the idiosyncratic error term.

The random effects u_i and the idiosyncratic errors ϵ_{it} are assumed to be normally distributed:

$$\begin{aligned} u_i &\sim \mathcal{N}(0, \sigma_u^2) \\ \epsilon_{it} &\sim \mathcal{N}(0, \sigma_\epsilon^2) \end{aligned}$$

The combined error term $\epsilon_{it} + u_i$ is therefore normally distributed with mean 0 and variance $\sigma_u^2 + \sigma_\epsilon^2$. The likelihood function for observing Y_{it} given X_{it} and β is:

$$L(\beta, \sigma_u^2, \sigma_\epsilon^2) = \prod_{i=1}^N \prod_{t=1}^T \frac{1}{\sqrt{2\pi(\sigma_u^2 + \sigma_\epsilon^2)}} \exp \left(-\frac{(Y_{it} - \beta_0 - \beta_1 X_{it} - u_i)^2}{2(\sigma_u^2 + \sigma_\epsilon^2)} \right)$$

Where β_0 and β_1 are the coefficients to be estimated, σ_u^2 is the variance of the random effect, and σ_ϵ^2 is the variance of the idiosyncratic error term.

To find the Maximum Likelihood Estimators (MLE) for β, σ_u^2 , and σ_ϵ^2 , we need to maximize the likelihood function or, equivalently, the log-likelihood function. The log-likelihood function is given by:

$$\begin{aligned} \ell(\beta, \sigma_u^2, \sigma_\epsilon^2) = & -\frac{NT}{2} \log(2\pi) - \frac{NT}{2} \log(\sigma_u^2 + \sigma_\epsilon^2) \\ & - \frac{1}{2(\sigma_u^2 + \sigma_\epsilon^2)} \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \beta_0 - \beta_1 X_{it} - u_i)^2 \end{aligned}$$

To find the MLE estimators, we need to:

- 1) Maximize the Log-Likelihood Function: This involves taking the partial derivatives of $\ell(\beta, \sigma_u^2, \sigma_\epsilon^2)$ with respect to $\beta_0, \beta_1, \sigma_u^2$, and σ_ϵ^2 , setting these derivatives to zero, and solving the resulting equations.
- 2) Estimate the Variance Components: The variance components σ_u^2 and σ_ϵ^2 are estimated by equating the sample moments to the moments implied by the model. These are typically estimated using the following approach:

$$\begin{aligned} \hat{\sigma}_u^2 &= \frac{1}{N} \sum_{i=1}^N (\hat{u}_i - \bar{\hat{u}})^2 \\ \hat{\sigma}_\epsilon^2 &= \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T (Y_{it} - \hat{\beta}_0 - \hat{\beta}_1 X_{it} - \hat{u}_i)^2 \end{aligned}$$

Where \hat{u}_i represents the estimated random effect for entity i , and $\hat{\beta}_0$ and $\hat{\beta}_1$ are the estimated coefficients from the likelihood maximization. MLE requires iterative optimization techniques such as the Expectation-Maximization (EM) algorithm or numerical optimization methods like Newton-Raphson or BFGS to maximize the log-likelihood function. The choice of method depends on the model's complexity and the optimization algorithm's convergence properties. MLE is a robust and flexible method for estimating parameters in Random Effects models. By

maximizing the likelihood function, MLE provides efficient estimates of both the regression coefficients and the variance components. The approach accommodates the correlation introduced by random effects and is preferred when the normal assumption of errors holds true. However, accurate estimation of variance components and appropriate numerical techniques are crucial for reliable MLE results (Greene, 2012; Wooldridge, 2010).

Restricted Maximum Likelihood (REML)

Restricted Maximum Likelihood (REML) is a refinement of the Maximum Likelihood Estimation (MLE) method that addresses the problem of bias in variance component estimation when the number of parameters to be estimated is large relative to the sample size. REML provides a more accurate estimation of variance components by focusing on the likelihood of the residuals, thereby correcting for the estimation of fixed effects that could bias the variance estimates. In the context of a Random Effects model, the model can be written as (Laird & Ware, 1982):

$$Y_{it} = X_{it}\beta + u_i + \epsilon_{it}$$

where Y_{it} is the dependent variable, X_{it} is the vector of regressors, β is the vector of regression coefficients, u_i represents the random effect specific to entity i , and ϵ_{it} is the idiosyncratic error term. The random effects u_i and the idiosyncratic errors ϵ_{it} are assumed to follow normal distributions:

$$\begin{aligned} u_i &\sim \mathcal{N}(0, \sigma_u^2) \\ \epsilon_{it} &\sim \mathcal{N}(0, \sigma_\epsilon^2) \end{aligned}$$

Thus, the total variance of Y_{it} is $\sigma_u^2 + \sigma_\epsilon^2$, and the covariance matrix of the errors $\epsilon_{it} + u_i$ can be expressed as (Greene, 2012):

$$\Omega = \sigma_u^2 J + \sigma_\epsilon^2 I$$

Where J is an $N \times N$ matrix with all elements equal to 1 (representing the covariance between random effects), and I is the identity matrix (representing the variance of the idiosyncratic errors). The likelihood function for Y_{it} given X_{it} and β is:

$$L(\beta, \sigma_u^2, \sigma_\epsilon^2) = \prod_{i=1}^N \prod_{t=1}^T \frac{1}{\sqrt{2\pi(\sigma_u^2 + \sigma_\epsilon^2)}} \exp \left(-\frac{(Y_{it} - \beta_0 - \beta_1 X_{it} - u_i)^2}{2(\sigma_u^2 + \sigma_\epsilon^2)} \right)$$

REML modifies the likelihood function to adjust for the estimation of the fixed effects, providing an unbiased estimate of the variance components. The REML likelihood focuses on the residuals after accounting for the fixed effects. The REML estimator is obtained by maximizing the restricted likelihood function, which is given by (Verbeke & Molenberghs, 2000):

$$L_{\text{REML}}(\sigma_u^2, \sigma_\epsilon^2) = \text{constant} \times |\Omega|^{-T/2} \exp \left(-\frac{1}{2} \text{tr} [\Omega^{-1}(Y - X\hat{\beta})'(Y - X\hat{\beta})] \right)$$

Where $\hat{\beta}$ represents the Ordinary Least Squares (OLS) estimates of the coefficients β , and tr denotes the trace of a matrix, which is the sum of its diagonal elements. The constant term in the REML function involves the determinant of Ω , which normalizes the likelihood function. To obtain the REML estimates, the following steps are typically taken:

- 1) Obtain OLS Estimates: First, estimate $\hat{\beta}$ using OLS. These estimates are used to compute the residuals $Y - X\hat{\beta}$.
- 2) Compute the Residual Covariance Matrix: The covariance matrix of the residuals is used to adjust for the fixed effects and estimate the variance components.
- 3) Maximize the Restricted Likelihood Function: The variance components σ_u^2 and σ_ϵ^2 are estimated by maximizing $L_{\text{REML}}(\sigma_u^2, \sigma_\epsilon^2)$ with respect to these parameters.

The REML estimator has desirable properties, including being unbiased for variance components and providing efficient estimates when the number of parameters is large relative to the sample size. Implementing REML requires numerical optimization techniques due to the complexity of the likelihood function. Common approaches include iterative algorithms such as the Expectation-Maximization (EM) algorithm or quasi-Newton methods. REML is particularly useful in large panel datasets where fixed effects estimation could bias variance component estimates. REML provides a refined estimation technique for Random Effects models by focusing on the likelihood of residuals rather than the full likelihood function. This approach corrects for biases in variance component estimation and is preferred in cases where the estimation of fixed

effects could otherwise lead to biased estimates. By accurately estimating the variance components, REML enhances the robustness of panel data analysis (Laird & Ware, 1982; Verbeke & Molenberghs, 2000).

3.4. Generalized Method of Moments (GMM)

Dynamic panel data have a large number of individuals (N) and a limited number of time points (T), and the model expressing these data can be written in the form of panel AR (1) as follows:

$$y_{it} = \lambda y_{i,t-1} + x_{it} \beta + u_{it} \quad (2.1)$$

Where:

$t = 2, \dots, T$ and $i = 1, \dots, N$.

λ : is a scalar.

x_{it} : is $1 \times K$ vector of K independent variables for specific individual (i) and at specific time (t).

β : is a $K \times 1$ vector of parameters.

u_{it} : is the error term that can be expressed as:

$$u_{it} = \mu_i + v_{it} \quad (2.2)$$

Where:

μ_i : unobservable individual specific effect.

v_{it} : idiosyncratic error that varies across individuals and through time.

Under the following assumptions:

B1. $\mu_i \sim IID(0, \sigma_\mu^2)$.

B2. $v_{it} \sim IID(0, \sigma_v^2)$.

B3. μ_i and v_{it} are independent.

B4. x_{it} is correlated with the unobservable individual specific effect μ_i , i.e., $E(x_{it}\mu_i) \neq 0$.

B5. x_{it} is strictly exogenous, i.e., $E(x_{it}v_{is}) = 0$; $s = 1, \dots, T$.

