# Section 13 Models for Pooled and Panel Data

## *Data definitions*

- **Pooled data** occur when we have a "time series of cross sections," but the observations in each cross section do not necessarily refer to the same unit.
    - o HGL is ambiguous about this and sometimes use pooled to refer to panel data
- **Panel data** refers to samples of the *same* cross-sectional units observed at multiple points in time. A panel-data observation has two dimensions: $x_{it}$, where $i$ runs from 1 to $N$ and denotes the cross-sectional unit and $t$ runs from 1 to $T$ and denotes the time of the observation.
    - o A **balanced panel** has every observation from 1 to $N$ observable in every period 1 to $T$.
    - o An **unbalanced panel** has missing data.
    - o Panel data commands in Stata start with **xt**, as in xtreg. Be careful about models and default assumptions in these commands.
- **Asymptotic properties**
    - o We need to be careful about large-sample properties of these estimators: are we talking about $N \rightarrow \infty$, about $T \rightarrow \infty$, or both?

## *Regression with pooled cross sections*

- The crucial question with pooled cross sections from different time periods is "Does the same model apply in each time period?"
    - o Has inflation changed the real values of some variables, requiring adjustment?
    - o Was the business cycle at different phases in different periods?
    - o Were there changes in technology or regulation that would cause behavior to be different?
    - o Are there other factors that might cause coefficients in one period to differ from those in others?
- This is a special case of the Assumption #0 question: Do all observations come from the same model?
- **Time dummy variables**
    - o A very general way of modeling (and testing for) differences in intercept terms or slope coefficients between periods is the use of time dummies.
    - o Including time dummies (for all but one, omitted date in the sample to avoid the dummy-variable trap) alone allows the intercept to have a different value in each period.

- The estimated intercept term in the model with time dummies is the estimated intercept in the period with the omitted dummy.
- The estimated coefficient on an included time dummy corresponding to a particular period is an estimate of the difference between the intercept in that period and the intercept in the omitted period.
- A joint test of whether all the dummies' coefficients are zero tests the hypothesis that the intercept does not vary at all over periods.
- The simple test of whether a particular dummy's coefficient is zero tests the hypothesis that the intercept in that dummy's period does not differ from that of the omitted period.
  - o Including interactions between time dummies and another variable $Z$ allows the coefficient on (effect of) $Z$ to vary across periods.
    - As before, the estimated coefficient on non-interacted $Z$ is the estimated effect in the period for which the dummy is omitted.
    - The estimated coefficient on the interaction between $Z$ and the dummy for period $t$ is the estimated difference between the effect of $Z$ in period $t$ and the effect in the omitted period.
    - The joint test of the interaction terms tests the hypothesis that the coefficients (effects) of $Z$ are the same in all periods.
    - The simple test of the interaction term for the period $t$ dummy tests whether the effect of $Z$ in period $t$ differs from the effect in the omitted period.
- **Using aggregate variables that vary only over time (not across units)**
  - o Suppose that we think that the reason for variation in either the intercept or a slope coefficient across periods is due to changes in one particular variable (the aggregate unemployment rate, for example).
  - o In this case, we can include that variable (for intercept effects) and perhaps interactions of that variable with some regressor $Z$ (to capture effects on unemployment on the marginal effect of $Z$).
  - o Interpretation of these coefficients is standard for continuous interaction variables.
- **Limitations on variables that vary only over time**
  - o If we include time dummies, we cannot include any other variables that vary only over time.
    - Any variable that varies only over time can be expressed as a linear function of the dummies.
      - If there are two periods with unemployment = 4 in the first period and 6 in the second, then $U = 4 + 2D_2$ , where $D_2$ is a dummy equal to one in the second period. Thus, including $U$, $D_2$, and a constant will result in perfect multicollinearity.

- - Same thing happens with more periods and/or more variables like $U$ that vary only over time (and not across units).
  - o If there are $T$ time periods represented in the data, there can be at most $T-1$ only-time-varying variables in the regression (assuming no dummies).
    - Again, there can be only $T$ distinct "observations" for any such variable, so just as $N$ must be at least $k+1$ in a standard regression, we can only identify the effects of $T-1$ such variables. Otherwise we have perfect multicollinearity.
    - We must also be careful about degrees of freedom here, because although we may have a large $N$, if we have only $T=2$, we don't really have much information about the effect of the one time-only-varying observation whose effect we can estimate.

