# Large Scale Longitudinal Experiments: Estimation and Inference

*October 30, 2024*

Apoorva Lal
*Paper (joint with Alex Fischer, Matthew Wardrop): arXiv:2410.09952*

# Why study panel data in recommender systems?

- We often have longitudinal data for users/devices/... : $\{\mathbf{Z}_{i,t}\}_{i=1,t=1}^{N,T}$
- In typical experiments, we *flatten the time dimension* and compute a difference in means in the post-treatment window (optionally adjusting for average outcome in pre-treatment window)
- This leads us to miss out on treatment-effect dynamics, which are intrinsic to recommender systems

# Why study panel data in recommender systems?

- We often have longitudinal data for users/devices/... : $\{\mathbf{Z}_{i,t}\}_{i=1,t=1}^{N,T}$
- In typical experiments, we *flatten the time dimension* and compute a difference in means in the post-treatment window (optionally adjusting for average outcome in pre-treatment window)
- This leads us to miss out on treatment-effect dynamics, which are intrinsic to recommender systems
    1. Compliance is often not immediate: typically experimental treatments are rolled out but compliance isn't measured; analyzed as *intent to treat* (ITT)

# Why study panel data in recommender systems?

- We often have longitudinal data for users/devices/... : $\{\mathbf{Z}_{i,t}\}_{i=1,t=1}^{N,T}$
- In typical experiments, we *flatten the time dimension* and compute a difference in means in the post-treatment window (optionally adjusting for average outcome in pre-treatment window)
- This leads us to miss out on treatment-effect dynamics, which are intrinsic to recommender systems
  1. Compliance is often not immediate: typically experimental treatments are rolled out but compliance isn't measured; analyzed as *intent to treat* (ITT)
  2. Timing of compliance is informative: **activity bias :=** estimates for early compliers might not be externally valid
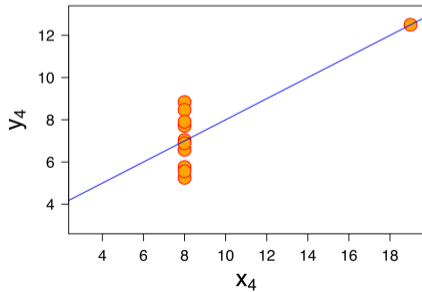
# Why study panel data in recommender systems?

- We often have longitudinal data for users/devices/... : $\{\mathbf{Z}_{i,t}\}_{i=1,t=1}^{N,T}$
- In typical experiments, we *flatten the time dimension* and compute a difference in means in the post-treatment window (optionally adjusting for average outcome in pre-treatment window)
- This leads us to miss out on treatment-effect dynamics, which are intrinsic to recommender systems
  1. Compliance is often not immediate: typically experimental treatments are rolled out but compliance isn't measured; analyzed as *intent to treat* (ITT)
  2. Timing of compliance is informative: **activity bias :=** estimates for early compliers might not be externally valid
  3. Dynamic effects are important: **novelty bias :=** estimates of early time periods might not be externally valid
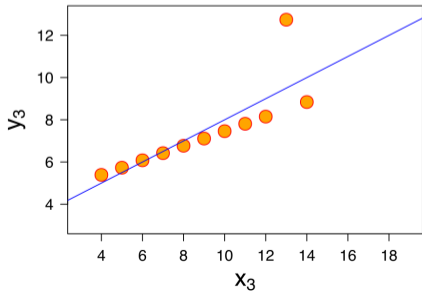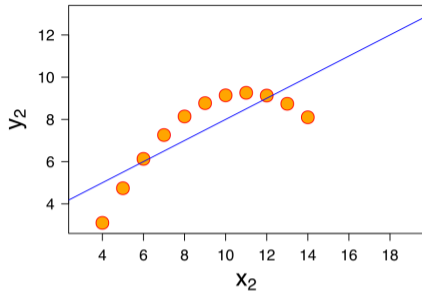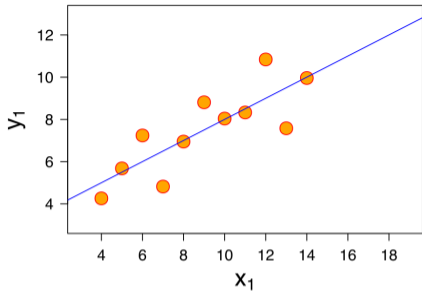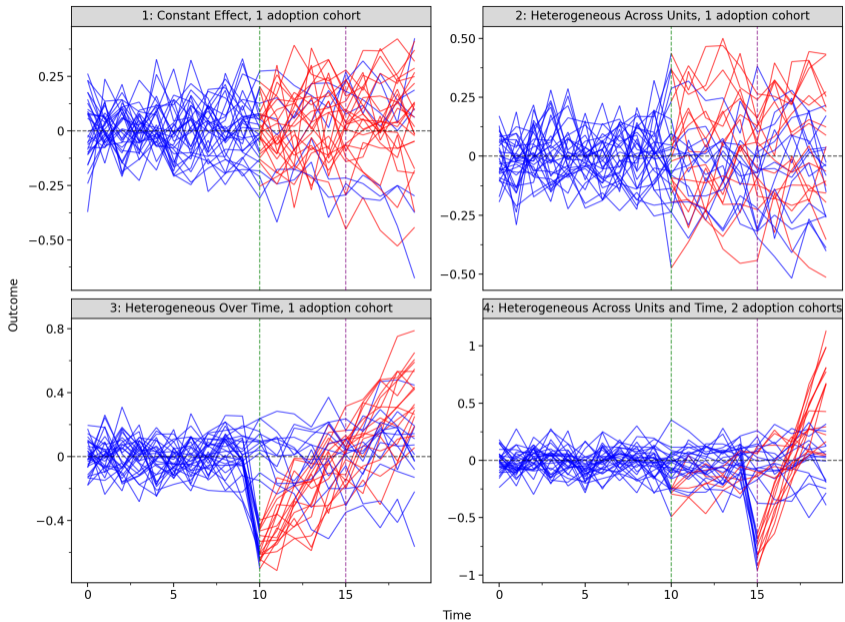
# Why study panel data in recommender systems?

