# ECON 672

## Week 5: Panel Data and Fixed Effects (Within) Estimator

Samuel Rowe, PhD 12/22/2022

# Overview

- Takeaway

- Fixed Effects DAG

- Estimation

- Pooled OLS

- Fixed Effects or Within Estimator

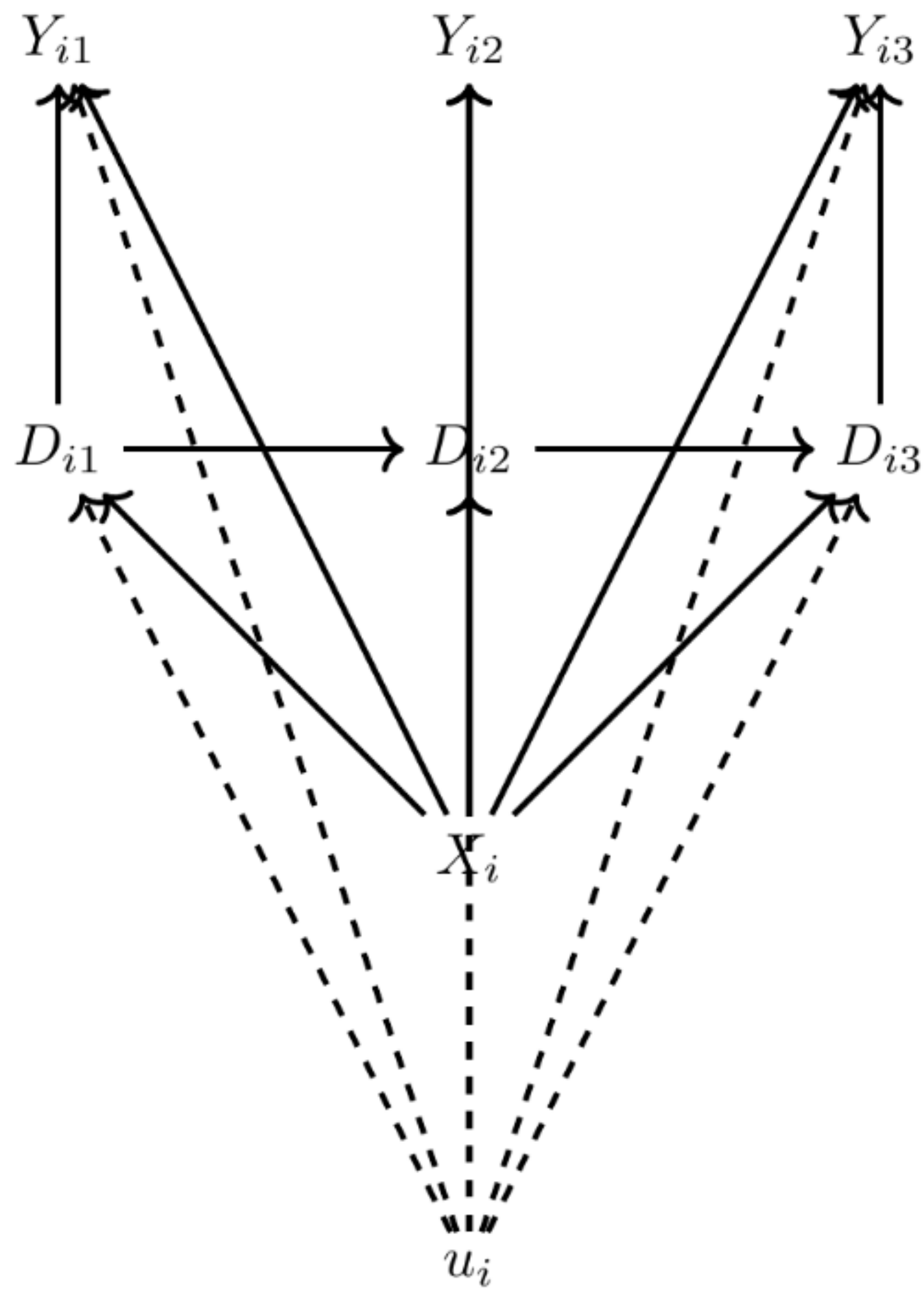- Caveats of Fixed Effects

- Example

# The Takeaway

- Panel Data

  - When we observe the same $i^{th}$ unit over time

  - With cross-sectional data, we do not observe the same $i^{th}$ over time

  - We can use the panel structure to control for unobserved or observed time-invariant heterogeneity

  - Time-invariant heterogeneity (observed or unobserved) are covariates that vary across units but do not vary within units