## Pooled estimation with panel data

- Simplest method is just to estimate by OLS with a sample of $NT$ observations, not recognizing panel structure of data
  - o Standard OLS would assume homoskedasticity and no correlation between unit $i$'s observations in different periods (or between different units in the same period)
  - o **Clustered standard errors** are standard errors that are robust to correlation between error terms of same unit and heteroskedasticity over time:
    - $$\text{cov}\left(e_{it}, e_{is}\right) = \psi_{ts},$$
      $$\text{var}\left(e_{it}\right) = \psi_{tt}.$$
    - Option vce(cluster) in Stata

## Unit (entity) fixed effects

- How might model vary across units?
  - o All coefficients might be different: $y_{it} = \beta_{1i} + \beta_{2i} x_{2it} + \beta_{3i} x_{3it} + e_{it}$
    - Implies need to estimate separate regressions for each $i$
    - Impractical if $T$ is small.
  - o Constant terms might be different but slopes the same:
    $y_{it} = \beta_{1i} + \beta_2 x_{2it} + \beta_3 x_{3it} + e_{it}$
    - This is **unit fixed-effects model**
    - Can be estimated in two ways: LSDV or with de-meaned data
      - Former is impractical with large $N$
    - Intercept (dummy variable) estimates only converge asymptotically as $T \to \infty$.
- **Least-squares with unit dummy variables**

- o Introduce a dummy variable for each unit (individual) in sample:

$$y_{it} = \sum_{j=1}^{N} \beta_{1j} D_{ji} + \beta_{2i} x_{2it} + \beta_{3i} x_{3it} + e_{it}$$

- o We now have $N + K - 1$ regressors with $NT$ observations. If $T$ is moderately large (> 2 or 3) then we can get reliable estimates, but note computational difficulty: moment matrix is $N + K - 1 \times N + K - 1$, which with $N > 100$ is probably computationally difficult.
  - ▪ Note that even if computing power is sufficient to blast through the problem, inverting very large matrices can be subject to more rounding error than simpler problems.
- o Can test the block of unit dummies to see if individual fixed effect is statistically significant.
- **Fixed-effects estimator using de-meaned data**
  - o $y_{it} = \beta_{1i} + \beta_2 x_{2it} + \beta_3 x_{3it} + e_{it}$
  - o Averaging across $T$ time periods for each unit $i$:

$$\frac{1}{T}\sum_{t=1}^{T} y_{it} = \beta_{1i} + \beta_2 \frac{1}{T}\sum_{i=1}^{T} x_{2it} + \beta_3 \frac{1}{T}\sum_{t=1}^{T} x_{3it} + \frac{1}{T}\sum_{t=1}^{T} e_{it}$$

$$\overline{y}_i = \beta_{1i} + \beta_2 \overline{x}_{2i} + \beta_3 \overline{x}_{3i} + \overline{e}_i$$

  - o Subtracting the means from the original equation yields

    $y_{it} - \overline{y}_i = \beta_2 \left( x_{2it} - \overline{x}_{2i} \right) + \beta_3 \left( x_{3it} - \overline{x}_{3i} \right) + \left( e_{it} - \overline{e}_i \right)$ or $\tilde{y}_{it} = \beta_2 \tilde{x}_{2it} + \beta_3 \tilde{x}_{3it} + \tilde{e}_{it}$ where the ~ indicates deviation from the unit mean
  - o This is the "within-unit estimator"
    - ▪ It uses only variation over time within each unit to identify the coefficients
    - ▪ If one $x$ varies mostly across units rather than over time, there is not going to be much information in the sample to allow it to be identified by fixed-effects estimation.
    - ▪ If a variable varies *only* across units (e.g., ethnicity or sex), then its effects cannot be identified at all in a fixed-effects model
      - All ~ values will be zero because each observation equals the unit mean.
      - This also happens in LSDV because the $x$ in question will be perfectly collinear with the unit dummies.
  - o The constant term is gone because both $\tilde{y}_{it}$ and $\tilde{x}_{kit}$ have zero means.
    - ▪ We can estimate the mean of the individual unit constant terms from "between unit estimator" $\overline{y}_i = \overline{\beta}_1 + \beta_2 \overline{x}_{2i} + \beta_3 \overline{x}_{3i} + \left( \overline{e}_i + \beta_{1i} - \overline{\beta}_1 \right)$
    - ▪ We can calculate the estimated intercept for any individual unit as $\hat{\beta}_{1i} = \overline{y}_i - \hat{\beta}_2 \overline{x}_{2i} - \hat{\beta}_3 \overline{x}_{3i}$ under the assumption that $E\left( \overline{e}_i \right) = 0$

- o Degrees of freedom
  - Although there are $NT$ observations on the ~ variables, only $N(T-1)$ of them are independent.
  - Should use $\hat{\sigma}^2 = SSE / (NT - N - K + 1)$ to reflect this: An FE estimator will correct this but if you de-mean yourself and use OLS it will not.
  - Note that using LSDV will do this automatically because there will be $N - 1$ additional coefficients being estimated.
- o An alternative to de-meaning the data is to subtract one time period from all the others.