- We often have longitudinal data for users/devices/... : $\{\mathbf{Z}_{i,t}\}_{i=1,t=1}^{N,T}$
- In typical experiments, we *flatten the time dimension* and compute a difference in means in the post-treatment window (optionally adjusting for average outcome in pre-treatment window)
- This leads us to miss out on treatment-effect dynamics, which are intrinsic to recommender systems
  1. Compliance is often not immediate: typically experimental treatments are rolled out but compliance isn't measured; analyzed as *intent to treat* (ITT)
  2. Timing of compliance is informative: **activity bias :=** estimates for early compliers might not be externally valid
  3. Dynamic effects are important: **novelty bias :=** estimates of early time periods might not be externally valid
- **this project:** propose scalable panel data estimators that help identify these
- temporal and cohort-level granularity is informative and important - don't flatten them with a T-test

Raw outcomes for 30 units from four DGPs
Cross-sectional ATE=0 for all four

1: Constant Effect, 1 adoption cohort

2: Heterogeneous Across Units, 1 adoption cohort

3: Heterogeneous Over Time, 1 adoption cohort

4: Heterogeneous Across Units and Time, 2 adoption cohorts

Outcome

Time

$$\widehat{\boldsymbol{\beta}} = \left( \underbrace{\mathbf{X}^\top}_{P \times N} \underbrace{\mathbf{X}}_{N \times P} \right)^{-1} \underbrace{\mathbf{X}^\top}_{P \times N} \underbrace{\mathbf{y}}_{N \times 1}$$

$$\mathbb{V}\left[\widehat{\beta}\right] = \left(\mathbf{X}^\top \mathbf{X}\right)^{-1} \mathbf{X}^\top \Omega \mathbf{X} \left(\mathbf{X}^\top \mathbf{X}\right)^{-1}$$

$$\widehat{\boldsymbol{\beta}} = \left( \underbrace{\mathbf{X}^\top}_{P \times N} \underbrace{\mathbf{X}}_{N \times P} \right)^{-1} \underbrace{\mathbf{X}^\top}_{P \times N} \underbrace{\mathbf{y}}_{N \times 1}$$

$$\mathbb{V}\left[ \widehat{\beta} \right] = \left( \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \Omega \mathbf{X} \left( \mathbf{X}^\top \mathbf{X} \right)^{-1}$$

- **X** includes intercept and binary treatment indicator $W_i$
- Suppose **X** takes on values in $\mathcal{X}$ with finite cardinality $C$
- Then, we compute sufficient statistics by strata, run $\widetilde{\mathbf{y}}/\widetilde{\mathbf{n}} \sim \widetilde{\mathbf{X}}\beta + \varepsilon$: $C$ observations

$$\widehat{\boldsymbol{\beta}} = \left( \underbrace{\mathbf{X}^\top}_{P \times N} \underbrace{\mathbf{X}}_{N \times P} \right)^{-1} \underbrace{\mathbf{X}^\top}_{P \times N} \underbrace{\mathbf{y}}_{N \times 1}$$

$$\mathbb{V}\left[ \widehat{\boldsymbol{\beta}} \right] = \left( \mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \Omega \mathbf{X} \left( \mathbf{X}^\top \mathbf{X} \right)^{-1}$$

- $\mathbf{X}$ includes intercept and binary treatment indicator $W_i$
- Suppose $\mathbf{X}$ takes on values in $\mathcal{X}$ with finite cardinality $C$
- Then, we compute sufficient statistics by strata, run $\widetilde{\mathbf{y}}/\widetilde{\mathbf{n}} \sim \widetilde{\mathbf{X}}\boldsymbol{\beta} + \varepsilon$: $C$ observations

Table 1: Example dataset and its compressed versions.

| (a) | | (b) | | | (c) | | | (d) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **M** | **y** | **M̂** | **ẏ** | **ṅ** | **M̄** | **ȳ** | **n̄** | **M̃** | **ỹ'** | **ỹ''** | **ñ** |
| A | 1 | A | 1 | 2 | A | 1.33 | 3 | A | 4 | 6 | 3 |
| A | 1 | A | 2 | 1 | B | 3.5 | 2 | B | 7 | 25 | 2 |
| A | 2 | B | 3 | 1 | C | 5 | 1 | C | 5 | 25 | 1 |
| B | 3 | B | 4 | 1 | | | | | | | |
| B | 4 | C | 5 | 1 | | | | | | | |
| C | 5 | | | | | | | | | | |

(a) Uncompressed data. (b) f-weights: $(\mathbf{y}, \mathbf{M})$-compressed records.
(c) Groups: $(\mathbf{M})$-compressed records. (d) Sufficient Statistics: $(\mathbf{M})$-compressed records.

ABlaze is built around this: Wong et al(2021)

# Compressing Least Squares

$$\widehat{\boldsymbol{\beta}} = \Big(\underbrace{\mathbf{X}^\top}_{P \times N}\underbrace{\mathbf{X}}_{N \times P}\Big)^{-1}\underbrace{\mathbf{X}^\top}_{P \times N}\underbrace{\mathbf{y}}_{N \times 1}$$

$$\mathbb{V}\big[\widehat{\boldsymbol{\beta}}\big] = \big(\mathbf{X}^\top\mathbf{X}\big)^{-1}\mathbf{X}^\top\Omega\mathbf{X}\big(\mathbf{X}^\top\mathbf{X}\big)^{-1}$$

- **X** includes intercept and binary treatment indicator $W_i$
- Suppose **X** takes on values in $\mathcal{X}$ with finite cardinality $C$
- Then, we compute sufficient statistics by strata, run $\widetilde{\mathbf{y}}/\widetilde{\mathbf{n}} \sim \widetilde{\mathbf{X}}\boldsymbol{\beta} + \varepsilon$: $C$ observations

Table 1: Example dataset and its compressed versions.

| (a) | | | (b) | | | | (c) | | | | (d) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **M** | **y** | | **M̀** | **ẏ** | **ṅ** | | **M̄** | **ȳ** | **n̄** | | **M̄** | **ỹ'** | **ỹ''** | **n̄** |
| A | 1 | | A | 1 | 2 | | A | 1.33 | 3 | | A | 4 | 6 | 3 |
| A | 1 | | A | 2 | 1 | | B | 3.5 | 2 | | B | 7 | 25 | 2 |
| A | 2 | | B | 3 | 1 | | C | 5 | 1 | | C | 5 | 25 | 1 |
| B | 3 | | B | 4 | 1 | | | | | | | | | |
| B | 4 | | C | 5 | 1 | | | | | | | | | |
| C | 5 | | | | | | | | | | | | | |

(a) Uncompressed data. (b) f-weights: $(\mathbf{y}, \mathbf{M})$-compressed records.