The model in (2.1) can be rewritten in matrix form as:

$$y = \lambda y_{-1} + X \beta + u \quad (2.3)$$

Where:

$y = [y_{12}, \dots, y_{N2}, \dots, y_{1T}, \dots, y_{NT}]$ is $N(T - 1) \times 1$ vector of dependent variable observations.

$y_{-1} = [y_{11}, \dots, y_{N1}, \dots, y_{1,T-1}, \dots, y_{N,T-1}]$ is $N(T - 1) \times 1$ vector of lagged observations of the dependent variable.

X : is $N(T - 1) \times K$ matrix of observations of K independent variables.

β : is $K \times 1$ vector of parameters.

$u = [u_{12}, \dots, u_{N2}, \dots, u_{1T}, \dots, u_{NT}]$ is a $N(T - 1) \times 1$ vector of error terms which can be written as:

$$u = z_{\mu} \mu + v$$

Where:

$z_{\mu} = I_N \otimes l_{T-1}$ is $N(T - 1) \times N$ matrix.

I_N : is identity matrix of order N .

l_{T-1} : is $(T - 1) \times 1$ vector of ones.

$v = [v_{12}, \dots, v_{N2}, \dots, v_{1T}, \dots, v_{NT}]$ is a $N(T - 1) \times 1$ vector of unobservable individual specific effects.

Moreover (2.3) can be rewritten as:

$$y = Z \delta + u \quad (2.4)$$

Where:

$\delta' = [\lambda \quad \beta']$ is $1 \times (K + 1)$ vector of parameters.

$Z = [y_{-1} \quad X]$ is $N(T - 2) \times (K + 1)$ matrix of observations of $(K + 1)$ regressors.

Since the dependent variable is a function of the individual specific effect, the lagged dependent variable also is a function of the cross-section-specific effect. In other words, they are correlated, and the endogeneity problem appears, resulting in inconsistent least squares estimators. Consequently, other estimation methods that eliminate the individual effect using an appropriate transformation are recommended.

The individual effect μ_i can be eliminated using the first difference transformation yielding:

$$\Delta y_{it} = \lambda \Delta y_{it-1} + \Delta x_{it} \beta + \Delta v_{it} \quad (2.5)$$

Where:

$i = 1, \dots, N$ and $t = 3, \dots, T$.

The differenced model in (2.5) might be written in matrix form as:

$$\Delta y = \lambda \Delta y_{-1} + \Delta X \beta + \Delta v$$

Where:

$\Delta y = [\Delta y_{13}, \dots, \Delta y_{1T}, \dots, \Delta y_{N3}, \dots, \Delta y_{NT}]'$ is $N(T-2) \times 1$ vector of differenced dependent variable observations.

$\Delta y_{-1} = [\Delta y_{12}, \dots, \Delta y_{1,T-1}, \dots, \Delta y_{N2}, \dots, \Delta y_{N,T-1}]'$ is $N(T-2) \times 1$ vector of lagged differenced dependent variable observations.

ΔX : is $N(T-2) \times K$ matrix of differenced observations of K independent variables.

$\Delta v = [\Delta v_{13}, \dots, \Delta v_{1T}, \dots, \Delta v_{N3}, \dots, \Delta v_{NT}]'$ is $N(T-2) \times 1$ vector of differenced unobservable individual specific effects.

Moreover

$$\Delta y = \Delta Z \delta + \Delta v \quad (2.6)$$

Where:

$\Delta Z = [\Delta y_{-1} \ \Delta X]$ is $N(T-2) \times (K+1)$ matrix of differenced observations of $(K+1)$ regressors.

First Difference Generalized Method of Moments Estimator

To find a consistent estimator, Arellano and Bond (1991) used instruments that are not correlated with the differenced error term in the differenced model (2.5), i.e., orthogonality conditions:

When $t = 3$:

$$y_{i3} - y_{i2} = \lambda (y_{i2} - y_{i1}) + (x_{i3} - x_{i2})\beta + (v_{i3} - v_{i2})$$

$y_{i1}, x_{i1}, x_{i2}, x_{i3}$ are valid instruments since they are not correlated with $(v_{i3} - v_{i2})$.

When $t = 4$:

$$y_{i4} - y_{i3} = \lambda (y_{i3} - y_{i2}) + (x_{i4} - x_{i3})\beta + (v_{i4} - v_{i3})$$

$y_{i1}, y_{i2}, x_{i1}, x_{i2}, x_{i3}, x_{i4}$ are valid instruments since they are not correlated with $(v_{i4} - v_{i3})$.

As a general case, for $t = T$:

$$y_{iT} - y_{i,T-1} = \lambda (y_{i,T-1} - y_{i,T-2}) + (x_{iT} - x_{i,T-1})\beta + (v_{iT} - v_{i,T-1})$$

$y_{i1}, y_{i2}, \dots, y_{i,T-2}, x_{i1}, x_{i2}, \dots, x_{iT}$ are valid instruments since they are not correlated with $(v_{iT} - v_{i,T-1})$.

The valid instruments for each cross-sectional unit i in the GMM method can be defined as:

W_i

$$= \begin{bmatrix} [y_{i1} & x_{i1} & x_{i2} & x_{i3}] & 0 & 0 & 0_k & 0_k & 0_k & 0_k & \dots & 0 & \dots & 0 & 0_k & \dots & 0_k \\ 0 & 0_k & 0_k & 0_k & [y_{i1} & y_{i2} & x_{i1} & x_{i2} & x_{i3} & x_{i4}] & \dots & 0 & \dots & 0 & 0_k & \dots & 0_k \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \dots & \vdots & \vdots & \dots & \vdots \\ 0 & 0_k & 0_k & 0_k & 0 & 0 & 0_k & 0_k & 0_k & 0_k & \dots & [y_{iT-1} & \dots & y_{iT-2} & x_{iT-1} & \dots & x_{iT}] \end{bmatrix} \quad (2.7)$$

W_i : is $(T - 2) \times P$ matrix, where $P = (T - 2)[(T - 1) + (T + 3)K]/2$.

0_k : is a $1 \times K$ vector of zeros.

For all cross-sectional units is:

$$W = [W_1', \dots, W_N']' \quad (2.8)$$

W: is $N(T - 2) \times P$ matrix of valid instruments for all cross-sectional units in DIF GMM.

P: is the number of columns of matrix W.

The orthogonality conditions for the instrumental variables in the DIF GMM method are:

$$E[y_{i,t-s}\Delta v_{it}] = 0 \quad \text{for } t = 3, \dots, T \text{ and } s = 2, \dots, T \quad (2.9)$$

$$E[x_{is}\Delta v_{it}] = 0 \quad \text{for } t = 3, \dots, T \text{ and } s = 1, \dots, T \quad (2.10)$$

Alternatively, using the matrix of valid instruments in (2.8), the moment conditions might be defined as:

$$M = W'\Delta v \quad (2.11)$$

M: is $P \times 1$ vector of moment conditions in case of DIF GMM.

4. Applied

Accurate inflation forecasting is a cornerstone of effective monetary policy, fiscal planning, and strategic investment in today's interconnected global economy. Traditional econometric techniques, particularly time series models, often face significant challenges in capturing both cross-sectional heterogeneity and dynamic patterns over time. These limitations can lead to biased or inefficient estimates, especially when analyzing macroeconomic phenomena such as inflation across countries. Panel data methodologies offer a powerful alternative by integrating both cross-sectional and time series dimensions, allowing for greater statistical power, better control for unobserved heterogeneity, and the ability to capture complex temporal dynamics (Baltagi, 2008; Wooldridge, 2010).

4.1. Dataset Overview

This study employs a rich macroeconomic panel dataset consisting of annual observations for 195 countries spanning the period 1980 to 2024. The data captures both developed and developing economies, enabling a robust comparative analysis of inflation dynamics across a wide range of economic structures.

The dataset is constructed primarily from the International Monetary Fund's World Economic Outlook (IMF WEO) database, which provides standardized, internationally comparable macroeconomic indicators, except TRWMA (Tariff Rate, Applied, Weighted Mean, All Products) is sourced from the World Bank to supplement the international trade perspective.

The dataset is designed to enable robust econometric modeling of inflation dynamics, with the core outcome variable being PCPIPCH—Inflation, average consumer prices. This variable is central to monetary policy research, reflecting cost-of-living changes and acting as a key macroeconomic target for central banks.

The regressors were chosen based on economic theory, policy relevance, and empirical findings from the literature on inflation determinants. They encompass a broad array of domains: public finance, economic output, exchange rate mechanisms, international trade, savings and investment behavior, and labor market indicators.

Each variable is coded consistently and standardized across countries and years, enabling comparability and efficient model estimation in both static and dynamic panel frameworks.

Below is a structured overview of the dataset’s variables categorized by thematic economic domains:

Table 1) Classification and Description of Variables:-

Country Metadata:

Code	Description	Unit
Country_Code	Unique numeric ID assigned to each country	Integer ID
Country	Country name (ISO)	String
Advanced_Country	Development status: advanced (1) or developing (0)	Boolean
Years	Year (range: 2000–2024)	Date

Inflation & Price Stability:

Code	Description	Unit
PCPIPCH	Target variable: Inflation rate (average consumer prices)	Percent change

This variable captures headline inflation and is a proxy for price stability. It reflects consumer cost changes and serves as a central dependent variable for modeling inflation forecasting.