## *Time fixed effects*

- If there are characteristics (especially unobserved ones) that are common to all units but vary across time, then we can use **time fixed effects**, which are just like the time dummies that we discussed in the pooling section.
- The model is then $Y_{it} = \beta_1 X_{it} + \lambda_t + u_{it}$. We omit the constant term if all $T$ dummies are used to avoid collinearity; alternatively, we can omit the dummy for one time period.
- The methods of estimation are identical to the unit fixed-effects model.
  - o We can, equivalently
    - Estimate the model with time dummies, or
    - Estimate with $y$ and $x$ expressed as deviations from time-period means across units.
- Any variable that varies only across time, and not across units, will be collinear with the dummy variables (or zero when de-meaned) and its effect cannot be estimated.
- We can also combine both unit and time fixed effects.
  - o Either LSDV with both unit and time dummies, or
  - o Demeaning the data both with respect to time and with respect to units.
    - To do this, we calculate

      $$\tilde{y}_{it} = (y_{it} - \overline{y}_i) - \frac{1}{N}\sum_{j=1}^{N}(y_{jt} - \overline{y}_j) = (y_{it} - \overline{y}_i) - (\overline{y}_t - \overline{\overline{y}}), \text{ where } \overline{\overline{y}} \text{ is the}$$

      overall mean across both units and time, and regress it on a similarly transformed $x$.
    - This is sometimes called the "differences-in-differences" estimator because it excludes the effects of changes that are strictly over time (taken out with time dummies or demeaning) *and* the effects of changes that are strictly across units (taken out with unit dummies or demeaning). This leaves only differences across units in how the variables change over time to estimate $\beta$.

## Random-effects models

- The fixed-effects model thinks of $\beta_{1i}$ as a fixed set of constants that differ across $i$.
- The random-effects model thinks of $\beta_{1i}$ as a random variable (with mean $\bar{\beta}_1$) that has one value for each $i$ drawn from a given probability distribution.
  - $$\beta_{1i} = \bar{\beta}_1 + u_i$$
    $$E(u_i) = 0$$
    $$\text{cov}(u_i, u_j) = 0, \, i \neq j$$
    $$\text{var}(u_i) = \sigma_u^2$$
  - $$y_{it} = (\bar{\beta}_1 + u_i) + \beta_2 x_{2it} + \beta_3 x_{3it} + e_{it}$$
    $$= \bar{\beta}_1 + \beta_2 x_{2it} + \beta_3 x_{3it} + (u_i + e_{it})$$
    $$= \bar{\beta}_1 + \beta_2 x_{2it} + \beta_3 x_{3it} + v_{it}$$
  - This leads to a particular pattern of correlation among the error terms $v_{it}$
    - Error terms of observations corresponding to the same $i$ will be correlated because they have $u_i$ in common
      - $$\text{corr}(v_{it} v_{is}) = \frac{\text{cov}(v_{it}, v_{is})}{\sqrt{\text{var}(v_{it}) \text{var}(v_{is})}} = \frac{\sigma_u^2}{\sigma_u^2 + \sigma_e^2} = \rho$$
    - Error terms of observations with different $i$ will be uncorrelated (by assumption)
  - We can estimate $\rho$ by looking at the correlation of error terms within units and use the estimated $\hat{\rho}$ to do feasible GLS: this is the **random-effects estimator**
- Assumptions of the random-effects estimator
  - $E(v_{it}) = 0$
  - $\text{var}(v_{it}) = \sigma_e^2 + \sigma_u^2$
  - $\text{cov}(v_{it}, v_{is}) = \sigma_u^2, \, t \neq s$
  - $\text{cov}(v_{it}, v_{js}) = 0, \, i \neq j$
  - $\text{cov}(e_{ij}, x_{kij}) = 0, \, k = 2, 3, ..., K$
  - $\text{cov}(u_i, x_{kit}) = 0, \, k = 1, 2, ..., K$
- Testing for presence of random effects vs. OLS assumption of independent errors:
  - Is $\sigma_u^2 = 0$?
    - If so, then there are no correlations because the $u$ error term is degenerate
    - A positive value of $u_i$ is a positive expected value of $v_{it}$ for that $i$.
      - We can estimate with $\hat{u}_i = \dfrac{1}{T} \sum_{t=1}^{T} \hat{v}_{it}$.