(c) Groups: $(\mathbf{M})$-compressed records. (d) Sufficient Statistics: $(\mathbf{M})$-compressed records.

ABlaze is built around this: Wong et al(2021)

$$\widehat{\mathbb{V}}(\boldsymbol{\beta}) = \overbrace{(\widetilde{\mathbf{X}}^\top\text{diag}(\widetilde{\mathbf{n}})\widetilde{\mathbf{X}})^{-1}}^{\text{Bread}}\ \overbrace{\widetilde{\mathbf{X}}^\top\text{diag}(\sum_j \underbrace{(\widetilde{y}_j^2\widetilde{n}_j - 2\widetilde{y}_j\widetilde{y}_j + \widetilde{y}_j'')}_{\text{RSS in } j}))\widetilde{\mathbf{X}}}^{\text{Meat}}\ \overbrace{(\widetilde{\mathbf{X}}^\top\text{diag}(\widetilde{\mathbf{n}})\widetilde{\mathbf{X}})^{-1}}^{\text{Bread}}$$

- Base untreated potential outcome : $Y_{it}^0 = \alpha_i + \gamma_t + \epsilon_{it}$
- Treated potential outcome under static, constant effects: $Y_{it}^1 = Y_{it}^0 + \tau$

# Building Blocks

- Base untreated potential outcome : $Y_{it}^0 = \alpha_i + \gamma_t + \epsilon_{it}$
- Treated potential outcome under static, constant effects: $Y_{it}^1 = Y_{it}^0 + \tau$
- Treated potential outcome under unit heterogeneity: $Y_{it}^1 = Y_{it}^0 + \tau_i$
- Treated potential outcome under time heterogeneity: $Y_{it}^1 = Y_{it}^0 + \sum_{k \geq 0} \tau_k 1\{t - T_i = k\}$

# Building Blocks

- Base untreated potential outcome : $Y_{it}^0 = \alpha_i + \gamma_t + \epsilon_{it}$
- Treated potential outcome under static, constant effects: $Y_{it}^1 = Y_{it}^0 + \tau$
- Treated potential outcome under unit heterogeneity: $Y_{it}^1 = Y_{it}^0 + \tau_i$
- Treated potential outcome under time heterogeneity: $Y_{it}^1 = Y_{it}^0 + \sum_{k \geq 0} \tau_k 1\{t - T_i = k\}$
- Treated potential outcome under time heterogeneity: $Y_{it}^1 = Y_{it}^0 + \sum_{k \geq 0} \tau_{ik} 1\{t - T_i = k\}$
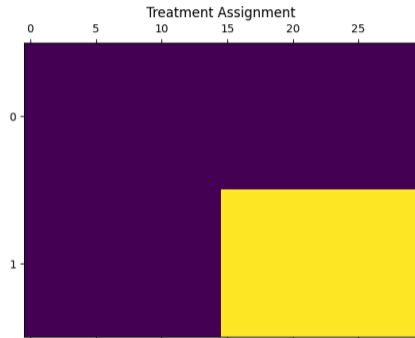
- Base untreated potential outcome : $Y_{it}^0 = \alpha_i + \gamma_t + \epsilon_{it}$
- Treated potential outcome under static, constant effects: $Y_{it}^1 = Y_{it}^0 + \tau$
- Treated potential outcome under unit heterogeneity: $Y_{it}^1 = Y_{it}^0 + \tau_i$
- Treated potential outcome under time heterogeneity: $Y_{it}^1 = Y_{it}^0 + \sum_{k \geq 0} \tau_k 1\{t - T_i = k\}$
- Treated potential outcome under time heterogeneity: $Y_{it}^1 = Y_{it}^0 + \sum_{k \geq 0} \tau_{ik} 1\{t - T_i = k\}$
- Observed outcome:

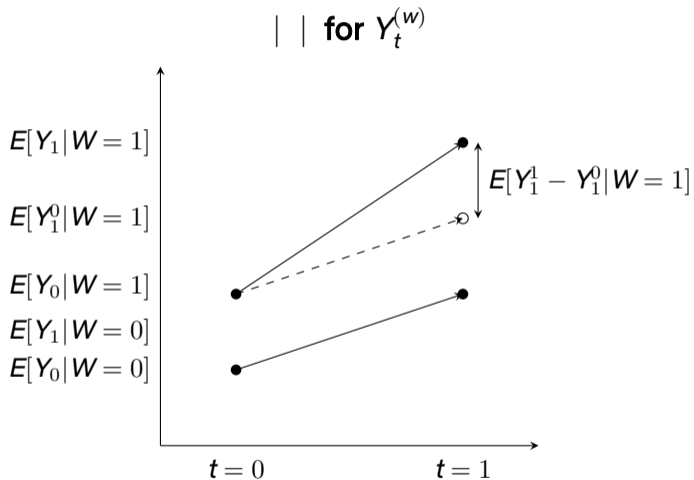$$Y_{it} = Y_{it}^0 (1 - W_{it}) + Y_{it}^1 W_{it}$$

- Under random assignment or parallel trends, we can form estimates of $Y_{it}^0$ and construct estimates of (reasonable averages of) $\tau_{it}$



When assignment time is endogenous, need generalized propensity score Arkhangelsky et al (2024)
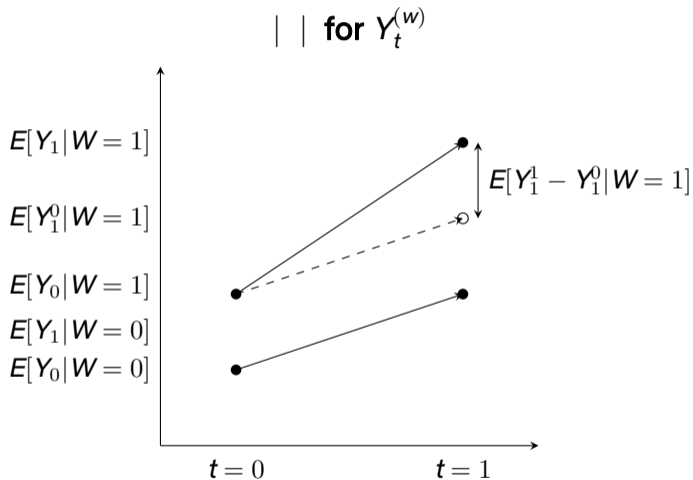
$| \ |$ **for** $Y_t^{(w)}$