Public Finance Indicators:

Code	Description	Unit
GGR_NGDP	Government revenue	% of GDP
GGX_NGDP	Total government expenditure	% of GDP
GGSB_NPGDP	Structural balance (cyclically adjusted)	% of potential GDP

GGXWDG_NGDP	Gross government debt	% of GDP
--------------------	-----------------------	----------

These indicators gauge fiscal stance and sustainability, with GGSB_NPGDP controlling for the business cycle—key in determining fiscal-inflationary pressures.

Economic Output & Productivity:

Code	Description	Unit
NGDP_RPCH	Real GDP growth (constant prices)	Percent change
PPPPC	GDP per capita at PPP-adjusted current prices	International dollars
PPPSH	Country's share in world GDP based on PPP	Percent

These variables help identify demand-side drivers of inflation and allow for convergence or divergence analysis across income groups.

Exchange Rate and Purchasing Power:

Code	Description	Unit
PPPEX	Implied PPP conversion rate	Local currency per intl. dollar

A lower PPPEX signals currency overvaluation, potentially dampening net exports and inflation. This variable acts as a proxy for real exchange rate distortions.

International Trade & Balance of Payments:

Code	Description	Unit
TX_RPCH	Export volume growth	Percent change
TM_RPCH	Import volume growth	Percent change
BCA_NGDPD	Current account balance	% of GDP
TRWMA	Weighted average applied tariff rate (goods)	Percent

Trade openness and competitiveness indicators, such as BCA_NGDPD and TRWMA, affect import-price pass-through, thus shaping inflation dynamics.

Savings & Investment Behavior:

Code	Description	Unit
NGSD_NGDP	Gross national savings	% of GDP
NID_NGDP	Total fixed investment	% of GDP

Savings-investment gaps are foundational in macroeconomic balance analysis and impact long-term inflation and growth trajectories.

Labor Market Conditions:

Code	Description	Unit
LUR	Unemployment rate	% of total labor force

High unemployment may suppress wage-push inflation, while tight labor markets may elevate inflationary pressures.

Together, these variables provide a holistic view of the macroeconomic environment in each country. Their inclusion supports the exploration of both short-run dynamics and long-run equilibrium relationships in inflation modeling, particularly under panel data econometric frameworks such as fixed effects, random effects, and GMM-based dynamic models.

This structure is specifically suited to identify structural inflation patterns, measure cross-country heterogeneity, and improve the predictive power of inflation forecasts under macroeconomic uncertainty.

In the next phase of the study, we proceed to data exploration, cleaning, and preprocessing, which includes handling missing values, detecting outliers, and performing feature selection to identify the most relevant predictors for inflation. These steps are essential to ensure the reliability and robustness of subsequent econometric estimations.

4.2. Dataset Preprocessing

This stage involves a comprehensive descriptive analysis and systematic data preprocessing, serving as the foundation for reliable econometric modeling. The primary goals are to address missing values, correct data inconsistencies, and identify outliers that may compromise the validity of statistical inferences. In parallel, we conduct an initial feature selection process to identify the most relevant macroeconomic indicators associated with inflation. These variables span key domains such as fiscal policy, international trade, labor market dynamics, and monetary conditions.

To uphold the statistical integrity of the econometric framework, we implement a structured data preparation pipeline. This includes the imputation or removal of missing observations, standardization of continuous variables, and binary encoding of categorical variables (e.g., the development status of countries). Additionally, we perform correlation analysis, Variance Inflation Factor (VIF) checks, and distributional assessments to detect multicollinearity, skewness, and other distortions that could bias model estimates.

Missing (Null) values:

Table 2) Number of Missing Values by Variable

Variable	Missing Values	Variable	Missing Values
WEO_Country_Code	0	Country	0
Advanced_Country	0	Year	0
BCA_NGDPD	1182	GGR_NGDP	2327
GGSB_NPGDP	6304	GGXWDG_NGDP	3049
GGX_NGDP	2380	LUR	4681
NGDP_RPCH	859	NGSD_NGDP	1859
NID_NGDP	1852	PCPIPCH	906
PPPEX	946	PPPPC	901
PPPSH	967	TM_RPCH	1907
TRWMA	5132	TX_RPCH	1863

Given the substantial number of missing values across several variables, the initial step in the data cleaning process involves excluding countries with more than 50% missing observations across the entire dataset or more than 70% missing values in any single variable within a country. These thresholds ensure that the retained countries contain sufficient information for meaningful econometric analysis, thereby reducing the risk of biased estimates and preserving the overall integrity of the panel structure.

Figure 2) Number of Missing Values by Country



After filtering the dataset, 70 countries were retained for further analysis. The next step was to define an appropriate time window that balances data availability and analytical relevance. Based on a preliminary visual inspection of missingness patterns over time, the period 2000 to 2024 was selected. This range offers sufficient temporal depth for panel-based macroeconomic analysis while minimizing the impact of missing data.

Iterative Imputation with a Bayesian Ridge Estimator (MICE-BR)

To address the issue of incomplete data, we employed Iterative Imputation using a Bayesian Ridge estimator—a method commonly known as Multivariate Imputation by Chained Equations (MICE). This technique treats each variable with missing values as a regression problem, where missing entries are estimated conditionally based on the other observed variables in the dataset, within a fully Bayesian regression framework. The imputation process is executed separately for each country, preserving heterogeneity in country-specific macroeconomic profiles. This method is particularly well-suited for panel data, as it allows for the imputation of missing values without discarding the rich inter-variable relationships.

Let the dataset $X \in \mathbb{R}^{n \times p}$, where n is the number of observations and p is the number of variables. For any variable X_j with missing data, and all remaining variables X_{-j} , the imputation proceeds as follows:

1. Initialization: All missing entries are filled with initial values (e.g., the column mean).
2. Iterative Regression:
 - Each incomplete variable X_j is regressed on the other variables X_{-j} .
 - Missing values in X_j are updated using predictions from the model.
 - The process is repeated for all variables with missing data.
3. Repeat: Steps are iterated multiple times (e.g., 20 iterations) until convergence.

Mathematical Formulation

At each iteration, the following regression is estimated:

$$X_j = X_{-j}\beta_j + \epsilon_j$$

Where:

- X_j : is the target variable (with missing values),
- X_{-j} : is the matrix of predictors,
- β_j : are the coefficients to be estimated,
- $\varepsilon_j \sim \mathcal{N}(0, \sigma^2 I)$ is the Gaussian error term.

In Bayesian Ridge Regression, the assumptions are:

$$\beta_j \sim \mathcal{N}(0, \lambda^{-1} I), \quad \varepsilon_j \sim \mathcal{N}(0, \alpha^{-1} I)$$

Where:

- λ : controls the prior variance of the coefficients (regularization),
- α : controls the noise level.

The posterior estimate of the coefficients is:

$$\hat{\beta}_j = (X_{-j}^T X_{-j} + \lambda I)^{-1} X_{-j}^T X_j$$

The posterior covariance is:

$$\Sigma_j = (\alpha X_{-j}^T X_{-j} + \lambda I)^{-1}$$

The missing values in X_j are then imputed from the predictive distribution:

$$\widehat{X_j^{\text{miss}}} \sim \mathcal{N}(X_{-j}^{\text{miss}} \hat{\beta}_j, \sigma^2)$$

Python Implementation

```
# Handle missing values for each country
def handle_missing_country(country):
    imputer = IterativeImputer(estimator=BayesianRidge(), max_iter=20,
random_state=0, verbose=2)
    country[cols_with_na] = imputer.fit_transform(country[cols_with_na])
    return country
# Apply the imputation process by country
df_interpolated = df_panel_b.groupby('Country',
group_keys=False).apply(handle_missing_country)
```

Why Bayesian Ridge Regression?

