# Testing whether group-level fixed effects are sufficient in panel data models

David Vincent

dvincent@dveconometrics.co.uk

September 2025

2025 UK Stata Conference

# Contents

## Introduction

- The fixed effects (FE) estimator is one of the most widely used methods for estimating coefficients on time varying variables in linear panel data models.

- Its main appeal is that it controls for correlation between the regressors and unobserved unit specific effects that would otherwise lead to omitted variable bias.

- The FE estimator works by applying OLS to the within transformed data, where all of the variables are expressed as deviations from their unit specific means.

- This transformation removes between-unit (i.e. cross sectional), variation in the data, such that coefficients are estimated solely from within unit variation.

## Limitations of Fixed Effects

- A well known limitation of FE is that when key variables of interest display little within-unit variation, the resulting estimates can be very imprecise.

- One approach is to turn to the random effects estimator, and decide whether this is valid by testing if the explanatory variables are uncorrelated with the unit-specific effects.

- The Hausman test is the conventional tool here, and as Wooldridge (2019) discusses, fully robust versions can be obtained from auxiliary regressions.

- Another possibility is to ask whether controlling for heterogeneity at a higher level of aggregation (e.g. group effects) might suffice to ensure regressors exogeneity.

## Papke and Wooldridge (2023)

- If regressors are only correlated with group effects, both unit and group FE are consistent, but group FE may be more efficient. If regressors are correlated with unit effects, only unit FE is consistent and the estimates will differ.

- This is examined by Papke and Wooldridge (2023), who propose a variable-addition test: estimate the model by group FE, augmented with unit-means of the time-varying regressors, and test the null that all coefficients on the means are zero.

- This is the Mundlak (1978) regression, where the coefficients on the means are differences in the between-effects (controlling for group dummies) and the unit FE slopes, while the joint test is for equality of unit and group FE.

## Papke and Wooldridge (2023)

- Although not discussed, the joint test is equivalent to a C-test of the orthogonality restrictions that justify group FE.

- Papke and Wooldridge (2023) also also propose a conventional Hausman test, which directly compares unit and group FE estimates for one regressors, and show how to obtain robust standard errors for the difference.

- While this allows one to focus on a specific coefficient of interest, the underlying estimators are unit and group FE.

- As such, the null can be rejected even when the variable of interest is uncorrelated with the unit-specific effect, provided it is correlated with other endogenous regressors.

## Motivation for xtfelevel

- I present xtfelevel, which builds on Papke and Wooldridge (2023) by comparing unit and group fixed effects methods in a more general IV framework.

- Instead of the all-or-nothing approach, the test allows some time varying regressors to be treated as always endogenous in both the consistent and efficient estimators.

- If the null is not rejected, the IV estimator can be substantially more efficient than full unit FE, and is equivalent to group FE estimation of the Mundlak regression, where unit-means of exogenous variables are partialled out from unit means of endogenous regressors.

- This differs from Papke and Wooldridge (2023) hybrid estimator, which drops insignificant means from the regression and may introduce bias when regressors are correlated.

## The Model

- To motivate the set-up, I use an instrumental variables (IV) framework that nests the standard unit and group fixed effects estimators as special cases.

- Let $y_{it}$ denote the outcome for individual $i = 1, \ldots, N$ observed in period $t = 1, \ldots, T_i$, and let $x_{1it}$ be a vector of time-varying regressors whose effects we are interested in.

- Individuals are assigned into $G$ group, with group membership indicated by $g(i)$. The model to be estimated is:

$$y_{it} = x_{1it}^{'}\beta_1 + x_{2it}^{'}\beta_2 + h_{g(i)} + u_i + \epsilon_{it}, \tag{1}$$

where $u_i$ is a unit-specific effect, $h_{g(i)}$ a group-specific effect, and $\epsilon_{it}$ an idiosyncratic error.

## The Model

- The regressors $x_{1it}$ are strictly exogenous with respect to $\epsilon_{it}$, and are divided into two sets: $x_{1Bit}$, which are always assumed to be correlated with $u_i$, and $x_{1Ait}$, which may or may not be correlated with $u_i$.

- The variables in $x_{2it}$ are typically time-invariant controls, but may also regressors that vary over time and are uncorrelated with $u_i$, such as macro-economic factors.