$E[Y_1|W=1]$

$E[Y_1^0|W=1]$

$E[Y_1^1 - Y_1^0|W=1]$

$E[Y_0|W=1]$
$E[Y_1|W=0]$
$E[Y_0|W=0]$

$t=0$ $\qquad$ $t=1$

$$| \; | \text{ for } Y_t^{(w)}$$



Easier to satisfy in experiments since $\mathbb{E}\left[Y_0 \mid W = 1\right] = \mathbb{E}\left[Y_0 \mid W = 0\right]$ under random assignment

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \varepsilon_i \qquad\qquad\qquad \textbf{Diff in Means}$$

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \beta \bar{Y}_{i,t<T_0} + \varepsilon_i \qquad\qquad \textbf{CUPED}$$

- $Z_{it} := \mathbb{1}\{\arg\min\{t : W_{it} = 1\} - t = k\}$ (treatment indicator $\times$ event-time)

# Regressions we'd like to run

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \varepsilon_i \qquad \text{\textbf{Diff in Means}}$$

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \beta \bar{Y}_{i,t<T_0} + \varepsilon_i \qquad \text{\textbf{CUPED}}$$

$$Y_{it} = \alpha_i + \gamma_t + \tau W_{it} + \varepsilon_{it} \qquad \text{\textbf{Static TWFE}}$$

$$Y_{it} = \alpha_i + \gamma_t + \sum_{k=0, k\neq -1}^{T} \tau_k Z_{it}^k + \varepsilon_{it} \qquad \text{\textbf{Dynamic TWFE (Event Study)}}$$

$$Y_{it} = \alpha + \gamma_t + \sum_{k\geq 0}^{T} \tau_k Z_{it}^k + \varepsilon_{it} \qquad \text{\textbf{Dynamic DiM}}$$

$$Y_{it} = \alpha_i + \gamma_t + \sum_{c=2}^{C} \sum_{k=0, k\neq -1}^{T} \mathbb{1}\{G_i = c\}\tau_{kc} Z_{it}^k + \varepsilon_{it} \qquad \text{\textbf{Disagg Dynamic TWFE (Staggered Event Study)}}$$

■ $Z_{it} := \mathbb{1}\{\text{argmin}\{t : W_{it} = 1\} - t = k\}$ (treatment indicator × event-time)

## Regressions we'd like to run

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \varepsilon_i \qquad \textbf{\textcolor{blue}{Diff in Means}}$$

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \beta \bar{Y}_{i,t<T_0} + \varepsilon_i \qquad \textbf{\textcolor{blue}{CUPED}}$$

$$Y_{it} = \alpha_i + \gamma_t + \tau W_{it} + \varepsilon_{it} \qquad \textbf{Static TWFE}$$

$$Y_{it} = \alpha_i + \gamma_t + \sum_{k=0, k\neq -1}^{T} \tau_k Z_{it}^k + \varepsilon_{it} \qquad \textbf{Dynamic TWFE (Event Study)}$$

$$Y_{it} = \alpha + \gamma_t + \sum_{k\geq 0}^{T} \tau_k Z_{it}^k + \varepsilon_{it} \qquad \textbf{\textcolor{cyan}{Dynamic DiM}}$$

$$Y_{it} = \alpha_i + \gamma_t + \sum_{c=2}^{C} \sum_{k=0, k\neq -1}^{T} \mathbb{1}\{G_i = c\} \tau_{kc} Z_{it}^k + \varepsilon_{it} \qquad \textbf{Disagg Dynamic TWFE (Staggered Event Study)}$$

- $Z_{it} := \mathbb{1}\{\text{argmin}\{t : W_{it} = 1\} - t = k\}$ (treatment indicator × event-time)
- $N \times T \approx 50m \times 95 = 4.75\text{b obs}$

# Regressions we'd like to run

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \varepsilon_i \qquad \text{**Diff in Means**}$$

$$\bar{Y}_{i,t>T_0} = \alpha + \tau W_i + \beta \bar{Y}_{i,t<T_0} + \varepsilon_i \qquad \text{**CUPED**}$$

$$Y_{it} = \alpha_i + \gamma_t + \tau W_{it} + \varepsilon_{it} \qquad \text{**Static TWFE**}$$

$$Y_{it} = \alpha_i + \gamma_t + \sum_{k=0, k \neq -1}^{T} \tau_k Z_{it}^k + \varepsilon_{it} \qquad \text{**Dynamic TWFE (Event Study)**}$$

$$Y_{it} = \alpha + \gamma_t + \sum_{k \geq 0}^{T} \tau_k Z_{it}^k + \varepsilon_{it} \qquad \text{**Dynamic DiM**}$$

$$Y_{it} = \alpha_i + \gamma_t + \sum_{c=2}^{C} \sum_{k=0, k \neq -1}^{T} \mathbb{1}\{G_i = c\} \tau_{kc} Z_{it}^k + \varepsilon_{it} \qquad \text{**Disagg Dynamic TWFE (Staggered Event Study)**}$$

- $Z_{it} := \mathbb{1}\{\arg\min\{t : W_{it} = 1\} - t = k\}$ (treatment indicator × event-time)
- $N \times T \approx 50m \times 95 = 4.75b$ obs
- $W_i, \alpha_i, \gamma_t$ jointly identify a single observation - cannot compress
- Unit intercepts $\alpha_i$: millions of distinct values

- We don't inherently care about $\alpha_i, \gamma_t$; they are nuisance parameters
- Partial them out (i.e. kick them out of $\mathbf{X'X}$)
    - Frisch-Waugh-Lovell Theorem / Gram-Schmidt Process

- We don't inherently care about $\alpha_i, \gamma_t$; they are nuisance parameters
- Partial them out (i.e. kick them out of $\mathbf{X'X}$)
  - Frisch-Waugh-Lovell Theorem / Gram-Schmidt Process
  - specialized methods for high dimensional categorical covariates: Method of alternating projections / Kaczmarz method (implemented in `areg`/`reghdfe`/ `fixest`/`felm`/`pyfixest/...`)
- Residualise RHS $W_{it} - \overline{W}_{i,\cdot} - \overline{W}_{\cdot,t} =: \ddot{W}_{it}$