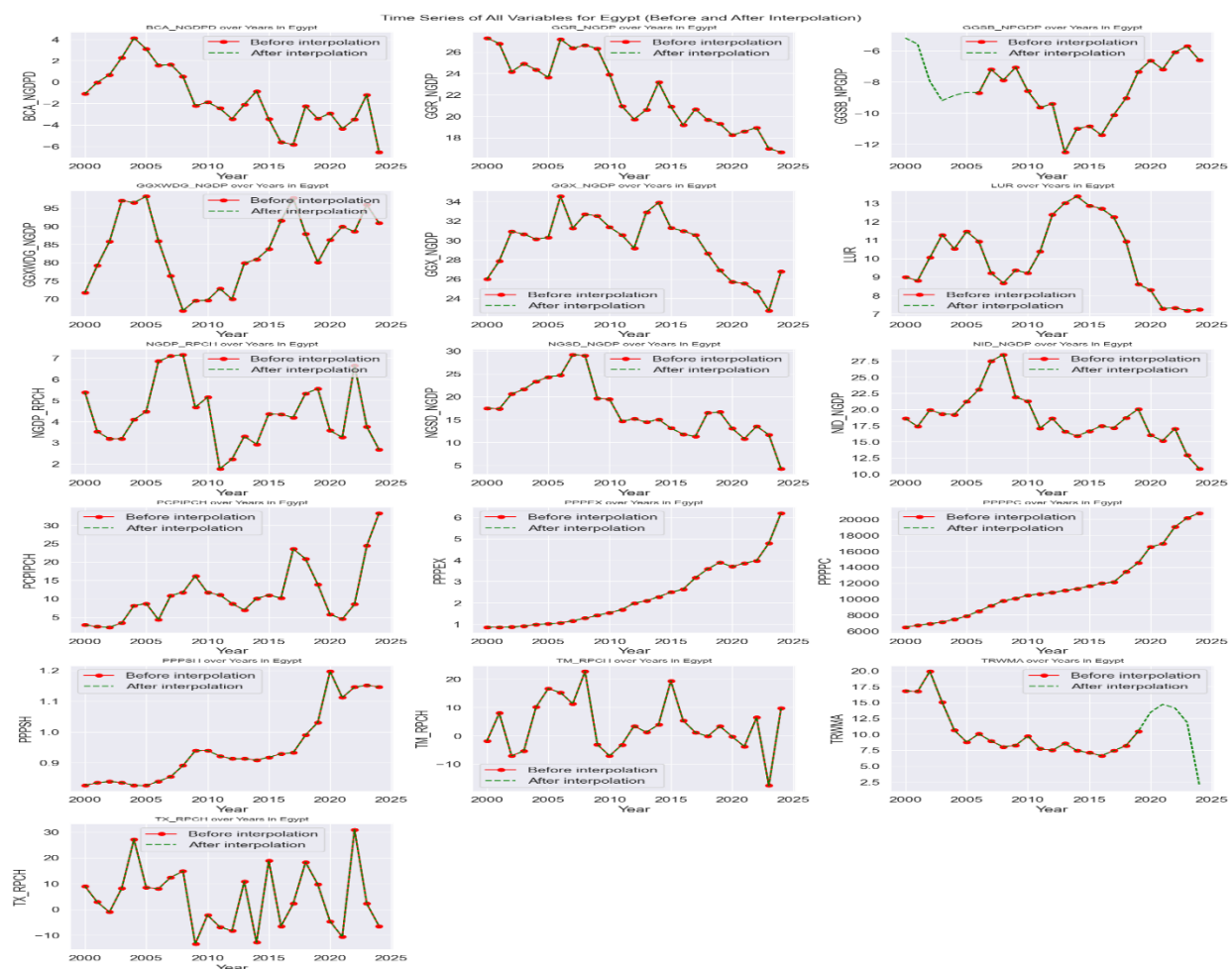
Bayesian Ridge Regression offers several advantages that make it particularly well-suited for imputing missing values in macroeconomic panel data.

First, it handles multicollinearity effectively through coefficient regularization, which reduces overfitting and ensures stable estimates when explanatory variables are highly correlated—a common feature in macroeconomic datasets.

Second, it is robust to small or noisy samples, making it appropriate for country-level time series where data quality and availability vary.

Third, as a Bayesian method, it provides distributional uncertainty estimates, allowing for probabilistic imputation and enhancing the credibility of inference.

Figure 3) Apply MICE-BR to missing values for Egypt



Descriptive Statistics

Table 3) Descriptive Statistics

Variable	Mean	Std	Min	25%	50%	75%	Max
Advanced_Country	0.4857	0.4999	0.0000	0.0000	0.0000	1.0000	1.0000
PCPIPC	5.1843	10.7675	-3.7230	1.5530	2.8615	5.5613	229.8240
GGR_NGDP	33.2174	11.1669	12.3520	23.3208	33.9800	41.5563	63.0530
GGX_NGDP	35.5736	11.4037	9.0150	25.9293	36.3735	44.5823	74.4160
GGSB_NPGDP	-2.4370	3.1935	-21.489	-4.3193	-2.2340	-0.3628	11.1370
GGXWDG_NGDP	57.3924	36.6984	0.0520	35.6383	49.7325	71.7825	258.3740
NGDP_RPCH	2.9226	3.8172	-28.759	1.3683	2.9860	4.9878	24.6160
PPPPC	29144.1583	19788.0113	3230.4700	13391.9858	24488.9180	40951.2930	148185.7760
PPPSH	0.9781	2.2242	0.0030	0.0970	0.3435	0.9235	20.5020
PPPEX	130.8898	550.0059	0.0160	0.6523	1.2665	9.8068	4917.8250
TX_RPCH	4.5983	8.9785	-57.613	0.7440	4.3065	8.8798	72.6580
TM_RPCH	4.8340	10.1164	-53.558	0.3045	4.9595	9.9955	68.4777
BCA_NGDPD	-0.8665	6.0990	-24.068	-4.3380	-1.3590	2.1575	30.1650
TRWMA	3.8135	3.6181	-4.0559	1.7200	2.3200	4.6152	26.3800
NGSD_NGDP	22.5413	8.1673	-76.451	17.8230	22.1465	26.6453	60.7880
NID_NGDP	23.4828	5.6027	8.9320	19.9790	22.8655	25.8663	68.2070
LUR	8.4863	5.3230	0.7000	5.0240	7.1735	10.0628	34.3000

Temporal Coverage: The dataset spans the years 2000 to 2024, with both the mean and median year equal to 2012, indicating balanced temporal coverage across the 25-year period. This makes it suitable for long-term macroeconomic trend analysis and panel modeling.

Economic Indicators: Current account balances (BCA_NGDPD) show a wide dispersion from -24.07% to 30.17% of GDP, with an average of -0.87%, suggesting that many countries operate under trade deficits. Government revenue (GGR_NGDP) and expenditure (GGX_NGDP) average around 33.2% and 35.6% of GDP, respectively, indicating moderate fiscal activity across the sample.

Fiscal Balance and Public Debt: The structural fiscal balance (GGSB_NPGDP) averages -2.44%, reflecting a general tendency toward deficit spending. Public debt (GGXWDG_NGDP) ranges from virtually zero to 258.37% of GDP, with a high standard deviation (36.69%), indicating substantial heterogeneity in debt levels across countries.

Unemployment and Economic Growth: The unemployment rate (LUR) varies between 0.7% and 34.3%, with an average of 8.49%, capturing a wide spectrum of labor market conditions. Both real GDP growth (NGDP_RPCH) and inflation (PCPIPCH) exhibit high variability, reflecting divergent macroeconomic cycles and price volatility across nations.

Purchasing Power and Price Indices: GDP per capita in PPP terms (PPPPC) shows extreme variation, with a standard deviation of 19,788, highlighting significant global income disparities. Other indicators like PPPEX (PPP conversion rate), PPPSH (share of world PPP GDP), and PCPIPCH (inflation) further confirm differences in price levels and economic scale among countries.

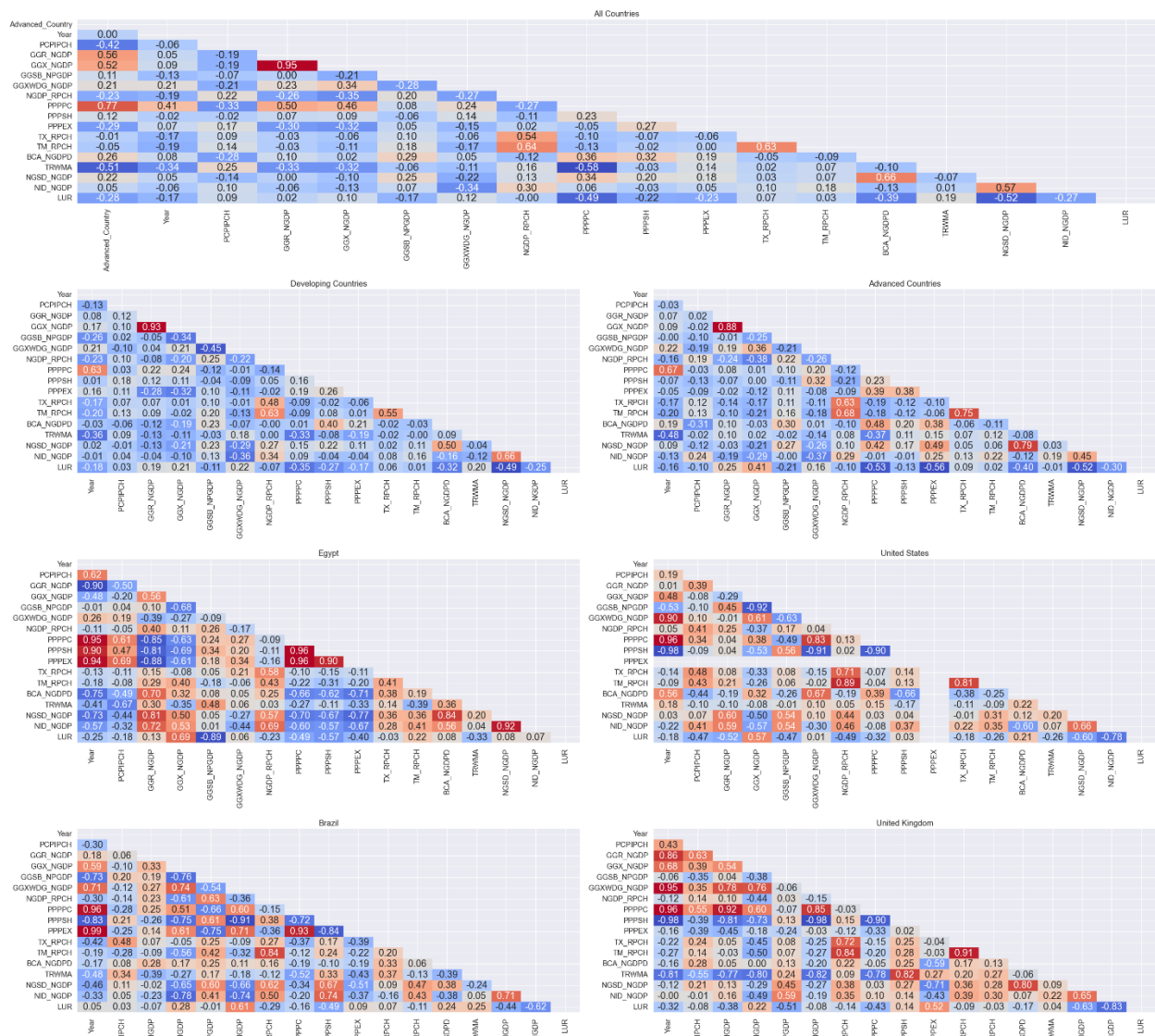
Trade Indicators: Import (TM_RPCH) and export (TX_RPCH) growth rates exhibit extensive variation, with both negative and positive extremes, indicating volatile trade performance. The tariff-weighted mean (TRWMA) reaches a maximum of 26.38%, though the average remains modest (3.81%), suggesting a few countries maintain high trade barriers, while most are relatively open.