- All variables, including $x_{2it}$ are allowed to be arbitrarily correlated with the group level heterogeneity $h_{g(i)}$ and the main focus on the and estimation of $\beta_1$.

- Note that even if $x_{1A}$ is uncorrelated with $u_i$, the IV estimator that treats $x_{1A}$ as exogenous may still differ from the one that treats it as endogenous when selection into the sample depends on both $u_i$ and $x_{1A}$.

## The Model

- It is helpful to consider the group-specific fixed effects as actual parameters to be estimated, which means that instruments are only needed for $x_1$.

- Let $h_{g(i)} = d_i^{'}\lambda$ where $d_i = \left[ d1_{g(i)}...dG_{g(i)} \right]^{'}$, is a vector of group dummy variables, such that $df_{g(i)} = 1$ if $f = g$. Then stacking (1) over all periods for individual $i$ gives:

$$y_i = X_{1i}^{'}\beta_1 + X_{2i}^{'}\beta_2 + D_i\lambda + e_i u_i + \epsilon_i \qquad (2)$$

where $e_i$ is a $T_i \times 1$ vector of 1's, and $D_i = e_i d_i^{'}$ is a $T_i \times G$ matrix of observations on the dummy variables.

## Hausman Test

- I consider two IV estimators that both treat the group-specific effects as fixed effects, but differ in how they control for correlation between $x_{1it}$ and the unit-heterogeneity $u_i$.

- The first, which I shall refer to as the **consistent** estimator $\widehat{\beta}_1^c$ uses the unit-demeaned variables $x_{1it} - \bar{x}_{1i}$ as instruments and is consistent when $x_{1it}$ is correlated with $u_i$ or $h_{g(i)}$.

- The second $\widehat{\beta}_1^e$ again uses $x_{1Bit} - \bar{x}_{1Bi}$ as instruments for $x_{1Bit}$, but now uses $x_{1Ait}$ as instruments for itself and will therefore be inconsistent when $x_{1it}$ is correlated with $u_i$.

- When $x_{1Ait}$ is uncorrelated with $u_i$, both estimators will have the same plim, but $\widehat{\beta}_1^e$ will tend to be more **efficient**, as the unit de-meaning employed by $\widehat{\beta}_1^c$ removes more of the variation in the data than is necessary.

## Hausman Test

- This suggests testing the joint endogeneity of $x_{1Ait}$ by testing for differences between the $\widehat{\beta}_1^c$ and $\widehat{\beta}_1^e$ estimators:

$$
\begin{aligned}
H_0 : \ \text{plim}\,\widehat{\beta}_1^c - \text{plim}\,\widehat{\beta}_1^e &= 0 \\
H_1 : \ \text{plim}\,\widehat{\beta}_1^c - \text{plim}\,\widehat{\beta}_1^e &\neq 0
\end{aligned}
\tag{3}
$$

- Under the null, the difference between the consistent and efficient IV estimators is also root-N consistent, with a mean 0 and a limit normal distribution, so that:

$$
\sqrt{N}\left(\widehat{\beta}_1^c - \widehat{\beta}_1^e\right) \longrightarrow N\left[0, V\right]
$$

- A Hausman test statistic is then:

$$
H = \left(\widehat{\beta}_1^c - \widehat{\beta}_1^e\right)^{'}\left(\widehat{V}/N\right)^{-1}\left(\widehat{\beta}_1^c - \widehat{\beta}_1^e\right) \xrightarrow{d} \chi^2_{K_{1A}}
$$

## Hausman Test

- Even though $x_1$ consists of $K_1$ variables, as only $x_{1A}$ regressors are being tested, $\widehat{V}$ will be of reduced rank, equal to $K_{1A}$.

- The Hausman test can also be applied to to test for differences in a single coefficient for any of the parameters in $x_1$:

$$H_k = \left[ \frac{\left( \widehat{\beta}_{1k}^c - \widehat{\beta}_{1k}^e \right)}{SE_k} \right]^2 \xrightarrow{d} \chi_1^2$$

where $SE_k = \sqrt{\widehat{V}[k,k]/N}$ is the standard error of the $k^{th}$ component of the difference. This could lead to a different conclusion from a test of all the parameters in $\beta_1$.

- The Hausman test is therefore best interpreted as a test of the consequences of different estimators on the same equation.

## Valid Instruments

- The consistent estimator treats all variables as potentially correlated with $h_{g(i)}$, and $x_1$ as correlated with the $u_i$.