- We don't inherently care about $\alpha_i, \gamma_t$; they are nuisance parameters
- Partial them out (i.e. kick them out of $\mathbf{X'X}$)
    - Frisch-Waugh-Lovell Theorem / Gram-Schmidt Process
    - specialized methods for high dimensional categorical covariates: Method of alternating projections / Kaczmarz method (implemented in `areg`/`reghdfe`/ `fixest`/`felm`/`pyfixest`/`...`)
- Residualise RHS $W_{it} - \overline{W}_{i,\cdot} - \overline{W}_{\cdot,t} =: \ddot{W}_{it}$
- Obviates the need to adjust for time-invariant member characteristics (all absorbed in FEs)

$$Y_{it} = \tau \ddot{W}_{it} + \varepsilon_{it}$$

This regression *can* be compressed. However, we don't have an equivalent representation for dynamic regressions (event studies, disaggregated event studies), etc.

# Trick 2: Mundlak (1978) + Wooldridge (2021) Trick

$$Y_{it} = \alpha_i + \gamma_t + \tau W_{it} + \varepsilon_{it}$$
$$Y_{it} = \delta + \tau W_{it} + \phi \bar{W}_{i,\cdot} + \psi \bar{W}_{\cdot,t} + \varepsilon_{it}$$

**Static TWFE**

**Mundlak-ed Static TWFE**

$$Y_{it} = \alpha_i + \gamma_t + \tau W_{it} + \varepsilon_{it}$$ **Static TWFE**

$$Y_{it} = \delta + \tau W_{it} + \phi \bar{W}_{i,\cdot} + \psi \bar{W}_{\cdot,t} + \varepsilon_{it}$$ **Mundlak-ed Static TWFE**

$$Y_{it} = \alpha_i + \gamma_t + \sum_{k=0, k\neq -1}^{T} \tau_k Z_{it}^* + \varepsilon_{it}$$ **Dynamic TWFE**

$$Y_{it} = \delta + \psi D_i + \sum_{k=1}^{T} \phi_t \mathbb{1}\{t = k\} + \sum_{k=1}^{T} \tau_k D_i \mathbb{1}\{t = k\} + \varepsilon_{it}$$

# Trick 2: Mundlak (1978) + Wooldridge (2021) Trick

$$Y_{it} = \alpha_i + \gamma_t + \tau W_{it} + \varepsilon_{it}$$ **Static TWFE**

$$Y_{it} = \delta + \tau W_{it} + \phi \bar{W}_{i,\cdot} + \psi \bar{W}_{\cdot,t} + \varepsilon_{it}$$ **Mundlak-ed Static TWFE**

$$Y_{it} = \alpha_i + \gamma_t + \sum_{k=0, k \neq -1}^{T} \tau_k Z_{it}^{\star} + \varepsilon_{it}$$ **Dynamic TWFE**

$$Y_{it} = \delta + \psi D_i + \sum_{k=1}^{T} \phi_t \mathbb{1}\{t = k\} + \sum_{k=1}^{T} \tau_k D_i \mathbb{1}\{t = k\} + \varepsilon_{it}$$

$$Y_{it} = \alpha_i + \gamma_t + \sum_{c=2}^{C} \sum_{k=0, k \neq -1}^{T} \mathbb{1}\{G_i = c\} \tau_{kc} Z_{it}^{\star} + \varepsilon_{it}$$ **Disagg Dynamic TWFE**

$$Y_{it} = \delta + \underbrace{\sum_{c=1}^{C} \psi_c \mathbb{1}\{D_i = k\}}_{\text{Cohort Dummies}} + \underbrace{\sum_{k=1}^{T} \phi_t \mathbb{1}\{t = k\}}_{\text{Time Dummies}} + \underbrace{\sum_{c=1}^{C} \sum_{k=1}^{T} \tau_{kc} \mathbb{1}\{D_i = c\} \mathbb{1}\{t = k\}}_{\text{Cohort} \times \text{Time interactions}} + \varepsilon_{it}$$
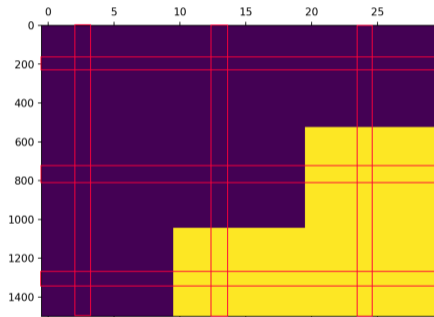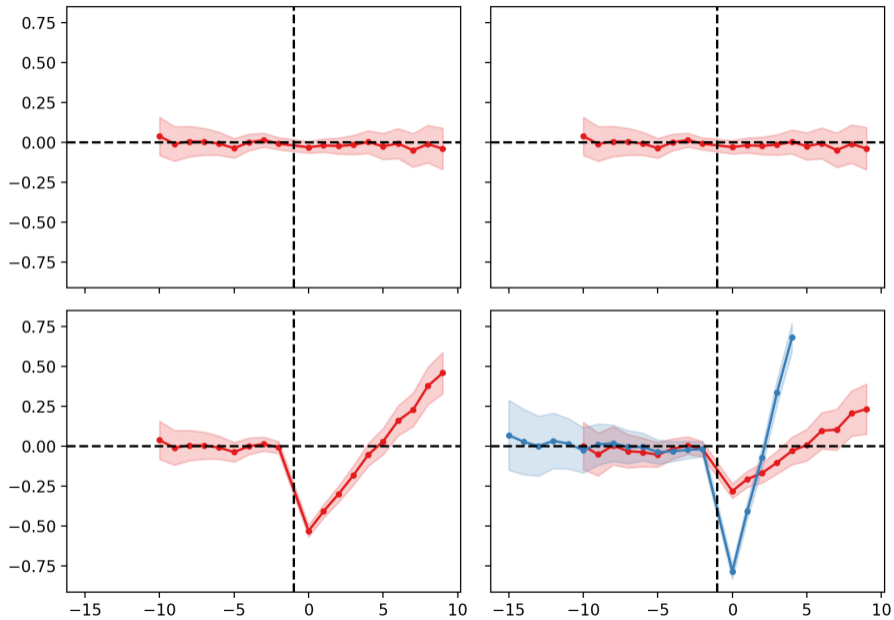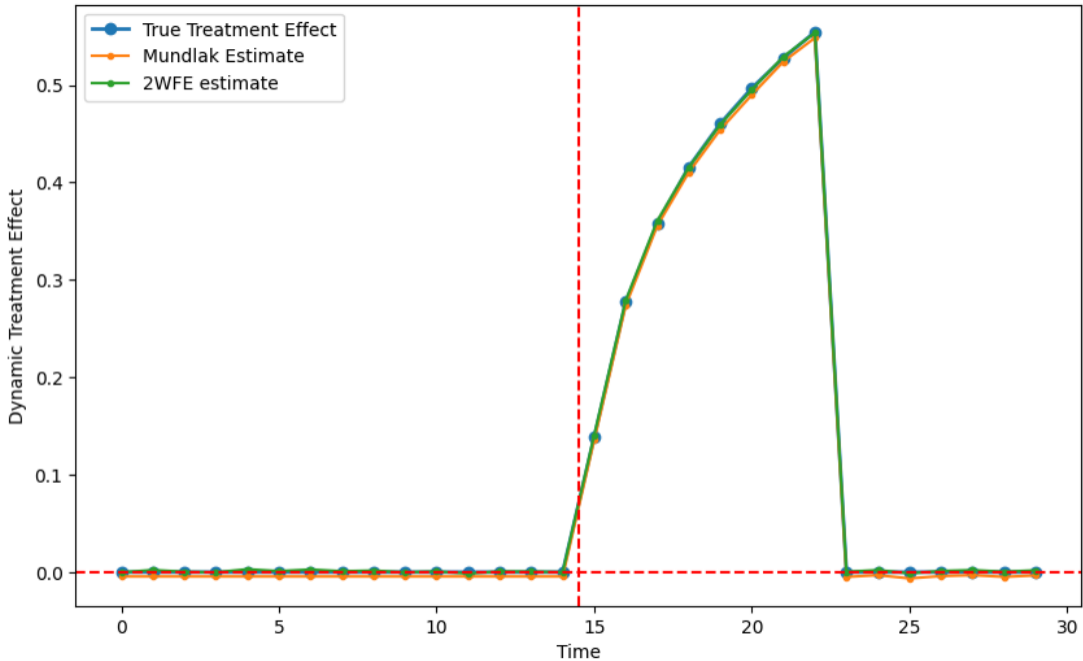
# Design Matrix Dimensions

