

Panel Regression: Fixed Effect and Demeaned Regression

Zahid Asghar

School of Economics, QAU, Islamabad

10/26/22

Types of Data I

- **Cross-sectional data:** compare different individual i 's at same time \bar{t}

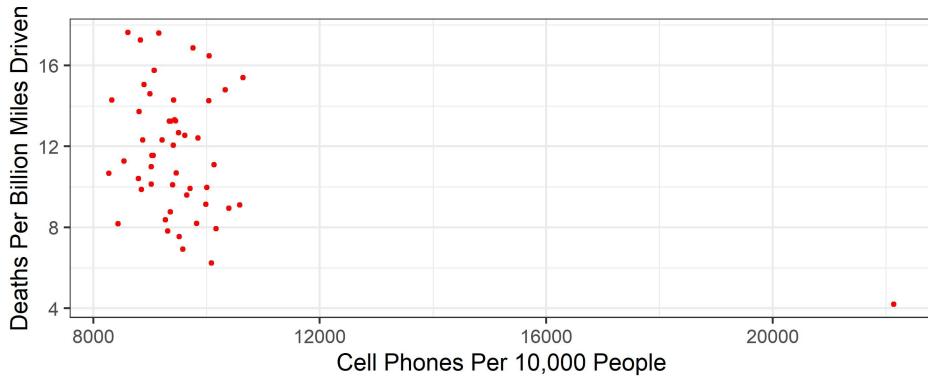
```
## # A tibble: 6 × 4
##   state     year  deaths cell_plans
##   <fct>    <dbl>   <dbl>      <dbl>
## 1 Alabama  2012    13.3     9434.
## 2 Alaska   2012    12.3     8873.
## 3 Arizona  2012    13.7     8811.
## 4 Arkansas 2012    16.5    10047.
## 5 California 2012    8.76    9362.
## 6 Colorado  2012    10.1     9403.
```

- **Time-series data:** track same individual \bar{i} over different times t

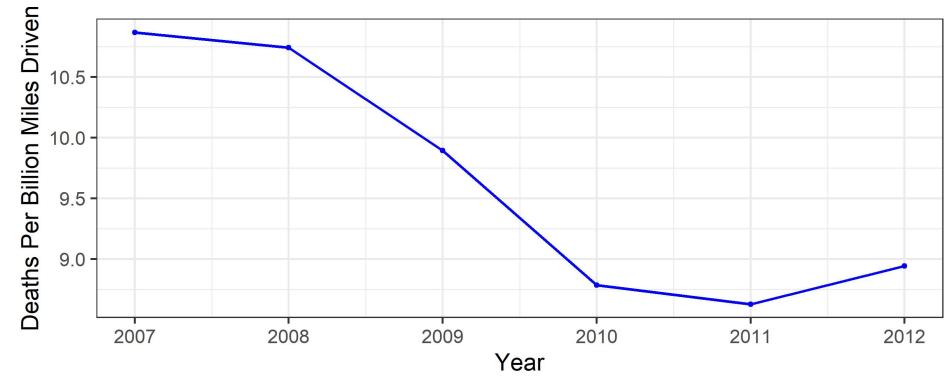
```
## # A tibble: 6 × 4
##   state     year  deaths cell_plans
##   <fct>    <dbl>   <dbl>      <dbl>
## 1 Maryland  2007    10.9     8942.
## 2 Maryland  2008    10.7     9291.
## 3 Maryland  2009    9.89     9339.
## 4 Maryland  2010    8.78     9630.
## 5 Maryland  2011    8.63     10336.
## 6 Maryland  2012    8.94     10393.
```

Types of Data I

- **Cross-sectional data:** compare different individual i 's at same time \bar{t}