- Let $Q_i^u = I_{T_i} - M_i^u$ denote the within transformation matrix, which subtracts unit level means $M_i^u = T_i^{-1} e_i e_i^{'}$, then valid instruments for $X_{1i}$ are $Q_i^u X_{1i}^{'}$, as then:

$$E\left[\left(Q_i^u X_{1i}\right)^{'}(e_i u_i + \epsilon_i)\right] = E\left[X_{1i}^{'} Q_i^u \epsilon_i\right] = 0$$

- Stacking over all individuals, the instrument set is:

$$W_2 = [Q^u X, X_2, D] \tag{4}$$

where $Q^u = \text{diag}\{Q_1^u ... Q_N^u\}$ and the within transformed instruments also include the exogenous controls.

## Consistent Estimation

- Basu (2023) shows that for IV estimators, exogenous variables can be partialled out of the instrument set and the model.

- For the set of $G$ group dummies, $D$, this simply requires applying the within transformation at the group level.

- Letting $Q = \text{diag}\{Q_1...Q_G\}$ denote the equivalent transformation matrix by group, then the within transformed model subtracting group means is:

$$\tilde{y} = \tilde{X}_1'\beta_1 + \tilde{X}_2'\beta_2 + \tilde{v} \tag{5}$$

where $\tilde{y} = Qy$, $\tilde{X} = QX$ and the error $\tilde{v} = Qv$, where for individual $i$, this is defined above as $v_i = e_i u_i + \epsilon_i$.

## Consistent Estimation

- As units are nested within groups, the group mean of the already unit-demeaned data is zero, hence $Q \times Q^u = Q^u$. Thus, after partialling out $D$ from the instruments:

$$W_2 = \left[ Q^u X, \tilde{X}_2 \right] \tag{6}$$

- Let $P_2 = W_2(W_2'W_2)^{-1}W_2'$ denote the projection matrix, and $M_{\tilde{X}_2}$ the usual residual maker, then by the FWL theorem, the IV estimator that treats all variables in $x_1$ as exogenous is:

$$\widehat{\beta_1^c} = (\tilde{X}_1'P_2 M_{\tilde{X}_2} P_2 \tilde{X}_1)^{-1}\tilde{X}_1 P_2 M_{\tilde{X}_2} \tilde{y} = A_c^{-1}\tilde{X}_1 P_2 M_{\tilde{X}_2} \tilde{y} \tag{7}$$

## Consistent Estimation

- It is helpful to note that the predicted values $P_2\tilde{X}_1$ can be decomposed into within and between-unit components:

$$P_2\tilde{X}_1 \ = \ Q^u X_1 \ + \ P_{M^u\tilde{X}_2} M^u \tilde{X}_1. \qquad (8)$$

- The first term is the within-unit variation in $X_1$, whereas the second term, are the predictions from separate between-effects regressions of $x_{1kit}$ on $x_{2it}$, including $G$ group dummies

- This is equivalent to running OLS on the unit-level means of these variables, with $G$ group dummies included, where each mean is replicated $T_i$ times in the dataset.

# Consistent Estimation: Unit Fixed Effects

- This shows that the variation used to estimate $\beta_1$ includes not only the within-unit variation in $X_1$, but also the between-unit variation that is predicted by the exogenous controls $X_2$.

- In most applications, the controls $X_2$ are time-invariant, and when this is the case, the additional between variation makes no contribution, and is eliminated by $M_{\tilde{X}_2}$.

- Also, as $M_{\tilde{X}_2} Q^u = Q^u$ this IV estimator that treats all $x_1$ as endogenous, reduces to standard unit fixed effects:

$$\widehat{\beta}_1^c = \widehat{\beta}_1^{FE(i)} \left( X_1' Q^u X_1 \right)^{-1} X_1' Q^u y.$$

## Efficient Estimation

- For the more efficient IV estimator, $X_{1A}$ is now treated as exogenous, with only $X_{1B}$ as endogenous, hence the instrument set becomes:

$$W_1 = \left[ Q^u X, \tilde{X}_{1A}, \tilde{X}_2 \right] \quad (9)$$

- Defining $P_1$ as the projection onto $W_1$ yields:

$$\widehat{\beta}_1^e = \left( \tilde{X}_1' P_1 M_{\tilde{X}_2} P_1 \tilde{X}_1 \right)^{-1} \tilde{X}_1' P_1 M_{\tilde{X}_2} \tilde{y} = A_e^{-1} \tilde{X}_1' P_1 M_{\tilde{X}_2} \tilde{y} \quad (10)$$