- Then estimate $\sigma_u^2 = \frac{1}{N}\sum_{i=1}^{N}\hat{u}_i$ .

  o One-tailed LM test for $H_0 : \sigma_u^2 = 0$ against $H_1 : \sigma_u^2 > 0$ in large samples is

  $$LM = \sqrt{\frac{NT}{2(T-1)}}\left[\frac{\sum_{i=1}^{N}\left(\sum_{t=1}^{T}\hat{v}_{it}\right)^2}{\sum_{i=1}^{N}\sum_{t=1}^{T}\hat{v}_{it}^2} - 1\right] \sim N(0,1)$$

  o If we reject $\sigma_u^2 = 0$ , then OLS is inefficient and we should use random effects (if assumptions are satisfied)

- Estimation in presence of random effects
  o OLS is unbiased, but inefficient. Can use clustered standard errors to correct bias in standard errors.
  o Efficient FGLS estimator:
    - First use OLS residuals to estimate $\sigma_u^2$ and $\sigma_e^2$ .

    - Calculate $\alpha = 1 - \dfrac{\sigma_e}{\sqrt{T\sigma_u^2 + \sigma_e^2}}$

    - Transform model by "quasi-de-meaning"
      $$y_{it}^* \equiv y_{it} - \alpha\bar{y}_i$$
      $$x_{1it}^* \equiv 1-\alpha$$
      $$x_{kit}^* \equiv x_{kit} - \alpha\bar{x}_{ki}, \ k=1,2,...,K$$
      $$v_{it}^* \equiv v_{ij} - \alpha\bar{v}_i$$

    - With this transformation, the $v^*$ error term is homoskedastic and not autocorrelated, so we can apply OLS to the transformed model to get efficient estimators
    - Note that the random-effects estimator = fixed-effects estimator if $\alpha = 1$, which is the limit as $T \rightarrow \bullet$ .

- Critical assumption of random-effects model is that the unit-specific error term (representing missing variables that are constant within each unit over time) are independent of the included regressors.
  o This is often not sensible.

## *Random effects or fixed effects?*

- Fixed-effects models have the advantage of not requiring cov($x$, $u$) = 0, which is often difficult to justify.
  o Fixed-effects models are fully efficient as $N$ gets large even if the true model is random effects.

- However, standard fixed-effects models cannot identify the effects of any variables that vary only across units (and has difficulty identifying effects if most of the meaningful variation is across units).
- Can do a Hausman test to examine whether the random-effects model is OK. (It is a nested sub-model of the fixed-effects model.)
  - The Hausman test is rejected if
    - The estimates are sufficiently different, **and**
    - The fixed-effects estimators are sufficiently precise.
  - Common decision rule: Use random-effects unless the Hausman test rejects it.
- **Hausman-Taylor estimator** (xthtaylor)
  - Applies IV methods to random-effects model to overcome the correlation between the regressors and the error term so that we can use RE estimation to estimate the effects of variables that do not vary across time within individuals
  - $y_{it} = \beta_1 + \beta_2 x_{it,exog} + \beta_3 x_{it,endog} + \beta_4 w_{i,exog} + \beta_5 w_{i,endog} + u_i + e_{it}$
  - The $w$ variables do not vary over time for any individual, the $x$ variables do
  - The number of $x_{it,exog}$ must be $\geq$ the number of $w_{i,endog}$ in order for identification to be possible.
  - This estimation uses the demeaned $x_{it,endog}$ as instruments for the quasi-demeaned values and the means across time of $x_{it,exog}$ as instruments for the $w_{i,endog}$

## SUR and panel data

- Panel estimation by random effects is a lot like SUR because we are taking account of the same kind of correlation across observations.
- With small $N$ and large $T$, it sometimes makes sense to treat the model as SUR and we can test whether the coefficients are the same across $i$ as cross-equation coefficient restrictions.

## Class demonstration?

- Dataset: S&W's Seatbelts.dta
  - Show dataset
  - Define as panel
    - xtset fips year
  - Discuss missing values problem
  - Note two state identifiers, one alpha and one numeric
- Following S&W's Empirical Exercise E10.2
  - Generate lnincome variable
  - Discuss expected results of regressing fatalityrate on sb_usage speed65 speed70 drinkage21 ba08 lnincome age
- OLS regression

- o sb_usage has "wrong" sign
    - ▪ Authors argue that this is omitted variable bias and might be corrected partially by including state fixed effects.
  - o Other effects are plausible
  - o Send to outreg2 using fatal , word ctitle(OLS)
- FE regression
  - o Now sb_usage has the expected sign
  - o Other variables decline in coefficient magnitude
  - o With and without the "robust" option, which gives clustered standard errors in this case.
  - o ctitle(FE)
- Adding time dummies
  - o What would time dummies control for?
    - ▪ Changes over time in such things as air bags, other auto or highway safety features
  - o xi: xtreg … i.year , fe
  - o Note loss in significance of all variables except age and drinkage, which become stronger.
  - o This is "differences in differences" estimator and has very high bar for significance.
  - o ctitle(FE & time)
- RE regression
  - o Results are similar to other methods
  - o Note requirement that the error term be uncorrelated with regressors
    - ▪ FE output has measure of the correlation between the fe dummies and the regressors.