| | **(1) Standard** | $M$ | **(2) Mundlak** | $\tilde{M}$ |
|---|---|---|---|---|
| Static | $Y_{it} = \alpha_i + \gamma_t + \tau W_{it} + \varepsilon_{it}$ | NT | $Y_{it} = \alpha + \tau W_{it} + \psi \overline{W}_{i,\cdot}$ $+ \phi \overline{W}_{\cdot,t} + \varepsilon_{it}$ | 2+(C+1) |
| Dyn | $Y_{it} = \alpha_i + \gamma_t$ $+ \sum_{k \neq -1} \tau_k Z_{it}^k + \varepsilon_{it}$ | NT | $Y_{it} = \alpha + \psi D_i + \sum_{k=1}^{T} \phi_t 1_{t=k}$ $+ \sum_{k=1}^{T} \tau_k D_i 1_{t=k} + \varepsilon_{it}$ | 2T |
| Dyn+Stagg | $Y_{it} = \alpha_i + \gamma_t$ $+ \sum_{c=1}^{C} \sum_{k \neq -1} \tau_{kc} 1_{G_i=c} Z_{it}^k$ $+ \varepsilon_{it}$ | NT | $Y_{it} = \alpha + \sum_{c=1}^{C} \psi_c 1_{D_i=c}$ $+ \sum_{k=1}^{T} \phi_t 1_{t=k}$ $+ \sum_{c=1}^{C} \sum_{k=1}^{T} \tau_{kc} 1_{D_i=c} 1_{t=k}$ $+ \varepsilon_{it}$ | CT |



- $N$ units, $T$ time periods
- $M$, $\tilde{M}$: number of required rows in design matrix
- $C$ unique adoption cohorts (including control)
- $2 + (C + 1) = 4$ obs in standard A/B test
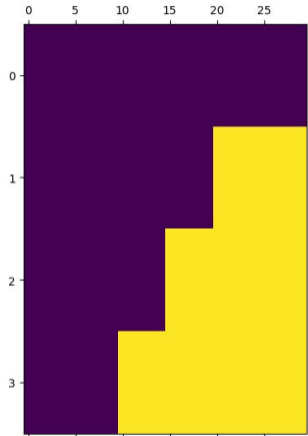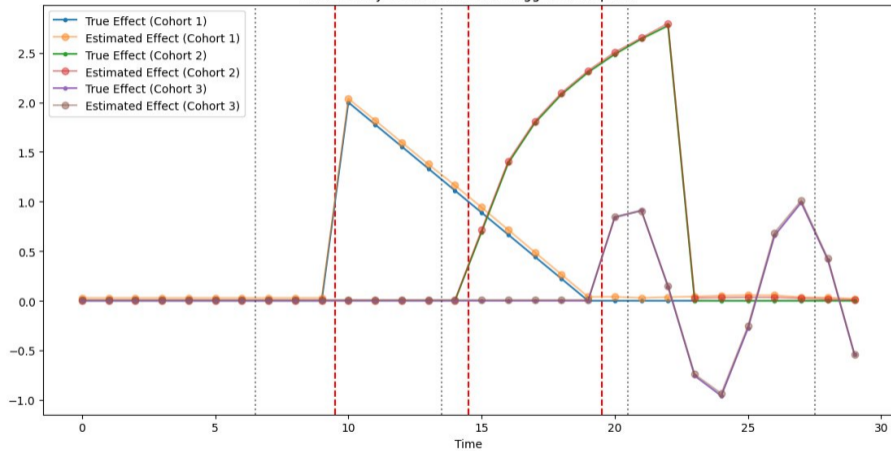- Cluster-robust inference with cluster-bootstrap