- When all variables are tested, $\tilde{X}_{1A} = \tilde{X}_1$, and the estimator reduces to group fixed effects:

$$\widehat{\beta}_1^e = \widehat{\beta}_1^{FE(g)} = \left( \tilde{X}_1' M_{\tilde{X}_2} \tilde{X}_1 \right)^{-1} \tilde{X}_1' M_{\tilde{X}_2} \tilde{y}$$

## Limiting Distribution

- The final component required to implement he Hausman test, is an consistent estimate of $V = \text{AsyVar}\left(\sqrt{N}(\widehat{\beta}_1^c - \widehat{\beta}_1^e)\right)$

- As $W_2 \subset W_1$, then $P_2 P_1 = P_2$, and the difference between the consistent and efficient estimators in (7) and (10) is:

$$\widehat{\beta}_1^c - \widehat{\beta}_1^e = A_c^{-1} \tilde{X}_1' P_2 M_{\tilde{X}_2} M^* \tilde{y} \qquad (11)$$

- The residual maker $M^*$ is defined as:

$$M^* = I - M_{\tilde{X}_2} P_1 \tilde{X}_1 A_e^{-1} \tilde{X}_1' P_1 M_{\tilde{X}_2}$$

# Limiting Distribution

- Given that all variables are group de-meaned, and $QQ = Q$, replacing $\tilde{y}$ with $v$ in (11) leads to the same result:

$$\sqrt{N}\left(\widehat{\beta}_1^c - \widehat{\beta}_1^e\right) = \left(\frac{1}{N}A_c\right)^{-1}\frac{1}{\sqrt{N}}\tilde{X}_1' P_2 M_{\tilde{X}_2} M^* v$$

- Letting $\Omega = E\left[vv' \mid W\right]$, then under $H_0$, the covariance matrix of the limiting distribution is:

$$V = \left(\operatorname{plim}\frac{1}{N}A_c\right)^{-1} S \left(\operatorname{plim}\frac{1}{N}A_c\right)^{-1} \tag{12}$$

where the middle term:

$$S = \lim_{N \to \infty} \frac{1}{N} E\left[\tilde{X}_1' P_2 M_{\tilde{X}_2} M^* \Omega M^* M_{\tilde{X}_2} P_2 \tilde{X}_1\right] \tag{13}$$

## Variance Estimation: IID Errors

- Define the estimator of $S$:

$$\widehat{S} = \frac{1}{N} \tilde{X}_1' P_2 M_{\tilde{X}_2} M^* \widehat{\Omega} M^* M_{\tilde{X}_2} P_2 \tilde{X}_1 \tag{14}$$

- Different choices of $\widehat{\Omega}$ lead to different estimators of $V$ in (12). Assuming errors are independent across units:

$$\widehat{\Omega} = \text{diag}(\widehat{\Omega}_1, \ldots, \widehat{\Omega}_N).$$

- Under the variance components model $u_i \sim \text{iid}(0, \sigma_u^2)$, $\epsilon_{it} \sim \text{iid}(0, \sigma_\epsilon^2)$:

$$\widehat{\Omega}_i = \widehat{\sigma}_\epsilon^2 I_{T_i} + \widehat{\sigma}_u^2 e_i e_i'.$$

- Estimates of $\sigma_\epsilon^2$ and $\sigma_u^2$ are obtained from residuals of the group-demeaned model, following xthtaylor.

# Variance Estimation: Cluster-Robust

- An alternative is to use an estimator that allows for clustering by unit or group and is robust to arbitrary serial correlation of errors within clusters and heteroskedasticty across clusters.

**Clustering at the unit level**

- Let $\widehat{v}_i$ be the $T_i \times 1$ vector of residuals for unit $i$. Then:

$$\widehat{\Omega}_i = \widehat{v}_i \widehat{v}_i'.$$

**Clustering at the group level**

- Let $\widehat{v}_g$ be the stacked residuals for all units in group $g$. Then:

$$\widehat{\Omega}_g = \widehat{v}_g \widehat{v}_g'.$$

## Auxiliary Regressions

- Following Davidson and MacKinnon (1993), the Hausman test can also be implemented using auxiliary regressions.