- Fixed Effects (Within) Estimator

  - Is an identification strategy that can be employed with panel data

  - It works when we have unobserved time-invariant confounders

# Takeaway (Fixed Effects)

- Strengths

  - We can control for unobserved confounders that do not vary over time with panel data and the fixed effects estimator

- Weaknesses

  - We cannot control for unobserved confounders that do vary over time

  - We cannot control for simultaneity or reverse causation

- Assumptions

  - Strict Exogeneity Assumption (Independence Assumption - not testable)

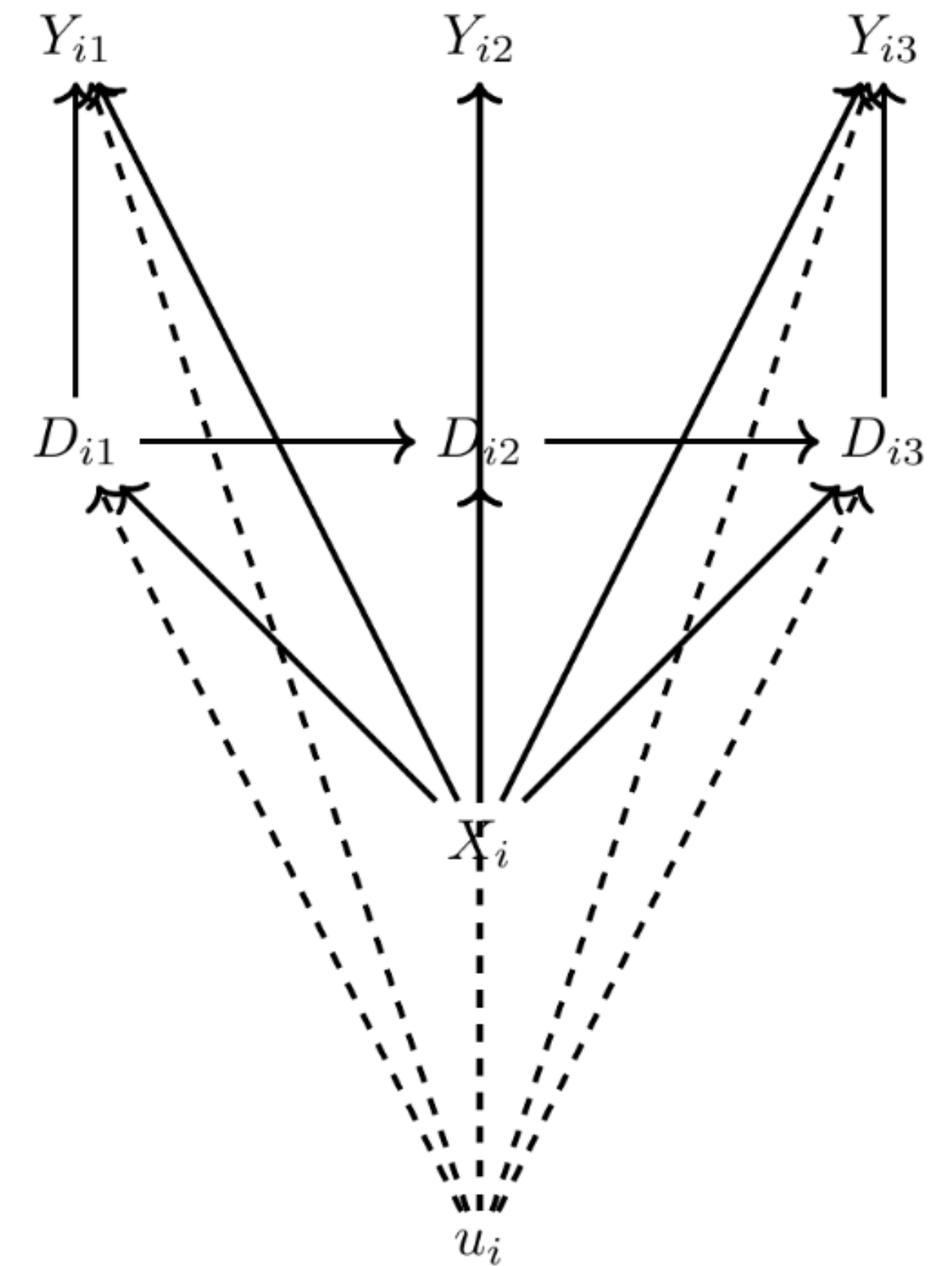  - Rank Assumption (covariates must vary - testable)