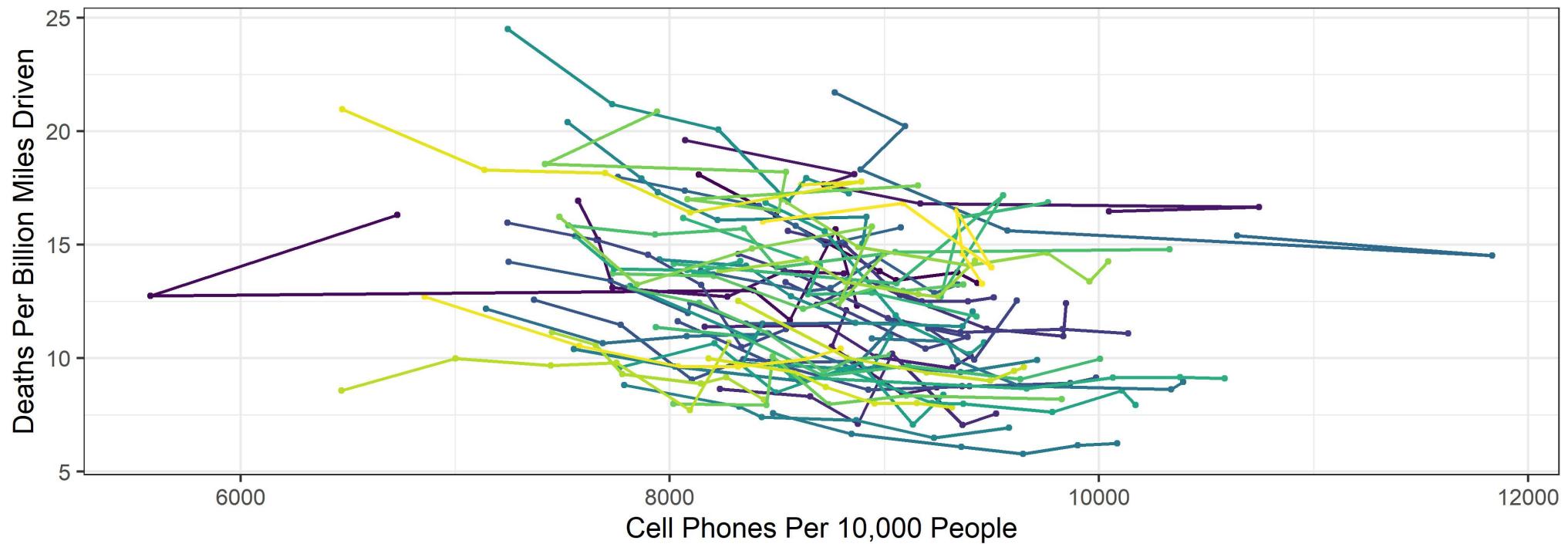


- **Time-series data:** track same individual i over different times t



- **Panel data:** combines these dimensions: compare all individual i 's over all time t 's

Panel Data I



Panel Data II

```
## # A tibble: 306 × 4
##   state    year  deaths cell_plans
##   <fct>   <fct> <dbl>      <dbl>
## 1 Alabama 2007    18.1     8136.
## 2 Alabama 2008    16.3     8494.
## 3 Alabama 2009    13.8     8979.
## 4 Alabama 2010    13.4     9055.
## 5 Alabama 2011    13.8     9341.
## 6 Alabama 2012    13.3     9434.
## 7 Alaska   2007    16.3     6730.
## 8 Alaska   2008    12.7     5581.
## 9 Alaska   2009    13.0     8390.
## 10 Alaska  2010   11.7     8561.
## # ... with 296 more rows
## # i Use `print(n = ...)` to see more rows
```

- **Panel** or **Longitudinal** data contains
 - repeated observations (t)
 - on multiple individuals (i)

Panel Data II

```
## # A tibble: 306 × 4
##   state    year  deaths cell_plans
##   <fct>   <fct> <dbl>      <dbl>
## 1 Alabama 2007    18.1     8136.
## 2 Alabama 2008    16.3     8494.
## 3 Alabama 2009    13.8     8979.
## 4 Alabama 2010    13.4     9055.
## 5 Alabama 2011    13.8     9341.
## 6 Alabama 2012    13.3     9434.
## 7 Alaska   2007    16.3     6730.
## 8 Alaska   2008    12.7     5581.
## 9 Alaska   2009    13.0     8390.
## 10 Alaska  2010   11.7     8561.
## # ... with 296 more rows
## # i Use `print(n = ...)` to see more rows
```

- **Panel or Longitudinal** data contains
 - .repeated observations (t)
 - on multiple individuals (i)
- Thus, our regression equation looks like:

$$\hat{Y}_{it} = \beta_0 + \beta_1 X_{it} + u_{it}$$

| for individual i in time t .