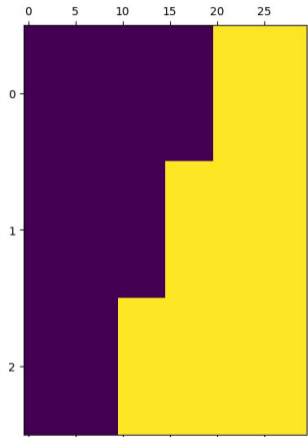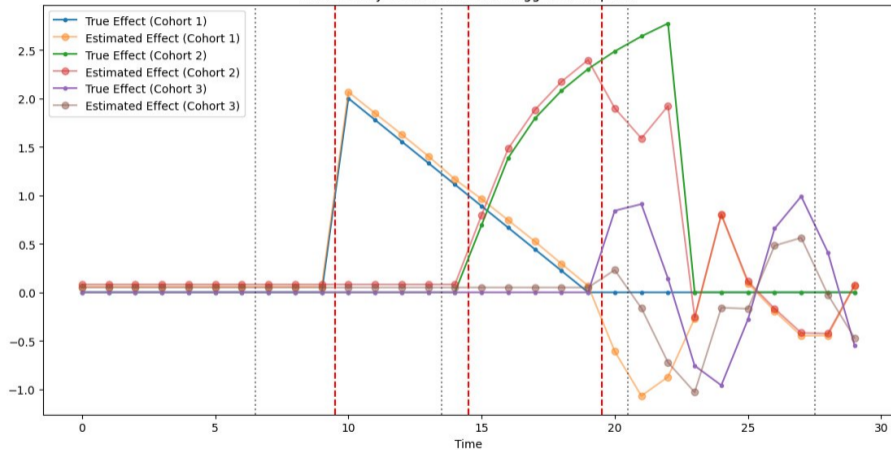Event Study figures for the four scenarios
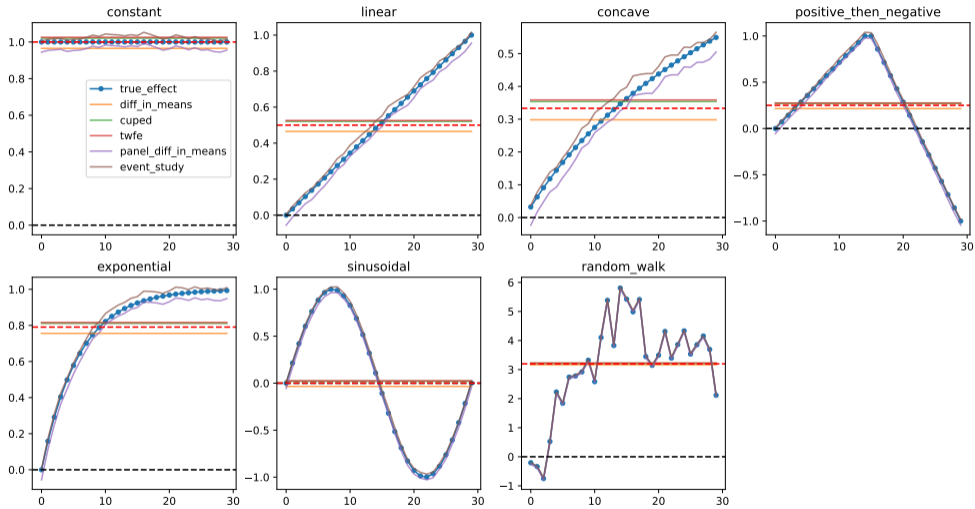
Event Study estimates under staggered adoption

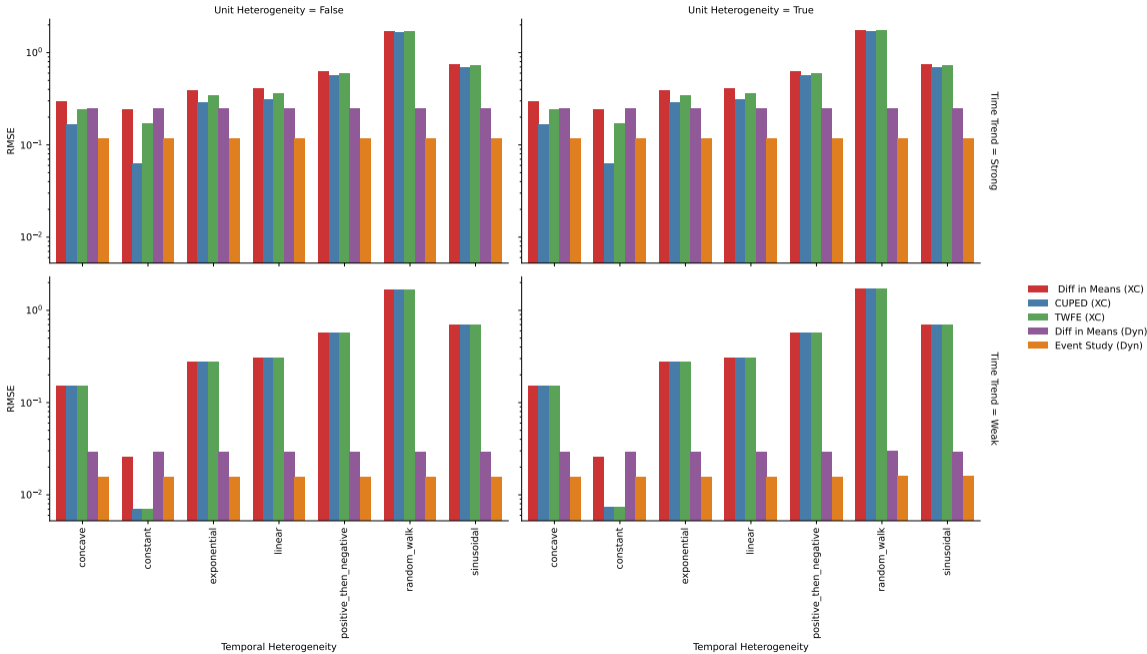Event Study estimates under staggered adoption

# Numerical Experiments

Estimation Accuracy for Different forms of Temporal Heterogeneity

RMSE By (Unit | Time) Heterogeneity and Misspecification

# Runtime Comparisons

$1000 \times \{14, 28, 42, 140, 280, 420\}$ units, $\{14, 28, 42\}$ days.