Country Representation and Panel Structure: Each of the 70 countries contributes exactly 25 annual records, confirming the presence of a balanced panel structure. This uniformity strengthens the validity of comparative econometric analysis and supports the use of panel data techniques that require consistency over time.

Correlation Analysis

Correlation analysis is conducted to examine the linear relationships among explanatory variables and their association with the target variable (inflation, PCPIPCH). This step provides valuable insights into inter-variable dependencies, helps identify potential issues of multicollinearity, and highlights which predictors have the strongest association with inflation. Understanding these correlations is essential for guiding variable selection, refining model specification, and ensuring the robustness of subsequent econometric analyses.

Figure 4) Correlation



Variable Selection Strategy and Dimensional Representation

In order to optimize the econometric specification and ensure both model parsimony and statistical validity, the analysis adopts a selective variable inclusion strategy. Specifically, only one or at most two representative indicators are selected from each macroeconomic dimension. These variables are chosen based on their ability to capture core variations in the respective domain, serve as theoretically sound proxies, and reduce potential issues associated with multicollinearity and overparameterization—particularly common when multiple variables measure overlapping economic concepts.

This selection is further supported by an extensive correlation matrix analysis, conducted both at the global level (all countries) and for selected individual countries. Variables were retained based on their significant and stable correlation with inflation (PCPIPCH), consistency with theoretical expectations, and absence of multicollinearity with other retained indicators.

This methodological choice not only increases the efficiency of parameter estimation but also preserves interpretability and identification power in the panel data framework. The selected representatives include:

GGSB_NPGDP and GGXWDG_NGDP — representing Public Finance Indicators, these variables reflect the fiscal stance and the evolution of general government gross debt, capturing long-term fiscal sustainability and debt burden dynamics.

PPPPC — representing Economic Output & Productivity, this variable reflects GDP per capita based on purchasing power parity. While related to income levels, it also encapsulates living standards and productivity, complementing fiscal indicators.

TX_RPCH and BCA_NGDPD — as proxies for International Trade and Balance of Payments, they capture external sector performance, export dynamics, and macroeconomic imbalances.

NID_NGDP — a core indicator of Savings and Investment Behavior, this variable tracks capital formation and demand-side pressures relevant to inflation and growth trade-offs.

By selecting these core variables, the model ensures that each economic factor is adequately represented without redundancy, enhancing both the robustness and predictive capability of the inflation forecasting models. This practice aligns with the recommendations in empirical Macroeconometrics literature for handling high-dimensional panel data (Baltagi, 2005; Wooldridge, 2010).

Table 4) Variance Inflation Factor (VIF)

Feature	VIF	Feature	VIF	Feature	VIF
GGSB_NPGDP	1.752526	GGXWDG_NGDP	3.403836	PPPPC	3.457779
TX_RPCH	2.082648	TM_RPCH	2.119305	NID_NGDP	4.473428

Following the variable selection phase, the dataset was re-evaluated to identify additional countries that met the completeness criteria under the reduced variable set. As a result, the number of countries increased from 70 to 83, reflecting improved data availability and allowing for a broader, yet statistically reliable, macroeconomic representation. To address the remaining missing values without removing any time periods, the dataset was then completed using Iterative Imputation with a Bayesian Ridge Estimator (MICE-BR). This method imputes missing entries for each country independently, leveraging inter-variable relationships while preserving cross-sectional and temporal structures critical to panel econometric modeling.

Outlier Detection in Panel Data

Following the completion of missing data imputation and variable selection, the dataset underwent a systematic outlier detection process, focusing on identifying countries with extreme inflation behavior that could bias pooled and GMM-based panel estimators. Due to the inherent sensitivity of pooled OLS, fixed-effects models, and particularly dynamic GMM estimators to outlier-driven variance and leverage effects, this step is essential for maintaining statistical robustness and valid inference (Blundell & Bond, 1998; Aghion et al., 2021; Bai et al., 2020).

Outliers were identified based on the distribution of the target variable (PCPIPCH) by computing both the country-level mean and standard deviation of inflation across the full time series. Countries were flagged as outliers if they exhibited both a high average inflation rate and large within-country inflation volatility, exceeding the upper bound of the interquartile range (IQR), defined as:

$$\text{Outlier Threshold} = Q_3 + 1.5 \times (Q_3 - Q_1)$$

This method aligns with standard practices in robust panel data cleaning and was applied to the distribution of country-level mean inflation.

Table 5) Outlier Detection by Country

Country	Mean (PCPIPCH)	Std. Dev. (PCPIPCH)
Peru	288.182333	1210.618708
Brazil	263.782711	627.631647

Angola	249.971733	749.301031
Ukraine	201.772865	708.277703
Belarus	167.894888	374.457601
Kazakhstan	114.168324	315.865945
Russia	60.175944	137.029847
Croatia	55.009429	225.017040
Bulgaria	46.114187	164.329857
Türkiye	40.748644	30.593876
Poland	34.157267	93.751501
Romania	32.933200	60.183799
Uruguay	28.177311	29.550033
Suriname	22.556933	26.559598
Mexico	22.521044	31.586555
Ecuador	20.915778	23.012659

Upon review, the following six countries were identified as statistical outliers in terms of both elevated inflation averages and unusually high volatility, typically associated with historical episodes of hyperinflation, regime shifts, or macroeconomic crises:

Peru, Brazil, Angola, Ukraine, Belarus, and Kazakhstan. These countries were manually flagged for exclusion based on their distortionary influence on the pooled distribution and model residual structure. This manual step also considers the known historical context in which inflation spikes may not be representative of general macroeconomic processes, consistent with approaches used in recent empirical macroeconomic panel studies (Ciccone, 2015; Becker et al., 2010; Wooldridge, 2010).

Other countries with moderately high but economically justifiable inflation rates—such as Türkiye, Romania, and Mexico—were retained, as their inflation behavior aligns with broader macroeconomic trends rather than extreme shock events.

After excluding the six identified outlier countries, the final dataset comprises 77 countries over a 45-year period, resulting in a balanced and refined panel with 3,735 total observations. This cleaned dataset will be used for the next stage: model training and inflation forecasting using panel econometric methods, including fixed effects, random effects, and dynamic GMM estimators.