Panel Data: Our Motivating Example

```
## # A tibble: 306 × 4
##   state    year  deaths cell_plans
##   <fct>   <dbl>   <dbl>
## 1 Alabama 2007    18.1    8136.
## 2 Alabama 2008    16.3    8494.
## 3 Alabama 2009    13.8    8979.
## 4 Alabama 2010    13.4    9055.
## 5 Alabama 2011    13.8    9341.
## 6 Alabama 2012    13.3    9434.
## 7 Alaska   2007    16.3    6730.
## 8 Alaska   2008    12.7    5581.
## 9 Alaska   2009    13.0    8390.
## 10 Alaska  2010    11.7    8561.
## # ... with 296 more rows
## # i Use `print(n = ...)` to see more rows
```

Example: Do cell phones cause more traffic fatalities?

- No measure of cell phones *used* while driving
 - `cell_plans` as a **proxy** for cell phone usage
- State-level data over 6 years

The Data I

```
1 glimpse(phones)
2 ## Rows: 306
3 ## Columns: 8
4 ## $ year           <fct> 2007, 2007, 2007, 2007, 2007, 2007, 2007, 200...
5 ## $ state          <fct> Alabama, Alaska, Arizona, Arkansas, California, Col...
6 ## $ urban_percent <dbl> 30, 55, 45, 21, 54, 34, 84, 31, 100, 53, 39, 45, 11...
7 ## $ cell_plans    <dbl> 8136, 6730, 7572, 8071, 8822, 8162, 8235, 8684, 159...
8 ## $ cell_ban       <fct> 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, ...
9 ## $ text_ban       <fct> 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, ...
10 ## $ deaths         <dbl> 18.08, 16.30, 16.93, 19.60, 12.10, 11.37, 8.64, 12...
11 ## $ year_num       <dbl> 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2007, 200...
```

The Data II

```
1 phones %>%
2   count(state)
3 ## # A tibble: 51 × 2
4 ##   state          n
5 ##   <fct>        <int>
6 ## 1 Alabama         6
7 ## 2 Alaska          6
8 ## 3 Arizona         6
9 ## 4 Arkansas         6
10 ## 5 California       6
11 ## 6 Colorado         6
12 ## 7 Connecticut      6
13 ## 8 Delaware         6
14 ## 9 District of Columbia 6
15 ## 10 Florida          6
16 ## # ... with 41 more rows
17 ## # i Use `print(n = ...)` to see more rows
```

```
1 phones %>%
2   count(year)
3 ## # A tibble: 6 × 2
4 ##   year     n
5 ##   <fct> <int>
6 ## 1 2007     51
7 ## 2 2008     51
8 ## 3 2009     51
9 ## 4 2010     51
10 ## 5 2011     51
11 ## 6 2012     51
```

The Data III

```

1 phones %>%
2   distinct(state)
3 ## # A tibble: 51 × 1
4 ##   state
5 ##   <fct>
6 ## 1 Alabama
7 ## 2 Alaska
8 ## 3 Arizona
9 ## 4 Arkansas
10 ## 5 California
11 ## 6 Colorado
12 ## 7 Connecticut
13 ## 8 Delaware
14 ## 9 District of Columbia
15 ## 10 Florida
16 ## # ... with 41 more rows
17 ## # i Use `print(n = ...)` to see more rows

```

```

1 phones %>%
2   distinct(year)
3 ## # A tibble: 6 × 1
4 ##   year
5 ##   <fct>
6 ## 1 2007
7 ## 2 2008
8 ## 3 2009
9 ## 4 2010
10 ## 5 2011
11 ## 6 2012

```

The Data IV

```
1 phones %>%
2   summarize(States = n_distinct(state),
3             Years = n_distinct(year))
4 ## # A tibble: 1 × 2
5 ##   States Years
6 ##     <int> <int>
7 ## 1      51      6
```

The Data: With plm

```
1 # install.packages("plm")
2 library(plm)
3
4 pdim(phones, index=c("state","year"))
5 ## Balanced Panel: n = 51, T = 6, N = 306
```

- **plm package** for panel data in R
- `pdim()` checks dimensions of panel dataset
 - `index=` vector of “group” & “year” variables
- Returns with a summary of:
 - `n` groups
 - `T` periods
 - `N` total observation