- When testing all variables in $x_1$ for endogeneity, this requires augmenting the original model with residuals from regressions of each $x_1$ variable on the instruments in $W_2$ together with the group dummies. The model is then estimated by pooled OLS

- As before, it is again easier to work with the model where the group dummies have been partialled out from the variables

$$\tilde{y} = \tilde{X}_1' \beta_1 + \tilde{X}_2' \beta_2 + M_2 \tilde{X}_1 \delta_1 + \tilde{\eta}. \qquad (15)$$

## Auxiliary Regressions

- We are interested in $\widehat{\delta}_1$. Applying the FWL theorem:

$$\widehat{\delta}_1 = -\left(\tilde{X}_1' P_2 M_{\tilde{X}} P_2 \tilde{X}_1\right)^{-1} \tilde{X}_1' P_2 M_{\tilde{X}_2} M^* \tilde{y}$$

- The second term is exactly what appears in equation (11) for the vector of contrasts $\widehat{\beta}_1^c - \widehat{\beta}_1^e$, hence by substitution:

$$\widehat{\delta}_1 = -A_\delta^{-1} A_c \left(\widehat{\beta}_1^c - \widehat{\beta}_1^e\right) \tag{16}$$

- As the $A$-terms in (16) cancel in the Wald statistic, testing $\delta_1 = 0$ is equivalent to testing $\text{plim}(\widehat{\beta}_1^c - \widehat{\beta}_1^e) = 0$. When the exogenous controls $x_2$ are all time-invariant, the consistent and efficient estimators reduce to unit and group FE.

## Auxiliary Regressions

- The auxiliary regression in (15) provides the control function representation of the IV estimator, and in this set-up the OLS estimator is unit-FE: $\widehat{\beta}_1 = \widehat{\beta}_1^{FE(i)}$.

- It can also be shown that $\widehat{\delta}_1$ equals the difference between a between-effects regression (including $G$ group dummies) and the unit fixed effects estimator:

$$\widehat{\delta}_1 = \widehat{\beta}_1^{BE} - \widehat{\beta}_1^{FE(i)}. \tag{17}$$

- Thus, while the joint test is the same as testing the equality of the full vector of unit and group FE contrasts, the $t$-tests on individual $\delta_{1k}$ components, compare BE with unit-FE.

- The comparison provided by xtfelevel is for unit-group fixed effects, and may therefore lead to different conclusions.

## Mundlak Regression

- It is easy to show that the control function regression in (15) is also the Mundlak device. The residual maker is:

$$M_2 \tilde{X}_1 = \tilde{X}_1 - P_2 \tilde{X}_1$$

- And from (8), we know that

$$P_2 \tilde{X}_1 = Q^u \tilde{X}_1 + P_{M^u \tilde{X}_2} M^u \tilde{X}_1$$

  hence by substitution:

$$M_2 \tilde{X}_1 = \left( I - P_{M^u \tilde{X}_2} \right) M^u \tilde{X}_1 \tag{18}$$

- These are the residuals from separate between-effects regressions of $\bar{x}_{1ki}$ on $\bar{x}_{2i}$ including $G$ group dummies, with each mean replicated $T_i$ times in the data.

## Mundlak Regression

- As the residuals are time-invariant for a given individual, for observation $it$, the residual for the $k$-th variable is

$$\left(M_2\tilde{X}_1\right)_{itk} = \bar{x}_{1ki} - \bar{x}'_{2i}\widehat{\pi}^{BE}_k - d'_i\widehat{\theta}_k$$

- Plugging these into the auxiliary model for $y_{it}$ and letting $d'_i\widehat{\theta}_k$ be subsumed into the overall group-effect yields:

$$y_{it} = x'_{1it}\beta_1 + x'_{2it}\beta_2 + \sum_{k=1}^{K_1}\left(\bar{x}_{1ki} - \bar{x}'_{2i}\widehat{\pi}^{BE}_k\right)\delta_k + d'_i\lambda + \eta_{it} \quad (19)$$

- The model can be estimated by OLS including $G$ group dummies, or by group fixed effects, which is OLS applied to the group within transformed model.

## Mundlak Equivalent of Efficient IV Estimator

- When all variables in $x_1$ are treated as endogenous, and the variables in $x_2$ are all time-invariant, there is no need to partial out their means (unless they are of direct interest), since their correlation with $x_1$ is absorbed in $\beta_2$.