# Directed Acyclic Graphs: Panel Data

# Directed Acyclic Graphs: Panel Data

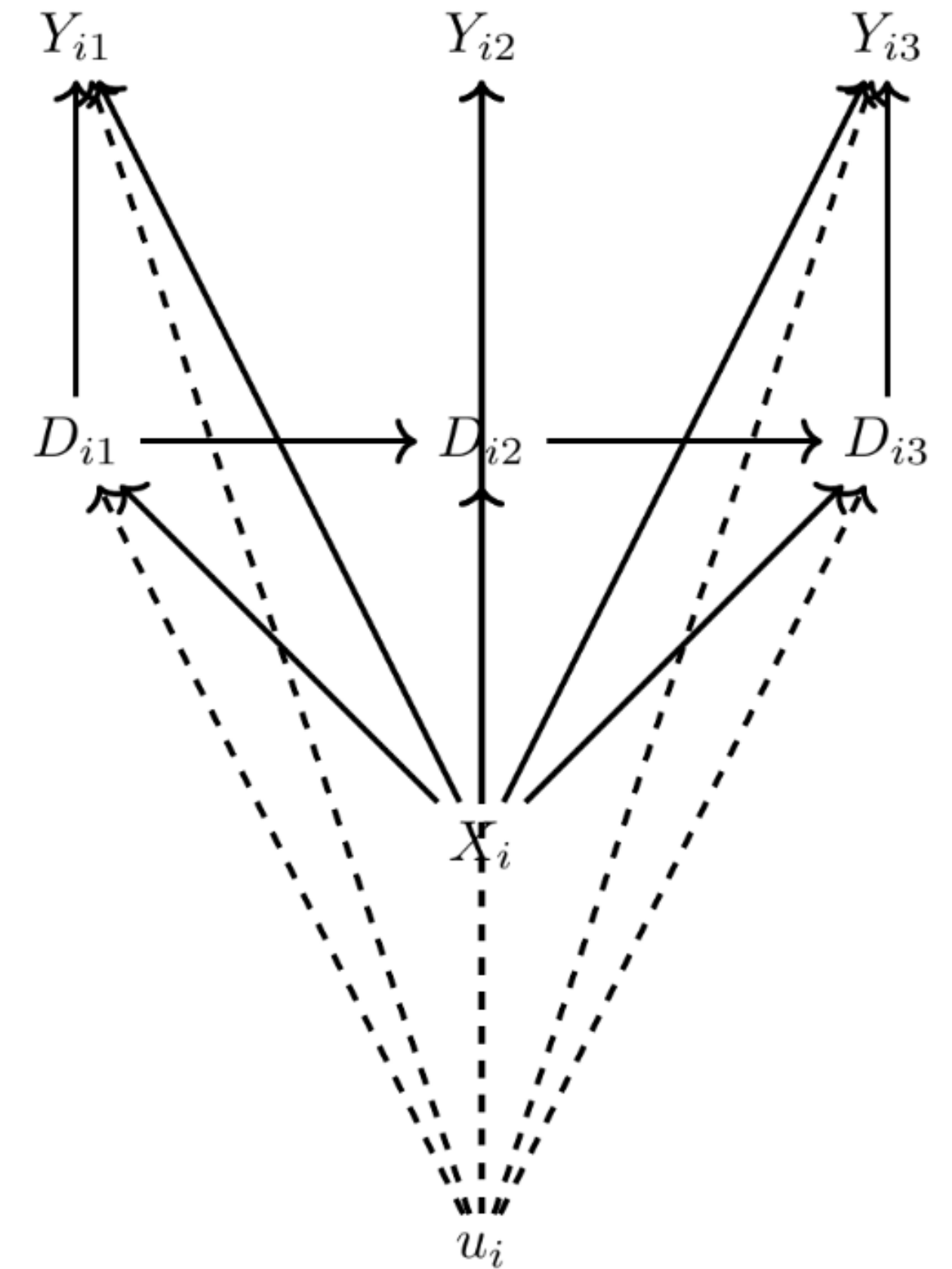# Directed Acyclic Graphs: Panel Data

- The most complex DAG yet

- $D_{it}$ can vary over time

  - $D_{i1} \to D_{i2}$

  - $D_{i2} \to D_{i3}$

- $Y_{it}$ can vary over time

  - $D_{i1} \to Y_{i1}$

  - $D_{i2} \to Y_{i2}$

  - $D_{i3} \to Y_{i3}$
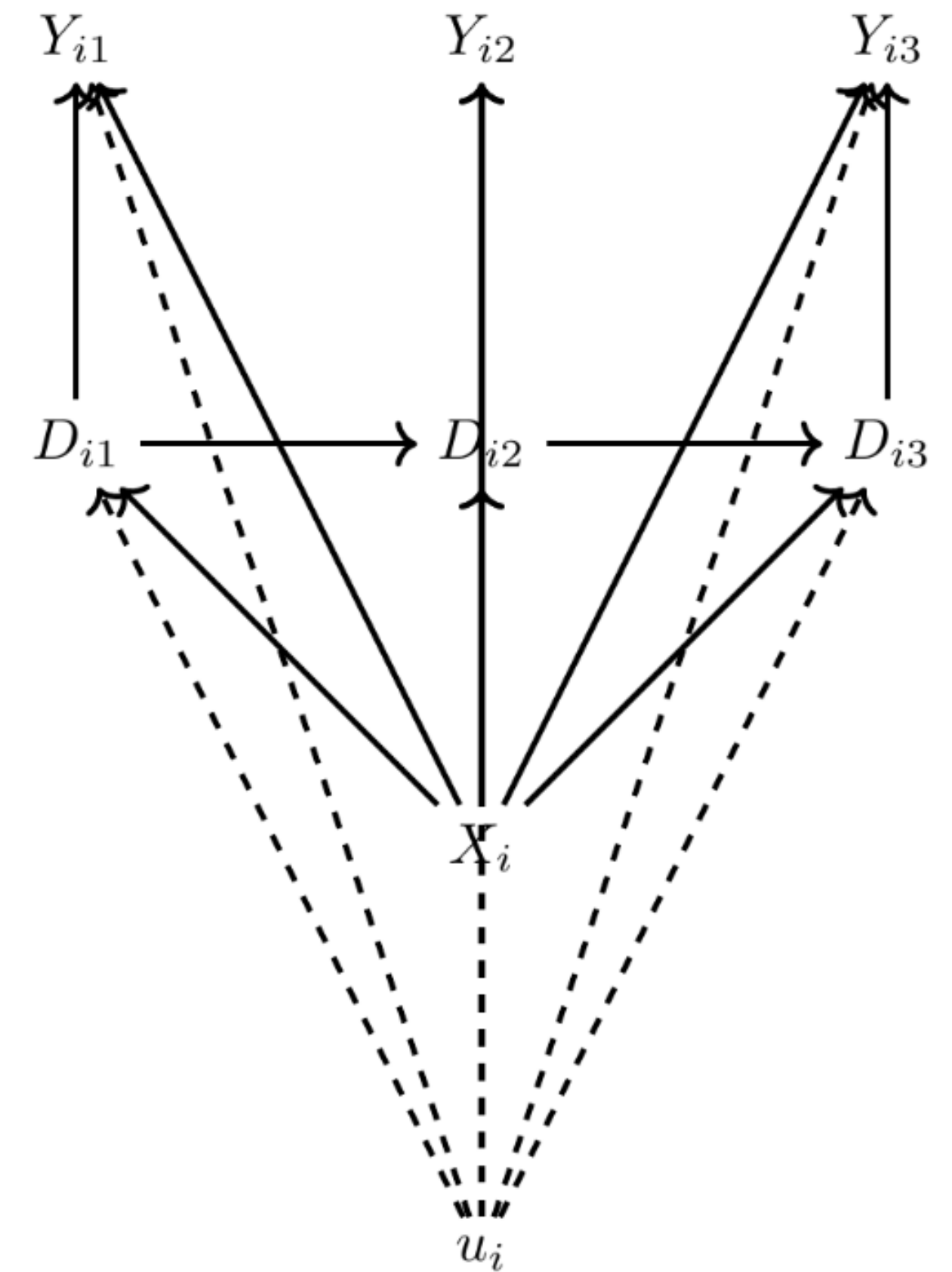
# Directed Acyclic Graphs: Panel Data

- $D_{it}$ affects $Y_{it}$ and $D_{i(t+1)}$

  - $Y_{it} \leftarrow D_{it} \rightarrow D_{i(t+1)}$

- We **assume** that outcomes are not affected by prior outcomes

  - $Y_{it} \perp Y_{i(t+1)}$

- We assume that prior treatment does not affect current outcomes ***directly***, but mediated

  - $D_{it} \rightarrow D_{i(t+1)} \rightarrow Y_{i(t+1)}$

# Directed Acyclic Graphs: Panel Data

- $X_i$ are observed confounders that do not vary over time

- $u_i$ are unobserved confounders that do not vary over time

- $D_{it}$ and $Y_{it}$ have t in the subscripts and $X_i$ and $u_i$ do not

- This means that D and Y vary over time and X and u are time-invariant

# Directed Acyclic Graphs: Panel Data

- Given that $X_i$ and $u_i$ are time-invariant

  - $Y_{it} \leftarrow X_i \rightarrow D_{it}$

  - $Y_{it} \leftarrow - u_i - \rightarrow D_{it}$

- $X_i$ and $u_i$ are time-invariant confounders that affect treatment and outcome in every period