Pooled Regression I

- What if we just ran a standard regression:

$$\hat{Y}_{it} = \beta_0 + \beta_1 X_{it} + u_{it}$$

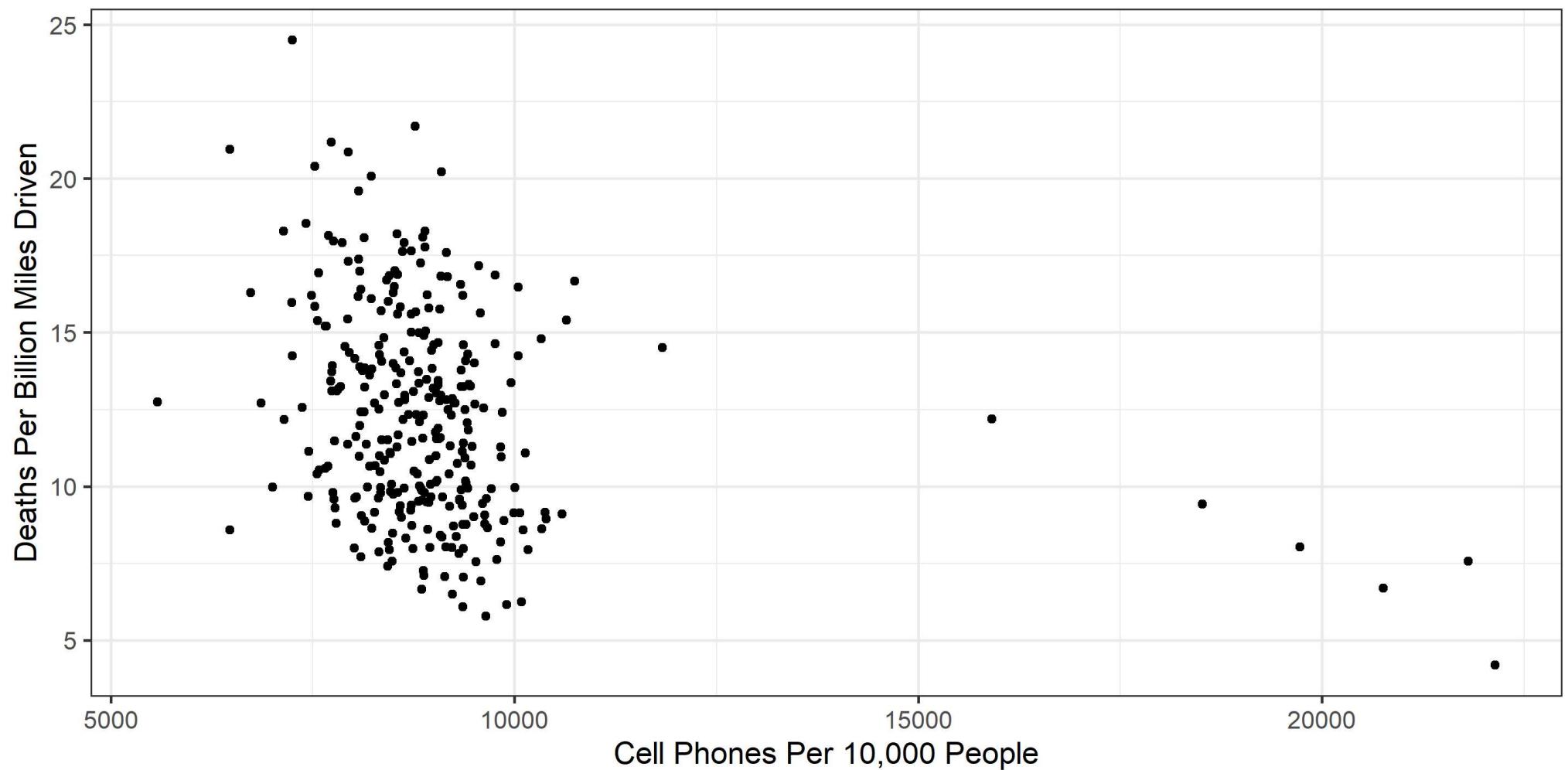
- N number of i groups (e.g. U.S. States)
- T number of t periods (e.g. years)
- This is a pooled regression model: treats all observations as independent

Pooled Regression II

```
1 pooled <- lm(deaths ~ cell_plans, data = phones)
2 pooled %>% tidy()
3 ## # A tibble: 2 × 5
4 ##   term      estimate std.error statistic p.value
5 ##   <chr>     <dbl>     <dbl>     <dbl>    <dbl>
6 ## 1 (Intercept) 17.3      0.975     17.8  5.82e-49
7 ## 2 cell_plans -0.000567  0.000107    -5.30 2.26e- 7
```