| Observations | Units | Periods | duckreg | pyfixest | pyfixest compressed | statsmodels |
|---|---|---|---|---|---|---|
| 14K | 1K | 14 | 0.03 | 0.18 | 0.17 | 2.75 |
| 28K | 1K | 28 | 0.03 | 0.19 | 0.18 | 6.27 |
| 42K | 1K | 42 | 0.04 | 0.20 | 0.21 | 9.43 |
| 140K | 10K | 14 | 0.05 | 0.22 | 0.26 | x |
| 280K | 10K | 28 | 0.06 | 0.26 | 0.36 | x |
| 420K | 10K | 42 | 0.04 | 0.31 | 0.47 | x |
| 1M | 100K | 14 | 0.07 | 0.64 | 1.27 | x |
| 3M | 100K | 28 | 0.21 | 1.00 | 2.28 | x |
| 4M | 100K | 42 | 0.24 | 1.41 | 3.43 | x |
| 14M | 1M | 14 | 1.03 | 4.88 | 22.12 | x |
| 28M | 1M | 28 | 3.41 | 8.07 | 62.07 | x |
| 42M | 1M | 42 | 10.92 | 13.60 | 117.63 | x |
| 140M | 10M | 14 | 19.70 | 123.88 | x | x |
| Mundlak | | | ✓ | - | ✓ | - |
| Compression | | | ✓ | - | ✓ | - |
| Out-of-Memory | | | ✓ | - | - | - |

# Extension: imbalanced panels, descriptive analysis - AKM Regressions

$$\underbrace{Y_{it}}_{\text{Stickiness}} \sim \underbrace{\alpha_i}_{\text{Member FE}} + \underbrace{\psi_{\mathbf{J}(i,t)}}_{\text{Title-FE}} + \underbrace{x'_{it}\beta}_{\text{Covariates}} + \varepsilon_{it}$$

# Extension: imbalanced panels, descriptive analysis - AKM Regressions

$$\underbrace{Y_{it}}_{\text{Stickiness}} \sim \underbrace{\alpha_i}_{\text{Member FE}} + \underbrace{\psi_{\mathbf{J}(i,t)}}_{\text{Title-FE}} + \underbrace{x'_{it}\beta}_{\text{Covariates}} + \varepsilon_{it}$$

- Classic tool in Labour Economics (Abowd et al 1999)
- $\alpha_i$ member FE: unit $i$'s baseline completion metric
- $J(i, t)$: $i$ watched $j$ at time $t$
- $\psi_j$ title FE: title $j$'s completion metric
- Requires connected user-title graph - plausible

$$Y_{it} \sim \underbrace{\alpha_i}_{\text{Member FE}} + \underbrace{\psi_{\mathbf{J}(i,t)}}_{\text{Title-FE}} + \underbrace{X'_{it}\beta}_{\text{Covariates}} + \varepsilon_{it}$$

$$\underbrace{Y_{it}}_{\text{Stickiness}}$$

- Classic tool in Labour Economics (Abowd et al 1999)
- $\alpha_i$ member FE: unit $i$'s baseline completion metric
- $J(i, t)$: $i$ watched $j$ at time $t$
- $\psi_j$ title FE: title $j$'s completion metric
- Requires connected user-title graph - plausible
- Admits to variance decomposition

$$\mathbb{V}\left[Y_{it} - X'_{it}\beta\right] = \underbrace{\mathbb{V}\left[\alpha_i\right]}_{\text{Member effects}} + \underbrace{\mathbb{V}\left[\psi_{j(i,t)}\right]}_{\text{Title Effects}} + \underbrace{\mathbb{V}\left[\mathsf{Cov}\left[\alpha_i, \psi_{j(i,t)}\right]\right]}_{\text{Sorting}} + \mathbb{V}\left[\varepsilon_{it}\right]$$

Naive estimator has problems, fixes use Jackknife (Kline et al 2020).

$$Y_{it} \underbrace{\sim}_{\text{Stickiness}} \underbrace{\alpha_i}_{\text{Member FE}} + \underbrace{\psi_{\mathbf{J}(i,t)}}_{\text{Title-FE}} + \underbrace{x'_{it}\beta}_{\text{Covariates}} + \varepsilon_{it}$$

- Classic tool in Labour Economics (Abowd et al 1999)
- $\alpha_i$ member FE: unit $i$'s baseline completion metric
- $J(i, t)$: $i$ watched $j$ at time $t$
- $\psi_j$ title FE: title $j$'s completion metric
- Requires connected user-title graph - plausible
- Admits to variance decomposition

$$\mathbb{V}\left[Y_{it} - X'_{it}\beta\right] = \underbrace{\mathbb{V}\left[\alpha_i\right]}_{\text{Member effects}} + \underbrace{\mathbb{V}\left[\psi_{j(i,t)}\right]}_{\text{Title Effects}} + \underbrace{\mathbb{V}\left[\text{Cov}\left[\alpha_i, \psi_{j(i,t)}\right]\right]}_{\text{Sorting}} + \mathbb{V}\left[\varepsilon_{it}\right]$$

Naive estimator has problems, fixes use Jackknife (Kline et al 2020). Variance components suggest different catalogue / recommendation strategies.

- Methods + Software for compressed, out-of-memory computation of commonly used least-squares panel data estimators
  - `duckreg` : powered by duckDB - out of memory
  - `pyfixest`: general purpose in-memory regression package, compression in `Polars`

# Extensions, Conclusion

- Methods + Software for compressed, out-of-memory computation of commonly used least-squares panel data estimators
  - `duckreg`: powered by duckDB - out of memory
  - `pyfixest`: general purpose in-memory regression package, compression in `Polars`
- Applicable to both scalable estimation of event studies in experiments (focus of this talk), and difference-in-differences methods in observational settings
- Applies to any saturated regression, balanced-or-unbalanced panel data, or any high-dimensional categorical covariate

- Methods + Software for compressed, out-of-memory computation of commonly used least-squares panel data estimators
  - duckreg : powered by duckDB - out of memory
  - pyfixest: general purpose in-memory regression package, compression in Polars
- Applicable to both scalable estimation of event studies in experiments (focus of this talk), and difference-in-differences methods in observational settings
- Applies to any saturated regression, balanced-or-unbalanced panel data, or any high-dimensional categorical covariate
- Extensions / Links
  - Extension to GLMs: Lumley (2018) - estimate MLE on subsample, perform one-step Fisher-scoring update
  - Sketching methods - randomized linear algebra tricks to approx $X'X$ 'well': Mahoney (2013), Pilanci et al (2018), Dobriban et al (2023)
- Other ideas?