# Directed Acyclic Graphs: Panel Data

- Under this scenario, we can use the Fixed Effects (Within) Estimator

- The Fixed Effects (Within) Estimator closes all of the confounding backdoors that do not vary over time

- Note: if u varied over time, then the Fixed Effects (Within) Estimator would not close all backdoor pathways

# Panel Data Estimators

- Panel data refers to data where an unit of observation (individual, firm, county, etc.) is observed longitudinally or over time (more than one time period)

- If observed or unobserved confounders do not vary across time for a unit (but vary across units)

  - We can use panel data to identify the causal effect

- Estimators with Panel Data

  - Pooled OLS

  - Fixed Effects (Within) Estimator

  - First-differencing (not our focus)

  - Random Effects Estimator (requires special assumption and not our focus)

# Notation

- We will utilize traditional notation instead of potential outcomes for panel data

- Let Y and D be random variables

  - Where $D \equiv D(D_i, D_2, \ldots, D_k)$

- Let u be an unobserved random variable

- We are interested in the partial effect of $D_j$ from

  - $E[Y \,|\, D_1, D_2, \ldots, D_k, u]$

# Notation

- We observe a sample $i = 1,2,...,N$ cross-sectional units for $t = 1,2,...,T$ time periods, which is a a balanced panel (no missing observations in matrix)

- Cross-sectional independence: individuals in the panel are identical and independent draws from the population $\{Y_i, D_i, u_i\}_{i=1}^{N} \sim i.i.d$

- 
$$
Y = \begin{bmatrix} Y_{11} & Y_{21} & \cdots & Y_{N1} \\ Y_{12} & Y_{22} & \cdots & Y_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ Y_{1T} & Y_{2T} & \cdots & Y_{NT} \end{bmatrix} \text{ and } D = \begin{bmatrix} D_{11} & D_{21} & \cdots & D_{N1} \\ D_{12} & D_{22} & \cdots & D_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ D_{1T} & D_{2T} & \cdots & D_{NT} \end{bmatrix}
$$

# Regression Notation

- $Y_{it} = \delta D_{it} + u_i + \varepsilon_{it}$ is our unobserved effects model

- Where $Y_{it}$ is our outcome of interest for $i = 1,2,...,N$ over $t = 1,2,...,T$

- $\delta$ is our treat effect of interest

- $D_{it}$ is our treatment for $i = 1,2,...,N$ over $t = 1,2,...,T$

- $u_i$ is the sum of all time-invariant person-specific characteristics, such as ability

- $\varepsilon_{it}$ is the idiosyncratic error - includes unobserved time-varying covariates

# Pooled OLS

- Pooled OLS is the simplest panel data estimator

- It does not account for the panel structure

  - $Y_{it} = \delta D_{it} + \eta_{it}; \ t = 1,2,...,T$

- Where the composite error

  - $\eta_{it} \equiv u_i + \varepsilon_{it}$

  - Where $u_i$ is time-invariant heterogeneity

  - $\varepsilon_{it}$ is time-varying heterogeneity

- We need to **assume** that $u_i$ is does not impact $D_{it}$ for all time periods

# Pooled OLS

- In order to identify the causal effect with Pooled OLS

  - We need to show that

  - $E[\eta_{it} | D_{i1}, D_{i2}, \ldots, D_{iT}] = E[\eta_{it} | D_{it}] = 0 \; \forall \; t = 1,2,...,T$

- TL;DR

  - We need to ignore omitted variable bias, which is unlikely to work well

  - This is likely not a credible assumption

# Fixed Effects (Within) Estimator

- The Fixed Effects (FE) Estimator is also called the Within Estimator

  - It accounts for the variation *within* a unit of observation over time

- The FE estimator will does a better job of controlling for observed and unobserved time-invariant confounders

# Fixed Effects (Within) Estimator

- Our unobserved effects model

  - $Y_{it} = \delta D_{it} + u_i + \varepsilon_{it}$

- We think of our fixed effect $u_i$ as a covariate to be estimated

- The OLS estimation for $\hat{\delta}$ and $\hat{u}_i$ under minimizing sum of squared (over i and t)

  - $$(\hat{\delta}, \hat{u}_1, \hat{u}_2, \ldots, \hat{u}_N) = \arg \min_{b, m_1, \ldots, m_N} \sum_{i=1}^{N} \sum_{t=1}^{T} (Y_{it} - D_{it}b - m_i)^2$$