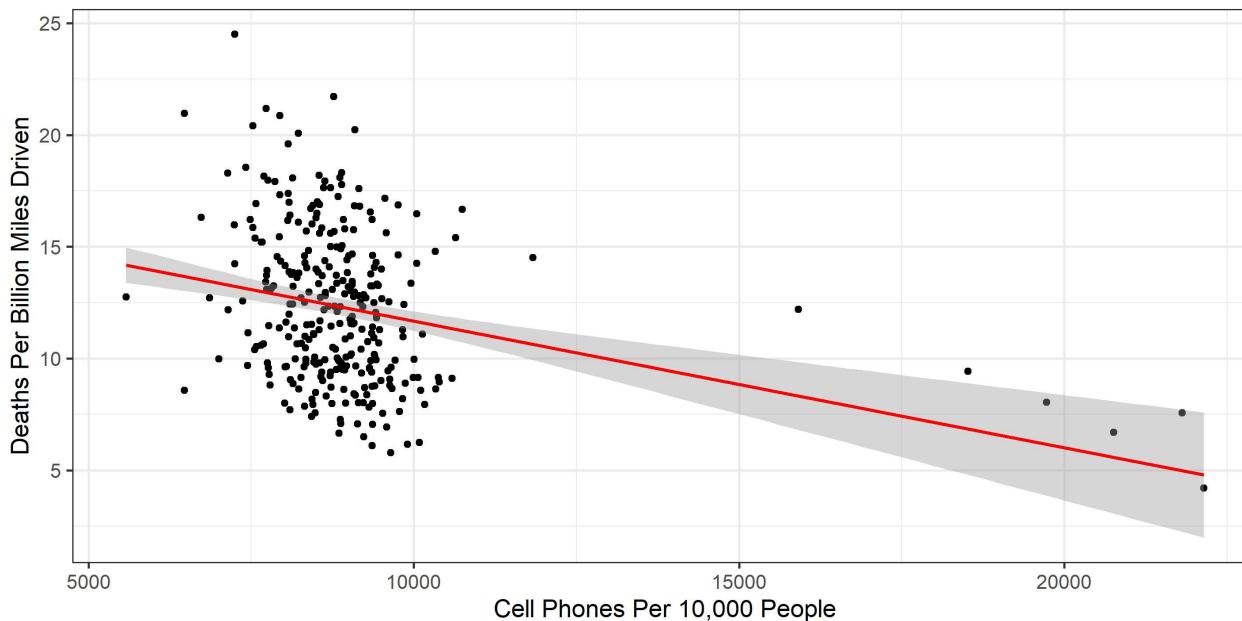
Pooled Regression III

```
1 ggplot(data = phones)+  
2   aes(x = cell_plans,  
3        y = deaths)+  
4   geom_point() +  
5   labs(x = "Cell Phones Per 10,000 People",  
6        y = "Deaths Per Billion Miles Driven") +  
7   theme_bw(base_family = "Fira Sans Condensed",  
8            base_size=14)
```



Pooled Regression III

```
1 ggplot(data = phones)+  
2   aes(x = cell_plans,  
3        y = deaths)+  
4   geom_point() +  
5   geom_smooth(method = "lm", color = "red") + #<<  
6   labs(x = "Cell Phones Per 10,000 People",  
7         y = "Deaths Per Billion Miles Driven") +  
8   theme_bw(base_family = "Fira Sans Condensed",  
9             base_size=14)
```



Zahid Asghar



Recap: Assumptions about Errors

Biases of Pooled Regression

$$\hat{Y}_{it} = \beta_0 + \beta_1 X_{it} + \epsilon_{it}$$

- Assumption 3: $\text{cor}(u_i, u_j) = 0 \quad \forall i \neq j$
- Pooled regression model is **biased** because it ignores:
 - Multiple observations from same group i
 - Multiple observations from same time t
- Thus, errors are serially or auto-correlated; $\text{cor}(u_i, u_j) \neq 0$ within same i and within same t