- The situation is different once some regressors are treated as exogenous. Consider the efficient IV estimator where $x_{1A}$ is an exogenous policy variable, and $x_{1B}$ is a single endogenous regressor (assume $x_2$ is empty).

- As illustrated in (19), the IV estimator is equivalent to a Mundlak regression in which the mean of $\bar{x}_{1A}$ is partialled out from the mean of the endogenous regressor $\bar{x}_{1B}$:

$$y_{it} = x_{1it}^{'}\beta_1 + x_{2it}^{'}\beta_2 + \left(\bar{x}_{1Bi} - \bar{x}_{1Ai}^{'}\widehat{\pi}^{BE}\right)\delta_{1B} + d_i^{'}\lambda + \eta_{it}$$

## Hybrid Estimator and Potential Bias

- The IV estimator therefore imposes a restriction on the implied coefficient of $\bar{x}_{1A}$ in the usual Mundlak model:

$$\delta_{1A} = -\delta_{1B}\pi^{BE}. \qquad (20)$$

- Papke and Wooldridge (2023) suggest that if the unrestricted estimate $\widehat{\delta}_{1A}$ in the usual Mundlak regression is not significant, it can be dropped, and this hybrid estimator may yield more precise estimates of $\beta_{1A}$ by reducing collinearity with $x_{1Ait}$.

- But if $x_{1A}$ has little within-unit variation, insignificance of $\bar{x}_{1A}$ may simply reflect low power rather than true exogeneity, and dropping it risks omitted variable bias.

## Hybrid Estimator and Potential Bias

- It can be shown that the hybrid estimator converges to a weighted average of the FE and between-effects estimators, where the weights depend on the share of within and between variation in the regressors.

- Since the between-effects estimator is inconsistent when $x_{1B}$ is endogenous (because $u_i$ is in the error term), the resulting mixture (i.e., the hybrid) will also be inconsistent.

- Thus, unless the excluded means are uncorrelated with the means that remain in the model (i.e. all $\pi_k^{BE} = 0$), omitting them rather than partialling them out risks potential omitted variables bias.

## xtfelevel

Robust Hausman tests for choosing the level of fixed effects in linear panel data models

<u>xtfelevel</u> *depvar varlist* $\big[$ *if* $\big]$ $\big[$ *in* $\big]$ *group(varlist)* $\big[$ tvendog(*varlist*) tvexog(*varlist*) sigmamless sigmamore vce(*vcetype*) $\big]$

---

group(*varlist*) variables that define groups that nest units.

tvendog(*varlist*) time-varying regressors always treated as correlated with the unit-specific effect.

tvexog(*varlist*) time-varying regressors always treated as uncorrelated with the unit-specific effect.

sigmamless use residuals from the consistent estimator to estimate the covariance matrix for the Hausman test.

sigmamore use residuals from the efficient estimator to estimate the covariance matrix for the Hausman test.

## xtfelevel: Defining Groups

- Groups do not have to be defined by a single identifier. They can be formed additively from several components, as long as units are nested within groups.

- Example (automotive data):

    - **Unit:** model–version–engine size–transmission–AC.

    - **Group (broad):** model–version only.

    - **Group (additive):** model–version + engine size + transmission + AC.

- xtfelevel implements both cases. The additive case uses the method of alternating projections

## Simulated Data

- Data is simulated from the following model, setting $G = 100$ groups, each containing 10 units observed over $T = 5$ periods.

$$
\begin{aligned}
y_{it} &= 1 + x_{1Ait} + x_{1Bit} + x_{2i} + u_i + h_{g(i)} + \epsilon_{it} \\
u_i &\sim N(0, 1) \\
\epsilon_{it} &\sim N(0, 1)
\end{aligned}
$$

- $x_{1Bit}$ is endogenous and correlated with $u_i$; whereas $x_{1Ait}$ is exogenous but correlated with $x_{1Bit}$.

- The policy variable of interest is $x_{1Ait}$, which shows little within-unit variation (time-varying in only 20% of units).

- Group fixed effects estimator is inconsistent for $\beta_{x_{1A}}$, because $x_{1Ait}$ is correlated with $x_{1Bit}$, and $x_{1Bit}$ is endogenous.