4.3. Models

Following the data preparation phase, the modeling process begins using a balanced panel of 77 countries over 45 years (1980–2024), totaling 3,735 observations. The objective is to estimate and evaluate alternative econometric models suited for panel data structures, in order to explain and forecast inflation dynamics effectively. The models employed in this study include:

Pooled Ordinary Least Squares (Pooled OLS) Model:

A basic benchmark model that pools all observations and ignores country-specific and time-specific heterogeneity. While easy to estimate, it is inappropriate when there is unobserved heterogeneity or endogeneity, which is often the case in macroeconomic panel data.

```
# Create formula for model
formula_str <- as.formula(paste("PCPIPCH ~ -1 +", paste(X_vars, collapse = " +
")))
# Fit Pooled OLS model
pooled_model <- plm(formula_str, data = train_model_df, model = "pooling", index
= c("Country", "Year"))
# Display model summary
summary(pooled_model)
```

Table 6) (Pooled OLS) Summary:

Adjusted R-squared	0.4516
F-statistic	454.253 (p < 0.001)

Table 7) (Pooled OLS) Coefficients:

Variable	Estimate	Std. Error	p-value
GGSB_NPGDP	-4.0484	0.3178	<0.001
GGXWDG_NGDP	0.3216	0.0262	<0.001
PPPPC	-0.00031	0.00016	0.0507
TX_RPCH	-0.0021	0.2151	0.9921
TM_RPCH	0.7339	0.2022	0.0003
NID_NGDP	0.1820	0.1514	0.2292
Advanced_Country	-19.344	6.2825	0.0021

Interpretation: The Pooled OLS model explains approximately 45% of the variation in inflation rates across the sample, as indicated by the R-squared value of 0.4525. The overall model is statistically significant ($F = 454.25$, $p < 0.001$), suggesting that the selected regressors jointly explain inflation variation across countries and time.

Fixed Effects (FE) Model:

This model accounts for unobserved, time-invariant country-specific characteristics that may correlate with the regressors. It is suitable in this context to capture structural differences across countries—such as policy regimes, geography, or institutional quality—that influence inflation but are not directly observed.

```
# Fit Fixed Effects (within)
fe_model <- plm(formula_str, data = train_model_df, model = "within", index =
c("Country", "Year"))
# Display model summary
summary(fe_model)
```

Table 8) (FE) Summary:

Adjusted R-squared	0.45582
F-statistic	535.948 ($p < 0.001$)

Table 9 (FE) Coefficients:

Variable	Estimate	Std. Error	p-value
GGSB_NPGDP	-4.2066	0.3314	<0.001
GGXWDG_NGDP	0.3169	0.0274	<0.001
PPPPC	-0.00062	0.00018	0.0007
TX_RPCH	0.0321	0.2212	0.8848
TM_RPCH	0.7695	0.2032	0.0002
NID_NGDP	-0.2718	0.4377	0.5347

Interpretation: The Fixed Effects model explains approximately 46.9% of the within-country variation in inflation, as indicated by the R-squared value. The model is statistically significant overall ($F = 535.95$, $p < 0.001$), confirming that the selected regressors jointly explain meaningful within-country variation in inflation rates.

Random Effects (RE) Model:

This model assumes that unobserved individual effects are uncorrelated with the regressors, allowing for the use of both within- and between-country variation. While more efficient under its assumptions, it can be biased if the assumption of orthogonality is violated, which is often tested using the Hausman test.

```
# Fit Random Effects model using plm
re_model <- plm(formula_str, data = train_model_df, model = "random", index =
c("Country", "Year"))
# Display model summary
summary(re_model)
```

Table 10) (RE) Summary

Adjusted R-squared	0.46448
Chisq:	3262.87(p < 0.001)

Table 11 (RE) Coefficients:

Variable	Estimate	Std. Error	p-value
GGSB_NPGDP	-4.1547	0.3254	<0.001
GGXWDG_NGDP	0.3200	0.0269	<0.001
PPPPC	-0.00059	0.00016	0.0002
TX_RPCH	0.0169	0.2178	0.9382
TM_RPCH	0.7598	0.2011	0.0002
NID_NGDP	0.0022	0.2301	0.9925

Interpretation: The Random Effects model explains approximately 46.5% of the variation in inflation across countries and over time. The model is statistically significant overall ($\chi^2 = 3,262.87$, $p < 0.001$), indicating strong joint explanatory power of the selected regressors.

Model Specification Tests: Pooled OLS vs. Fixed Effects vs. Random Effects:

To determine the most appropriate model among Pooled OLS, Fixed Effects, and Random Effects, we applied two standard specification tests: the Wald Test and the Hausman Test.

Hausman Test: Fixed Effects vs Random Effects

The Hausman test evaluates whether the Random Effects (RE) estimator is consistent. It tests the assumption that individual effects are uncorrelated with the regressors. If this assumption is violated, the RE model becomes inconsistent, and the Fixed Effects (FE) model is preferred.

Hypotheses:

H_0 : All individual (country) effects are jointly equal to zero \rightarrow Pooled OLS is sufficient.

H_1 : At least one individual effect is non-zero \rightarrow Fixed Effects model is preferred.

Test Statistic:

$$F = \frac{(RSS_{\text{Pooled}} - RSS_{\text{FE}})/k}{RSS_{\text{FE}}/(N - k - 1)}$$

Where:

RSS_{Pooled} , RSS_{FE} = Residual Sum of Squares, k = number of constraints, N = observations

```
# Wald test: Pooled OLS vs Fixed Effects  
waldtest(pooled_model, fe_model, test = "F")
```

Results: F-statistic = 9.48 (p-value = 0.0021)

Decision: Since the p-value < 0.05 , we reject the null hypothesis.

Conclusion: The test confirms that individual (country-specific) effects are statistically significant, and therefore, the Fixed Effects model is preferred over Pooled OLS. This supports the presence of unobserved heterogeneity across countries, which must be accounted for in the modeling process.

Hausman Test: Fixed Effects vs Random Effects

The Hausman test evaluates whether the Random Effects (RE) estimator is consistent. It tests the assumption that individual effects are uncorrelated with the regressors. If this assumption is violated, the RE model becomes inconsistent, and the Fixed Effects (FE) model is preferred.

Hypotheses:

H_0 : Random Effects model is consistent and efficient $\rightarrow \text{Cov}(\alpha_i, X_{it}) = 0$

H_1 : Random Effects model is inconsistent $\rightarrow \text{Cov}(\alpha_i, X_{it}) \neq 0$

Test Statistic:

$$\chi^2 = (\widehat{\beta}_{FE} - \widehat{\beta}_{RE})' [\text{Var}(\widehat{\beta}_{FE}) - \text{Var}(\widehat{\beta}_{RE})]^{-1} (\widehat{\beta}_{FE} - \widehat{\beta}_{RE})$$

Where:

$\widehat{\beta}_{FE}$, $\widehat{\beta}_{RE}$: Estimated coefficients, $\text{Var}(\cdot)$: Estimated variance-covariance matrices

```
# Hausman test: fixed vs random effects
```

```
phptest(fe_model, re_model)
```

Results: Chi-squared = 2.58 (p-value = 0.8597)

Decision: Since the p-value > 0.05, we fail to reject the null hypothesis.

Conclusion: There is no statistical evidence of correlation between the country-specific effects and the regressors, implying that the Random Effects model is consistent and preferred over the Fixed Effects model in this case. This result justifies the use of the RE specification, especially when leveraging both within-country and between-country variation.

Given these results, the Random Effects model is preferred for the baseline specification, as it is more efficient and consistent under the test outcomes, and is retained for interpretation and further comparison with dynamic models. However, the Fixed Effects model remains valuable for robustness checks, especially when the assumption of no correlation between unobserved effects and regressors is questioned in further applications.

Diagnostic Testing and Model Validity Assessment

In this section, we assess the statistical properties of the residuals from the Random Effects (RE) model, which was selected based on model specification tests. These tests examine key assumptions including homoskedasticity, serial independence, cross-sectional independence, and stationarity of the dependent variable. Ensuring that the assumptions of classical panel estimation are not violated is critical to obtaining valid standard errors and efficient estimators.

Heteroskedasticity Tests (Breusch–Pagan and White)

To test whether the error variances are constant (homoskedastic) across observations. Heteroskedasticity can lead to biased standard errors and invalid statistical inference.

Hypotheses (for both tests):

H_0 : Homoskedasticity (constant variance). H_1 : Presence of heteroskedasticity.

```
# Breusch-Pagan Test
bptest(re_model)
# White Test
fe_residuals <- residuals(re_model)
bptest(lm(fe_residuals^2 ~ ., data = train_model_df[, X_vars]))
```

Results:

Breusch–Pagan Test: Statistic = 13.62 (p-value = 0.0182)

White Test: Statistic = 18.01 (p-value = 0.0062)

Decision: Since both p-values are below 0.05, we reject the null hypothesis of homoskedasticity.

Conclusion: Heteroskedasticity is present, indicating that standard errors may be biased. It is recommended to use robust (heteroskedasticity-consistent) standard errors for inference and hypothesis testing.

Serial Correlation Test Breusch–Godfrey (Wooldridge version)

To detect the presence of autocorrelation in the residuals across time for a given country, which violates the independence assumption of classical linear models.

Hypotheses: H_0 : No serial correlation. H_1 : Serial correlation exists.

```
Serial Correlation test (Breusch-Godfrey) residuals
pbgttest(re_model)
```

Results: Chi-squared = 751.38, p-value < 2.2e-16

Decision: With a p-value near zero, we reject the null hypothesis.

Conclusion: There is strong evidence of serial correlation in the residuals, requiring the use of cluster-robust or serial-correlation-robust standard errors.

Pesaran's Cross-sectional Dependence (CD) Test

To test for cross-sectional dependence, i.e., whether shocks affecting one country also affect others, which is a common feature in macroeconomic panel datasets.