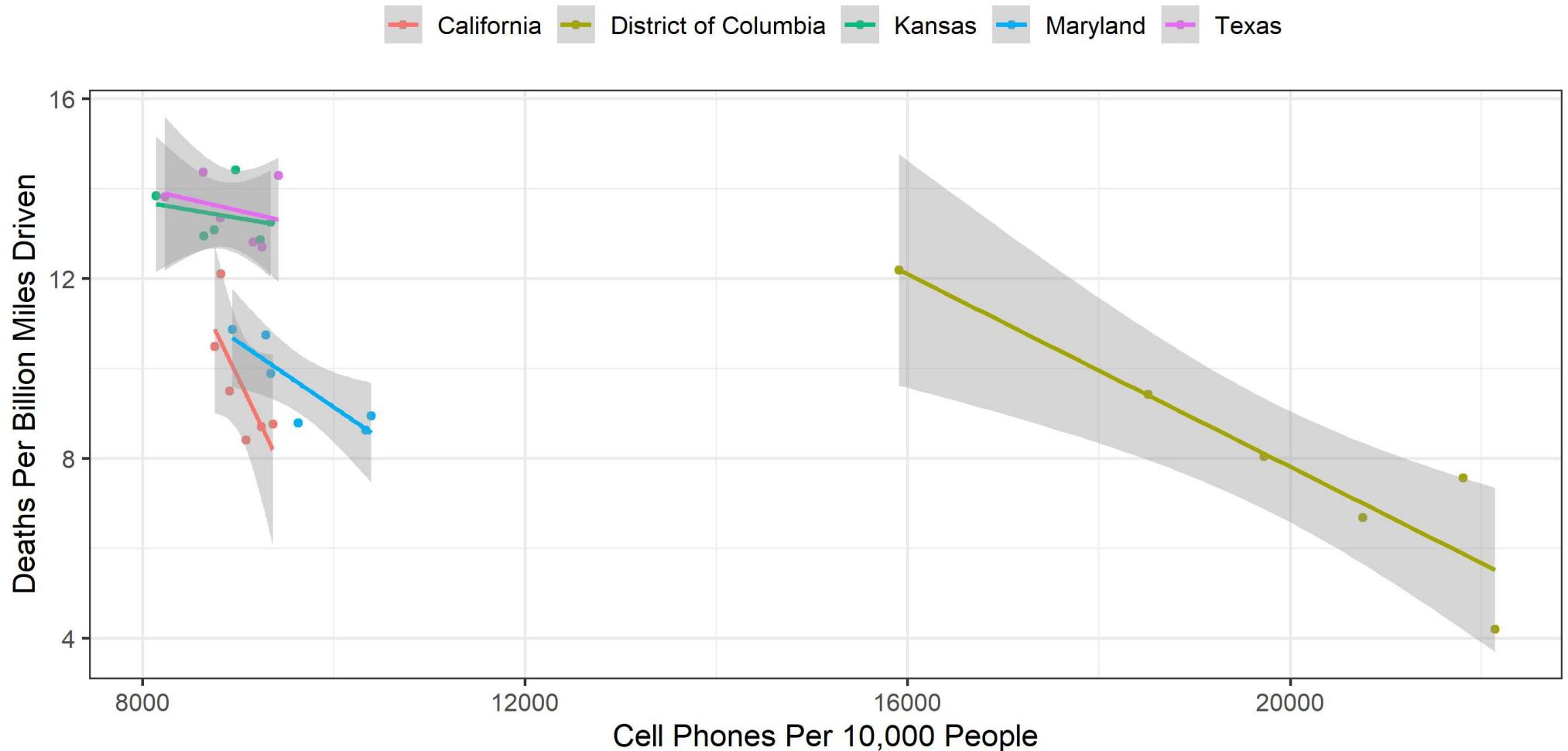
Biases of Pooled Regression: Our Example

$$\widehat{\text{Deaths}}_{it} = \beta_0 + \beta_1 \text{Cell Phones}_{it} + u_{it}$$

- Multiple observations from same state i
 - Probably similarities among u for obs in same state
 - Residuals on observations from same state are likely correlated
- Multiple observations from same year t
 - Probably similarities among u for obs in same year
 - Residuals on observations from same year are likely correlated