# Testing $X_{1A}$ and $X_{1B}$ ($H_0$ : Unit FE = Group FE)

```
. xtfelevel y x1A x1B x2,  group(idg) vce(cluster idg)
                                         Number of obs     =        5,000
Unit variable:   idu                     Number of units   =        1,000
Group variable:  idg                     Number of groups  =          100
```

|  | Coef. | Robust Std. Err. | z | P>\|z\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| **Efficient** | | | | | | |
| x1A | 1.459548 | .02643 | 55.22 | 0.000 | 1.407746 | 1.51135 |
| x1B | 1.470631 | .0138896 | 105.88 | 0.000 | 1.443408 | 1.497855 |
| x2 | .9712746 | .0221096 | 43.93 | 0.000 | .9279406 | 1.014609 |
| **Consistent** | | | | | | |
| x1A | .6029016 | .3161743 | 1.91 | 0.057 | -.0167886 | 1.222592 |
| x1B | 1.015079 | .0160868 | 63.10 | 0.000 | .9835499 | 1.046609 |
| x2 | .9629337 | .0361778 | 26.62 | 0.000 | .8920266 | 1.033841 |
| **Diff_Tested** | | | | | | |
| x1A | -.8566462 | .3622874 | -2.36 | 0.018 | -1.566717 | -.1465759 |
| x1B | -.4555519 | .0491084 | -9.28 | 0.000 | -.5518026 | -.3593013 |

```
Endogeneity Test:                                                    86.2027
Variables Tested:  x1A x1B                     Chi-sq(2) P-val   =     0.0000
```

# Testing $X_{1A}$ ($H_0$ : Unit FE $=$ Group FE, $X_{1B}$ Endogenous)

```
. xtfelevel y x1A x1B x2, tvendog(x1B) group(idg) vce(cluster idg)

                                             Number of obs    =      5,000
Unit variable:  idu                          Number of units  =      1,000
Group variable: idg                          Number of groups =        100

                         Robust
          y      Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]

Efficient
        x1A    1.007809   .0355045    28.39   0.000    .9382216    1.077397
        x1B    1.014763   .0161035    63.02   0.000    .9832011    1.046326
         x2    .9662914   .0315063    30.67   0.000    .9045402    1.028043

Consistent
        x1A    .6029016   .3161743     1.91   0.057   -.0167886    1.222592
        x1B    1.015079   .0160868    63.10   0.000    .9835499    1.046609
         x2    .9629337   .0361778    26.62   0.000    .8920266    1.033841

Diff_Tested
        x1A   -.4049075   .3176995    -1.27   0.202   -1.027587    .2177721

Diff_Endog
        x1B    .0003161    .000248     1.27   0.202     -.00017    .0008023

Endogeneity Test:                                                    1.6243
Variables Tested:  x1A                        Chi-sq(1) P-val   =    0.2025
```

# Mundlak Representation: Testing $X_{1A}$ and $X_{1B}$

```
.  *estimate mundlak regression - estimates are unit FE´s
. reghdfe y x1A x1B x2 mui_x1A mui_x1B, a(idg) vce(cluster idg)

HDFE Linear regression                          Number of obs    =       5,000
Absorbing 1 HDFE group                          F(   5,     99)  =     3463.00
Statistics robust to heteroskedasticity         Prob > F         =      0.0000
                                                R-squared        =      0.9096
                                                Adj R-squared    =      0.9077
                                                Within R-sq.     =      0.8264
Number of clusters (idg)     =         100      Root MSE         =      1.0662

                             (Std. Err. adjusted for 100 clusters in idg)
```

| y | Coef. | Robust Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---:|---:|---:|---:|---:|---:|---:|
| x1A | .6029016 | .3162692 | 1.91 | 0.060 | -.0246452 | 1.230448 |
| x1B | 1.015079 | .0160916 | 63.08 | 0.000 | .9831502 | 1.047009 |
| x2 | .9749439 | .020648 | 47.22 | 0.000 | .9339739 | 1.015914 |
| mui_x1A | 1.189611 | .3181731 | 3.74 | 0.000 | .5582867 | 1.820936 |
| mui_x1B | .7901667 | .0242087 | 32.64 | 0.000 | .7421315 | .838202 |
| _cons | 1.040708 | .0058963 | 176.50 | 0.000 | 1.029009 | 1.052408 |

```
. testparm mui_x1A mui_x1B

 ( 1)  mui_x1A = 0
 ( 2)  mui_x1B = 0

       F(  2,    99) =  538.21
            Prob > F =    0.0000
```