  - This means we include $N$ number of individual dummy variables in the regression of $Y_{it}$ on $D_{it}$

# Fixed Effects (Within) Estimator

- We'll use our two first order conditions

  - Recall that $E[u \,|\, x] = 0$ and $E[ux] = 0$

- $$\sum_{i=1}^{N} \sum_{t=1}^{T} D'_{it}(Y_{it} - D_{it}\hat{\delta} - \hat{u}_i) = 0$$

- $$\sum_{t=1}^{T} (Y_i t - D_{it}\hat{\delta} - \hat{u}_i) = 0 \text{ for } i = 1,2,...,N$$

# Fixed Effects (Within) Estimator

- Therefore, for $i = 1,2,...,N$

$$\hat{u}_i = \frac{1}{T}\sum_{t=1}^{T}(Y_{it} - D_{it}\hat{\delta}) = \bar{Y}_i - \bar{D}_i\hat{\delta}$$

- Where $\bar{D}_i \equiv \frac{1}{T}\sum_{t=1}^{T}D_{it}$ and $\bar{Y}_i \equiv \frac{1}{T}\sum_{t=1}^{T}Y_{it}$

- Plug in the results into the first first-order condition: $\displaystyle\sum_{i=1}^{N}\sum_{t=1}^{T}D_{it}'(Y_{it} - D_{it}\hat{\delta} - \hat{u}_i) = 0$

# Fixed Effects (Within) Estimator

- $\hat{\delta} = \dfrac{\sum_{i=1}^{N} \sum_{t=1}^{T} (D_{it} - \bar{D}_i)'(Y_{it} - \bar{Y}_i)}{\sum_{i=1}^{N} \sum_{t=1}^{T} (D_{it} - \bar{D}_i)'(D_{it} - \bar{D}_i)} = \dfrac{\sum_{i=1}^{N} \sum_{t=1}^{T} \ddot{D}_{it}' \ddot{Y}_{it}}{\sum_{i=1}^{N} \sum_{t=1}^{T} \ddot{D}_{it}' \ddot{D}_{it}}$

  - Where $\ddot{D}_{it} \equiv D_{it} - \bar{D}_i$ and $\ddot{Y}_{it} \equiv Y_{it} - \bar{Y}_i$

- Recall $\hat{\delta} = \dfrac{C(Y, D)}{V(D)}$

# Fixed Effects (Within) Estimator

- TL;DR

  - Using time-demeaned variables $\ddot{D}_{it} \equiv D_{it} - \bar{D}_i$ and $\ddot{Y}_{it} \equiv Y_{it} - \bar{Y}_i$ is equivalent to a regression of $Y_{it}$ on $D_{it}$ with unit-specific dummy variables

- This is why it is called the Within Estimator, since we are utilizing the variation within a specific-unit

- When we include unit-specific fixed effects and year-specific fixed effects, this is called the "two-way fixed effects" estimator

  - We'll cover this later

# Fixed Effects (Within) Estimator

- Using time-demeaned variables, the time-invariant confounders zeros out

  - Time-invariant confounders do not vary, such that $E[c] = c$

  - Demeaning eliminates time-invariant observed and unobserved confounders, such that $u_i - \bar{u}_i = 0$

- $Y_{it} = \delta D_{it} + u_i + \varepsilon_{it}$

  - Demean across $T$: $(Y_{it} - \bar{Y}_i) = (\delta D_{it} - \delta \bar{D}_i) + (u_i - \bar{u}_i) + (\varepsilon_{it} - \bar{\varepsilon}_{it})$

- $\ddot{Y}_{it} = \delta \ddot{D}_{it} + \ddot{\varepsilon}_{it}$

# Implement Fixed Effects

- There are three ways that we can implement fixed effects in our regression

- 1) Demean and regression $\ddot{Y}_{it}$ on $\ddot{D}_{it}$ (and need to correct for degrees of freedom)

- 2) Regress $Y_{it}$ on $D_{it}$ and unit-specific dummy variables (dummy variable regression)

- 3) Regress $Y_{it}$ on $D_{it}$ with canned fixed effects routine in Stata or R

# Assumptions

- There are a couple of necessary identifying assumption we need in order for the Fixed Effects (Within) Estimator to identify the causal effect

- Strictly exogenous assumption (independence assumption and not testable)

  - $E[\varepsilon_{it} | D_{i1}, D_{i2}, \ldots, D_{iT}, u_i] = 0$ for $t = 1,2,...,T$