Hypotheses: H_0 : Cross-sectional independence. H_1 : Cross-sectional dependence.

```
pcdtest(re_model, test = "cd") # Test Pesaran CD
```

Results: z-statistic = 22.65, p-value < 2.2e-16

Decision: Reject H_0 significant cross-sectional dependence exists.

Conclusion: There is strong cross-sectional dependence in the dataset. This implies potential bias in standard panel models that ignore these interdependencies. Consideration of cross-sectional robust estimators or factor-augmented panel models may be warranted in future extensions.

Levin, Lin, and Chu (LLC) Panel Unit Root (Stationarity) Test

To test whether the dependent variable (inflation rate: PCPIPCH) is stationary, a prerequisite for many panel regressions and especially dynamic models like GMM.

Hypotheses: H_0 : Non-stationarity (unit root present). H_1 : Stationarity.

```
# test Levin, Lin & Chu
llc_test <- purtest(panel_df$PCPIPCH, test = "levinlin")
summary(llc_test)
```

Results: Overall statistic = -20.16, p-value = 0.0000

Decision: Reject H_0 the variable is stationary.

Conclusion: The inflation rate series (PCPIPCH) is stationary, supporting its suitability for inclusion in the panel model without differencing. This validates its use in both static and dynamic models.

Correcting Standard Errors: Driscoll–Kraay Robust Estimation

Given the results from the diagnostic tests—where heteroskedasticity, serial correlation, and cross-sectional dependence were all detected—the standard errors of the Random Effects model were corrected using the Driscoll–Kraay (SCC) estimator. This method adjusts standard errors to remain consistent in the presence of multiple sources of correlation in macro-panel datasets.

Table 12) Corrected Coefficient Significance – RE Model with (SCC) Standard Errors:

Variable	Estimate	Robust Std. Error	t-value	p-value	Significance
GGSB_NPGDP	-4.1547	4.3543	-0.95	0.3401	
GGXWDG_NGDP	0.3200	0.3864	0.83	0.4076	
PPPPC	-0.00059	0.00020	-2.91	0.0037	**

TX_RPCH	0.0169	0.2415	0.07	0.9442	
TM_RPCH	0.7598	0.2488	3.05	0.0023	**
NID_NGDP	0.0022	0.6100	0.0035	0.9972	

Interpretation of Results

Below is a detailed interpretation of each coefficient in light of economic theory and its corrected statistical significance:

GGSB_NPGDP (General Government Structural Balance): The estimated effect (−4.15) suggests that improved structural fiscal balance is associated with reduced inflation. However, the p-value (0.34) indicates no statistical significance, which may result from high variability across countries and time, especially after correcting for serial and spatial dependence.

GGXWDG_NGDP (Gross Public Debt as % of GDP): The positive estimate (0.32) implies a potential inflationary effect from rising public debt, consistent with debt-driven inflation theory. Nonetheless, the result is not statistically significant ($p = 0.408$) after robust error correction.

PPPPC (GDP per Capita, PPP): A negative and statistically significant coefficient (−0.00059, $p = 0.0037$) suggests that higher income per capita is associated with lower inflation, reflecting stronger monetary institutions and more stable economies in high-income countries.

TX_RPCH (Export Volume Growth): The estimate is near zero (0.017), with an extremely high p-value (0.944), indicating no significant effect of export growth on inflation, once model assumptions are fully corrected.

TM_RPCH (Import Volume Growth): The coefficient (0.76) is positive and statistically significant ($p = 0.0023$). This supports the interpretation that increased imports can contribute to inflation through imported inflation effects or increased aggregate demand.

NID_NGDP (Total Investment as % of GDP): Despite a positive point estimate (0.0022), the result is statistically insignificant ($p = 0.997$). This may suggest that investment-driven inflation effects are not uniform or substantial once corrected for country-level dependencies.

Given the results of previous models, significant heteroskedasticity, serial correlation, and cross-sectional dependence were detected in the panel structure. Additionally, inflation (PCPIPCH) exhibited persistence and potential endogeneity with respect to macroeconomic predictors such as public debt and investment. The Two-Step Difference Generalized Method of Moments (Difference GMM) estimator proposed by Arellano and Bond (1991) is well-suited for this context.

Two-Step Difference Generalized Method of Moments (Difference GMM, Arellano–Bond):

This dynamic panel model is particularly appropriate when addressing endogeneity, autocorrelation, and omitted variable bias, especially due to the inclusion of lagged dependent variables (inflation persistence). In this study, inflation is known to exhibit inertia, and explanatory variables be endogenous. Difference GMM uses past levels of variables as instruments for their first differences, providing consistent estimates even when regressors are not strictly exogenous. The two-step estimator improves efficiency by incorporating a robust weighting matrix, which is crucial in macroeconomic panel settings with a relatively large cross-section (countries) and a moderate time dimension. This model effectively addresses:

- Accounts for endogeneity by instrumenting endogenous variables (e.g., lagged inflation).
- Accounts for autocorrelation by differencing and testing for higher-order serial correlation.
- Corrects for heteroskedasticity using a robust weighting matrix.
- Accounts for omitted variable bias by differencing away unobserved fixed effects.
- Accounts for efficiency by via robust two-step weighting.

```
# Build the dynamic panel formula:
# The left-hand side of the formula includes:
# - lagged dependent variable: lag(PCPIPCH, 1)
# - contemporaneous exogenous variables: fiscal, trade, and income predictors
# The right-hand side of the formula provides instruments:
# - lagged inflation from period 2 to 7 as instruments
dynamic_formula <- as.formula(paste0("PCPIPCH ~ lag(PCPIPCH, 1) + ",
                                     paste(X_vars, collapse = " + "), " | lag(PCPIPCH, 2:7)"))
# -----
# Estimate the GMM model (difference GMM (Arellano–Bond))
# -----
gmm_model <- pgmm(
  formula      = dynamic_formula,
  data         = panel_df,
```

```

effect      = "individual",      # panel-level fixed effect
model       = "twosteps",        # two-step GMM estimation
transformation = "ld",          # first-differenced transformation
collapse    = TRUE              # instrument collapsing for efficiency
)

# Summarize model output
summary(gmm_model)

```

Table 13) (GMM) Coefficients

Variable	Estimate	Std. Error	z-value	p-value	Significance
lag(PCPIPCH, 1)	0.2741	0.0936	2.93	0.0034	**
GGSB_NPGDP	-3.2433	1.4750	-2.20	0.0279	*
GGXWDG_NGDP	0.2870	0.1456	1.97	0.0487	*
PPPPC	-0.00022	0.00019	-1.13	0.2593	
TX_RPCH	0.2467	0.3035	0.81	0.4162	
TM_RPCH	0.5278	0.2599	2.03	0.0423	*
NID_NGDP	-0.5358	0.3329	-1.61	0.1075	

Interpretation of Coefficients:

lag(PCPIPCH, 1): A significant positive coefficient (0.2741, $p = 0.0034$) confirms inflation persistence, meaning past inflation significantly predicts current inflation—justifying the dynamic specification.

GGSB_NPGDP (Structural Fiscal Balance): A significant negative effect (-3.24 , $p = 0.0279$) implies that fiscal tightening is associated with reduced inflation, supporting counter-cyclical fiscal policy effects.

GGXWDG_NGDP (Government Debt): Positive and significant (0.287, $p = 0.0487$), indicating that higher debt burdens may elevate inflation, possibly through expectations or monetization risk.

PPPPC (Per Capita Income): Not significant, suggesting that income level (after controlling for dynamics and endogeneity) is not a direct driver of short-term inflation.

TX_RPCH (Export Volume): Statistically insignificant, implying exports do not meaningfully influence inflation rates in the short term.

TM_RPCH (Import Volume): Significant positive effect (0.528, $p = 0.0423$), consistent with imported inflation dynamics.

NID_NGDP (Investment): Negative but not significant, suggesting weak inflationary effect from investment in this dynamic setting.

Diagnostic Tests for the GMM Model:-

Wald Test for Joint Significance of Coefficients: To test the joint significance of all model coefficients. This confirms whether the explanatory variables collectively have explanatory power.

Hypotheses: H_0 : All slope coefficients are jointly zero. H_1 : At least one coefficient is non-zero.

Test Statistic:

$$W = \widehat{\beta}^\top \left(\text{Var}(\widehat{\beta}) \right)^{-1} \widehat{\beta}$$

Results: $\chi^2(7) = 660.07$, $p\text{-value} < 0.001$

Decision: Reject H_0 significant cross-sectional dependence exists.

Conclusion: The set of regressors is jointly significant, confirming that the model provides meaningful information about inflation dynamics.

Sargan Test for Overidentifying Restrictions : To assess the validity of instruments used in the GMM estimation. The test checks whether the instruments are exogenous and uncorrelated with the error term. Hypotheses:

H_0 : All instruments are valid (uncorrelated with residuals).

H_1 : At least one instrument is invalid (correlated with residuals).

Test Statistic:

$$J = u^\top Z(Z^\top WZ)^{-1}Z^\top u$$

Where: u = residual vector. Z = matrix of instruments. W = weighting matrix

```
# (i) Hansen/Sargan test for overid  
sargan(gmm_model)
```

Results: $\chi^2 = 11.86$, p-value = 0.457

Decision: Fail to reject H_0 , Instruments are valid.