# Mundlak Equivalent of Efficient IV Estimator

```
.  *partial out mean of x1A and x2 from x1B
. reghdfe mui_x1B mui_x1A x2, a(idg) resid
. predict mui_x1B_partial, res


.  *estimate adjusted mundlak
. reghdfe y x1A x1B x2 mui_x1B_partial, a(idg) vce(cluster idg)
```

HDFE Linear regression                          Number of obs   =      5,000
Absorbing 1 HDFE group                          F(  4,    99) =    4328.09
Statistics robust to heteroskedasticity         Prob > F        =     0.0000
                                                R-squared       =     0.9096
                                                Adj R-squared   =     0.9077
                                                Within R-sq.    =     0.8264
Number of clusters (idg)      =        100      Root MSE        =     1.0662

                            (Std. Err. adjusted for 100 clusters in idg)

| y | Coef. | Robust Std. Err. | t | P>|t| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| x1A | 1.007809 | .0252821 | 39.86 | 0.000 | .957644 | 1.057974 |
| x1B | 1.014763 | .0161051 | 63.01 | 0.000 | .9828074 | 1.046719 |
| x2 | .9662914 | .0206452 | 46.80 | 0.000 | .9253269 | 1.007256 |
| mui_x1B_partial | .7904828 | .0242593 | 32.58 | 0.000 | .742347 | .8386187 |
| _cons | 1.185564 | .0056486 | 209.89 | 0.000 | 1.174356 | 1.196772 |

# Hybrid Estimator: No Partialling Out (Illustration Only)

```
. reghdfe y x1A x1B x2 mui_x1B, a(idg) vce(cluster idg)

HDFE Linear regression                              Number of obs   =      5,000
Absorbing 1 HDFE group                              F(  4,     99) =    4346.80
Statistics robust to heteroskedasticity             Prob > F        =     0.0000
                                                    R-squared       =     0.9094
                                                    Adj R-squared   =     0.9075
                                                    Within R-sq.    =     0.8260
Number of clusters (idg)      =       100           Root MSE        =     1.0672

                              (Std. Err. adjusted for 100 clusters in idg)
```

|           y |      Coef. | Robust Std. Err. |      t | P>\|t\| | [95% Conf. Interval] |          |
|------------:|-----------:|-----------------:|-------:|--------:|---------------------:|---------:|
|         x1A |   1.788236 |        .0267531  |  66.84 |   0.000 |             1.735152 |  1.84132 |
|         x1B |   1.016747 |        .0159982  |  63.55 |   0.000 |             .9850031 | 1.048491 |
|          x2 |   .9749033 |        .0206466  |  47.22 |   0.000 |              .933936 | 1.01587  |
|    mui_x1B  |   .7865974 |        .0240587  |  32.69 |   0.000 |             .7388598 |  .834335 |
|       _cons |   1.041618 |        .0058695  | 177.46 |   0.000 |             1.029971 | 1.053264 |

## Conclusion

- This presentation has set out a Hausman-type test for deciding whether group-level fixed effects are sufficient, compared to unit fixed effects.

- The approach builds on Papke and Wooldridge (2023) to a general IV framework, allowing some regressors to be treated as endogenous under both the null and the alternative.

- The test can be implemented using the Stata command xtfelevel, which provides robust and cluster-robust covariance estimators.

- The efficient IV estimator is equivalent to a Mundlak form, where the unit-means of exogenous variables are partialled out of the means of the endogenous regressors.

## References

Basu, D. 2023. The Yule-Frisch-Waugh-Lovell theorem for linear instrumental variables estimation. *arXiv preprint arXiv:2307.12731.* .

Davidson, R., and J. G. MacKinnon. 1993. *Estimation and inference in econometrics*, vol. 63. New York: Oxford.

Mundlak, Y. 1978. On the pooling of time series and cross section data. *Econometrica* 46: 69–85.

Papke, L. E., and J. M. Wooldridge. 2023. A simple, robust test for choosing the level of fixed effects in linear panel data models. *Empirical Economics* 64(6): 2683–2701.

Wooldridge, J. 2019. Correlated random effects models with unbalanced panels. *Journal of Econometrics* 211(1): 137–150.