- Rank assumption (variation is required for at least some units and is testable)

  - $rank\left( \sum_{t=1}^{K} E[\ddot{D}_{it}' \ddot{D}_{it}] \right) = K$

  - Recall that this is just the $V(D_{it})$ from $\hat{\delta} = \dfrac{C(Y_{it}, D_{it})}{V(D_{it})}$

# Assumptions

- Our strict exogeneity assumption will be similar to our independence assumption

  - Our time-invariant confounders can be related to $D_{it}$ since we control for them with fixed effects

  - We need to be concerned about unobserved time-varying confounders, which will violate this assumption

- Our rank assumption requires that there be variation in treatment over time for at least some units of observation

  - Otherwise it will be 0 and violate the assumption

# Assumptions

- If our two main assumption hold, then the fixed effects estimator identifies the causal effect

- The fixed effect estimate is consistent $(p \lim_{N \to \infty} \hat{\delta}_{FE,N} = \delta)$ and unbiased

  - This holds as long as the number of clusters is large enough

    - This will be an issue we'll run into with Diff-in-Diff

# Caveats of Fixed Effects (Within) Estimator

- There are two key caveats we need to be aware of

- 1) Fixed effects cannot resolve reverse causality

  - $Y \rightarrow D$

- 2) Fixed effects cannot control for time-varying unobserved confounders

  - We have to assume that there are no time-varying unobserved confounders

  - $E[\varepsilon_{it} | D_{i1}, D_{i2}, \ldots, D_{iT}, u_i] = 0$ for $t = 1,2,...,T$

# Caveats

- Fixed Effects Estimator cannot handle reverse causality

  - Cornwell and Trumbell (1994) find positive correlations between policing and crime rates using panel data from North Carolina

Table 8.1: Panel estimates of police on crime

| Dependent variable | Between | Within | 2SLS (FE) | 2SLS (No FE) |
|---|---|---|---|---|
| Police | 0.364 | 0.413 | 0.504 | 0.419 |
| | (0.060) | (0.027) | (0.617) | (0.218) |
| Controls | Yes | Yes | Yes | Yes |

North Carolina county level data. Standard errors in parenthesis.

# Caveats

- Does this mean that police cause more crime?

  - It is likely a reverse causality problem and they did not identify the causal effect

  - Police spending is a function of crime rates and crimes rates

  - **Simultaneity bias** is creating bias for their estimated treatment effects

- Becker (1986) predicts the police spending per capita will theoretically reduce crime

  - What is the relationship between the outcome, treatment, and covariates of interest?

# Caveats

- We need to assume no time-varying unobserved confounders

  - The presence of any time-varying unobserved confounders will prevent the backdoor criterion from being satisfied

  - This is just **omitted variable bias**

- You will need another research design to handle time-varying unobserved confounders

  - Demean time-varying confounder is just moved to the composite error term and the strict exogeneity assumption (independence assumption) is violated

# Example: Returns to Marriage

- Cornwell and Rupert (1997) attempt to estimate the returns to marriage on earnings

  - The SDO shows that married men earn more than unmarried men

  - The SDO is likely biased from confounders and selection into marriage

- We'll use panel data estimators to see the returns to marriage

  - What does the Feasible Generalized Least Squared model compared to Fixed Effects?

# Example: Returns to Marriage

- Let's set up the model for individuals $i$ observed over four periods $t$

    - $Y_{it} = \alpha + \delta M_{it} + \beta X_{it} + A_i + \gamma_i + \varepsilon_{it}$

- Where

    - $Y_{it}$ is earnings for individual $i$ earnings in time period $t$

    - $M_{it}$ is the outcome of interest and $\delta$ is our treatment effect of interest

    - $X_{it}$ is a set of observable covariates for individual $i$ earnings in time period $t$

    - $A_i$ is unobserved time-invariant ability for individual $i$

    - $\gamma_i$ is unobserved time-invariant confounders for individual $i$

    - $\varepsilon_{it}$ is our idiosyncratic error or unobserved determinants of wage that are assumed to be unrelated to $M_{it}$

# Example: Returns to Marriage

| Dependent variable | FGLS | Within | Within | Within |
|---|---|---|---|---|
| Married | 0.083 | 0.056 | 0.051 | 0.033 |
|  | (0.022) | (0.026) | (0.026) | (0.028) |
| Education controls | Yes | No | No | No |
| Tenure | No | No | Yes | Yes |
| Quadratics in years married | No | No | No | Yes |

# Example: Returns to Marriage

- Cornwell and Rupert (1997) find that the Feasible Generalized Least Squares model is upward biased