Conclusion: The instruments used in the GMM estimation are valid. The model is not overfitted with invalid instruments, supporting the consistency of the GMM estimator.

Arellano–Bond Test for First-Order Autocorrelation: To test for first-order serial and second-order serial correlation in the differenced residuals. First-order correlation is expected in differences. The absence of AR(2) is crucial for the validity of instruments in Difference GMM.

Hypotheses: H_0 : No serial correlation. H_1 : serial correlation.

Test Statistic:

$$z = \frac{\widehat{m}_1}{\sqrt{\text{Var}(\widehat{m}_1)}}$$

Where: \widehat{m}_1 : Sample moment. $\text{Var}(\widehat{m}_1)$: Estimated variance of the moment.

```
# (ii) Arellano–Bond tests for autocorrelation in residuals  
# AR(1)  
print(mtest(gmm_model, order = 1))  
# AR(2)  
print(mtest(gmm_model, order = 2))
```

Results [AR (1)]: $z = -0.99$, p-value = 0.3199

Results [AR (2)]: $z = -0.58$, p-value = 0.5607

Decision Rule: Fail to reject H_0 , No serial correlation

Conclusion: The test does not detect significant first-order serial correlation, which is unexpected (since AR (1) is often mechanically present after differencing). However, the result does not invalidate the model as the AR (2) result is more critical for GMM consistency. There is no evidence of second-order serial correlation. This result validates the use of lagged instruments and confirms the internal consistency of the GMM estimator.

4.4. Conclusion: Model Comparison and Final Selection

After estimating and evaluating four model specifications—Pooled OLS, Fixed Effects, Random Effects, and Two-Step Difference GMM—we summarize the findings:

- Pooled OLS failed to account for country heterogeneity and produced biased estimates.
- Fixed Effects controlled unobserved heterogeneity but was inconsistent under endogeneity.
- Random Effects passed the Hausman test and provided a consistent structure but ignored dynamics and endogeneity.
- Difference GMM successfully addressed autocorrelation, heteroskedasticity, and endogeneity, and passed all diagnostic checks (Sargan, AR(1), AR(2)).

In Final, the Two-Step Difference GMM model proves to be statistically valid, having passed all key diagnostic tests—including the Sargan test for instrument validity and the Arellano–Bond test for second-order autocorrelation (AR(2)). It emerges as the most appropriate and robust framework for modeling inflation dynamics across countries over time. By effectively addressing endogeneity, serial correlation, cross-sectional dependence, and inflation persistence, this model delivers estimates that are both theoretically consistent and empirically reliable, making it the preferred specification for macroeconomic panel analysis in this context.

5. Recommendations and Future Work

Recommendations:

Based on the findings, it is recommended that applied macroeconomic research and policy forecasting in inflation contexts prioritize dynamic panel estimators—particularly Two-Step Difference GMM—when dealing with persistent variables and potential endogeneity. Researchers should also incorporate robust standard errors (e.g., Driscoll–Kraay) to account for common violations in macro-panel structures, such as heteroskedasticity and cross-sectional dependence. Emphasis should be placed on careful variable selection and pre-estimation diagnostics to preserve model consistency and efficiency.

Future Work:

- **Expansion of Dataset Coverage:** Increasing the number of countries, expanding the historical range, or adding more frequent observations (e.g., quarterly data) can enhance the richness of the panel, improve estimator efficiency, and allow for sub-sample analysis by region or development level.
- **Extended Diagnostic Testing:** Future analyses should apply more advanced diagnostic tools, such as tests for structural breaks, cross-sectional heteroskedasticity heterogeneity, and interactive fixed effects, to further validate the assumptions and robustness of panel model specifications.
- **Inclusion of Energy Price Indices:** Incorporating a panel variable that reflects energy price fluctuations (e.g., Brent crude or energy import indices) would allow for better modeling of supply-side inflation shocks, especially in energy-dependent or importing economies.
- **Integration of Machine Learning (ML) and Deep Learning (DL):** Future research can explore hybrid approaches by integrating ML and DL techniques—such as LSTM networks or panel-aware boosted trees—to capture complex, nonlinear patterns in panel inflation data while maintaining the structure and advantages of panel frameworks.

Appendix (Codes)

(B): SQL Codes:-

Dataset: https://github.com/1145267383/Panal_Data_Inflation/tree/main/02-Dataset

Database: https://github.com/1145267383/Panal_Data_Inflation/blob/main/03-Clean_Organize_EDA/04-SQL.ipynb

(B): Python Codes:-

Clean and organize and EDA:

https://github.com/1145267383/Panal_Data_Inflation/tree/main/03-Clean_Organize_EDA

Descriptive and Correlation Analysis:

https://github.com/1145267383/Panal_Data_Inflation/tree/main/04-Descriptive_Correlation_Analysis

Models: https://github.com/1145267383/Panal_Data_Inflation/blob/main/05-Models_Panel_Data/01-Models_Python.ipynb

(C): R Codes:-

Models: https://github.com/1145267383/Panal_Data_Inflation/blob/main/05-Models_Panel_Data/02-Models_R.ipynb

References

- Aghion, P., Bergeaud, A., Lequien, M., & Melitz, M. (2021). R&D and productivity growth: Addressing endogeneity and cross-sectional dependence with bias-corrected GMM. *Journal of Econometrics*, 220(2), 413–431.
- Ahn, S. C., Lee, Y. J., & Schmidt, P. (2019). Jackknife bias reduction for dynamic panels with persistent data. *Econometric Reviews*, 38(8), 849–872.
- Bai, J., Liao, Y., & Shi, S. (2020). Robust dynamic panel estimation with latent factors. *Econometric Reviews*, 39(1), 1–27.
- Baltagi, B. H. (2008). *Econometric Analysis of Panel Data* (4th ed.). John Wiley & Sons.
- Barro, R. J., & Sala-i-Martin, X. (2004). *Economic Growth*. MIT Press.
- Becker, S. O., Fetzer, T., & Novy, D. (2010). Who voted for Brexit? A factor-augmented panel analysis of policy preferences. *Economic Journal*, 130(628), F272–F297.
- Blundell, R., & Bond, S. (1998). Initial conditions and moment restrictions in dynamic panel data models. *Journal of Econometrics*, 87(1), 115–143.
- Breitung, J., & Das, S. (2013). Panel cointegration for emerging markets’ inflation–output nexus. *Journal of Applied Econometrics*, 28(2), 374–394.
- Chudik, A., & Pesaran, M. H. (2018). Common correlated effects in large panels. *Econometric Reviews*, 37(3), 405–445.
- Ciccone, A. (2015). Financial crises and growth: A system GMM panel analysis. *Review of Economics and Statistics*, 97(2), 282–297.
- Driscoll, J. C., & Kraay, A. C. (2013). Consistent covariance matrix estimation with spatially dependent panel data. *Review of Economics and Statistics*, 95(4), 1144–1155.
- Fischer, S. (1993). The role of macroeconomic factors in growth. *Journal of Monetary Economics*, 32(3), 485–512.
- Hausman, J. A. (1978). Specification tests in econometrics. *Econometrica*, 46(6), 1251–1271.

- Hsiao, C. (2003). *Analysis of Panel Data*. Cambridge University Press.
- Huang, W.-Q., & Pesaran, M. H. (2022). Spatially augmented interactive effects models for panel data. *Journal of Econometrics*, 229(1), 112–131.
- Levin, A., Lin, C.-F., & Chu, C.-S. J. (2002). Unit root tests in panel data: Asymptotic and finite-sample properties. *Journal of Econometrics*, 108(1), 1–24.
- Moulton, B. R. (1990). An illustration of a pitfall in estimating the effects of aggregate variables on micro units. *Review of Economics and Statistics*, 72(2), 334–338.
- Moon, H. R., & Weidner, M. (2015). Identifying global dynamics in large panels with interactive fixed effects. *Econometric Theory*, 31(6), 838–863.
- Mundlak, Y. (1978). On the pooling of time series and cross section data. *Econometrica*, 46(1), 69–85.
- Nickell, S. (1997). Unemployment and labor market rigidities: Europe versus North America. *Journal of Economic Perspectives*, 11(3), 55–74.
- Pesaran, M. H. (2004). General diagnostic tests for cross-sectional dependence in panels. *CESifo Working Paper* No. 1229.
- Sun, Y., & Kim, T.-H. (2023). High-dimensional instrument selection for dynamic panels via LASSO. *Journal of Business & Economic Statistics*, 41(4), 1072–1087.
- Verbeek, M., & Nijman, T. (1992). Testing for selectivity bias in panel data models. *International Economic Review*, 33(3), 681–703.
- Wooldridge, J. M. (2010). *Econometric Analysis of Cross Section and Panel Data* (2nd ed.). MIT Press.
- Zhang, Y., & Lee, C.-J. (2024). Machine learning augmented GMM for panel data with structural breaks. *Journal of Applied Econometrics*, 39(2), 255–274.