  - After controlling for time-varying covariates, such as education, tenure, and years of marriage

  - The treatment effect is not statistically significant

# Example

- Stata Exercise (Cunningham and Kendall, 2011, 2014, 2016)

Table 8.3: POLS, FE and Demeaned OLS Estimates of the Determinants of Log Hourly Price for a Panel of Sex Workers

| Depvar: | POLS | FE | Demeaned OLS |
|---|---|---|---|
| Unprotected sex with client of any kind | 0.013 | 0.051* | 0.051* |
| | (0.028) | (0.028) | (0.026) |
| Ln(Length) | $-0.308$*** | $-0.435$*** | $-0.435$*** |
| | (0.028) | (0.024) | (0.019) |
| Client was a Regular | $-0.047$* | $-0.037$** | $-0.037$** |
| | (0.028) | (0.019) | (0.017) |
| Age of Client | $-0.001$ | 0.002 | 0.002 |
| | (0.009) | (0.007) | (0.006) |
| Age of Client Squared | 0.000 | $-0.000$ | $-0.000$ |
| | (0.000) | (0.000) | (0.000) |
| Client Attractiveness (Scale of 1 to 10) | 0.020*** | 0.006 | 0.006 |
| | (0.007) | (0.006) | (0.005) |

# Example

- Stata Exercise (Cunningham and Kendall, 2011, 2014, 2016)

| | | | |
|---|---|---|---|
| Second Provider Involved | 0.055 | 0.113* | 0.113* |
| | (0.067) | (0.060) | (0.048) |
| Asian Client | −0.014 | −0.010 | −0.010 |
| | (0.049) | (0.034) | (0.030) |
| Black Client | 0.092 | 0.027 | 0.027 |
| | (0.073) | (0.042) | (0.037) |
| Hispanic Client | 0.052 | −0.062 | −0.062 |
| | (0.080) | (0.052) | (0.045) |
| Other Ethnicity Client | 0.156** | 0.142*** | 0.142*** |
| | (0.068) | (0.049) | (0.045) |
| Met Client in Hotel | 0.133*** | 0.052* | 0.052* |
| | (0.029) | (0.027) | (0.024) |
| Gave Client a Massage | −0.134*** | −0.001 | −0.001 |
| | (0.029) | (0.028) | (0.024) |

# Example

- Stata Exercise (Cunningham and Kendall, 2011, 2014, 2016)

| | | | |
|---|---|---|---|
| Age of provider | 0.003 | 0.000 | 0.000 |
| | (0.012) | (.) | (.) |
| Age of provider squared | −0.000 | 0.000 | 0.000 |
| | (0.000) | (.) | (.) |
| Body Mass Index | −0.022*** | 0.000 | 0.000 |
| | (0.002) | (.) | (.) |
| Hispanic | −0.226*** | 0.000 | 0.000 |
| | (0.082) | (.) | (.) |
| Black | 0.028 | 0.000 | 0.000 |
| | (0.064) | (.) | (.) |
| Other | −0.112 | 0.000 | 0.000 |
| | (0.077) | (.) | (.) |
| Asian | 0.086 | 0.000 | 0.000 |
| | (0.158) | (.) | (.) |

# Example

- Stata Exercise (Cunningham and Kendall, 2011, 2014, 2016)

| | | | |
|---|---|---|---|
| Imputed Years of Schooling | 0.020** | 0.000 | 0.000 |
| | (0.010) | (.) | (.) |
| Cohabitating (living with a partner) but unmarried | −0.054 | 0.000 | 0.000 |
| | (0.036) | (.) | (.) |
| Currently married and living with your spouse | 0.005 | 0.000 | 0.000 |
| | (0.043) | (.) | (.) |
| Divorced and not remarried | −0.021 | 0.000 | 0.000 |
| | (0.038) | (.) | (.) |
| Married but not currently living with your spouse | −0.056 | 0.000 | 0.000 |
| | (0.059) | (.) | (.) |
| N | 1,028 | 1,028 | 1,028 |
| Mean of dependent variable | 5.57 | 5.57 | 0.00 |

Heteroskedastic robust standard errors in parenthesis clustered at the provider level. $^{*}p < 0.10$, $^{**}p < 0.05$, $^{***}p < 0.01$

# Concluding Remarks

- Fixed Effects (Within) Estimator can be a powerful tool

- May not be utilized enough in federal evaluations

  - Unlike PSM, fixed effects can control for some unobservable confounders

- As long as

  - 1) There are no time-varying confounders, you have a straight-forward methodology that can identify the causal effect

  - 2) There is no reverse causality

    - This will require an instrumental variable strategy