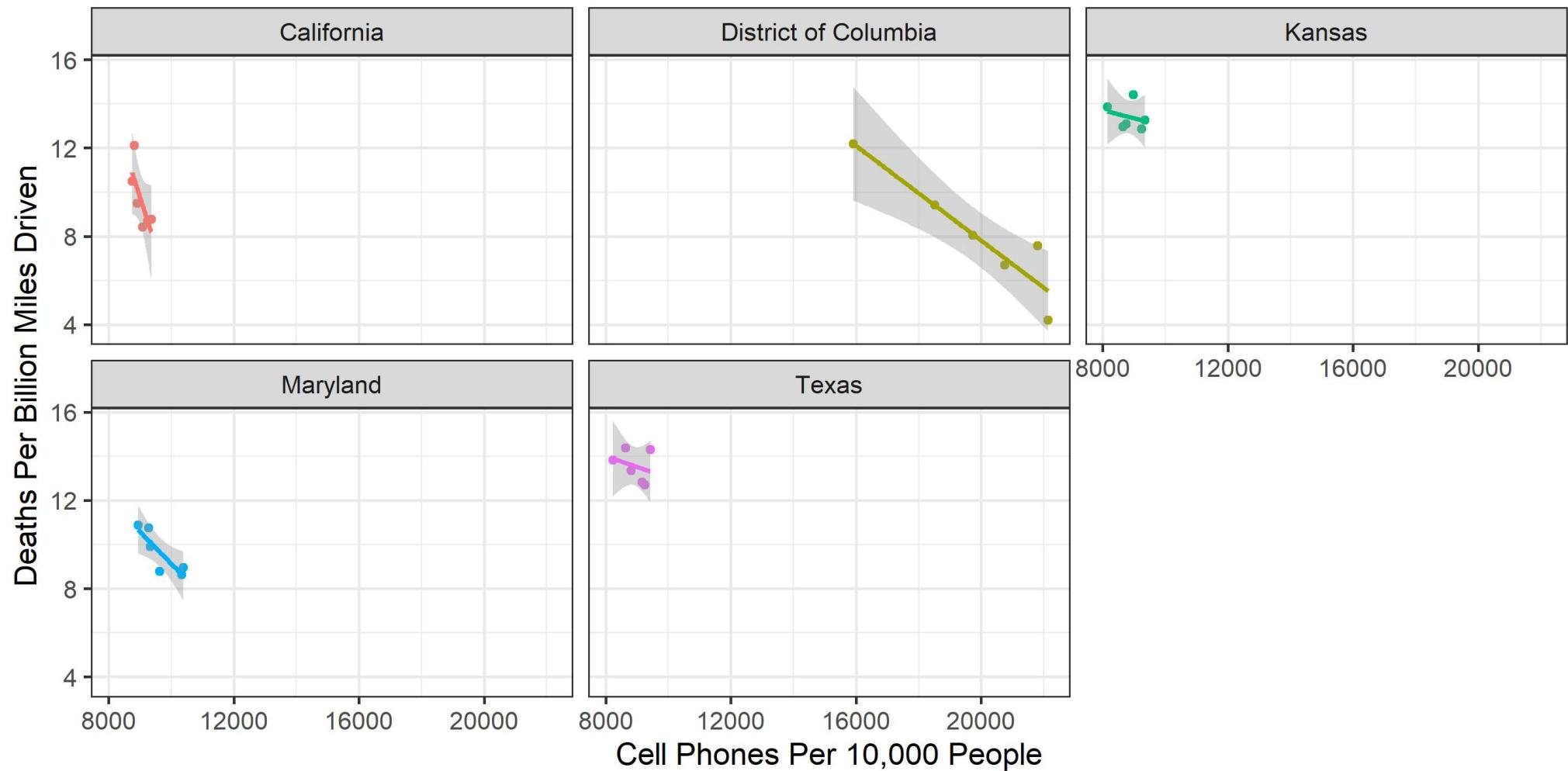
Example: Consider Just 5 States

```
1 phones %>%
2   filter(state %in% c("District of Columbia",
3                       "Maryland", "Texas",
4                       "California", "Kansas")) %>%
5   ggplot(data = .) +
6   aes(x = cell_plans,
7        y = deaths,
8        color = state) + #<<
9   geom_point() + #<<
10  geom_smooth(method = "lm") + #<<
11  labs(x = "Cell Phones Per 10,000 People",
12        y = "Deaths Per Billion Miles Driven",
13        color = NULL) +
14  theme_bw(base_family = "Fira Sans Condensed",
15            base_size=14) +
16  theme(legend.position = "top")
```



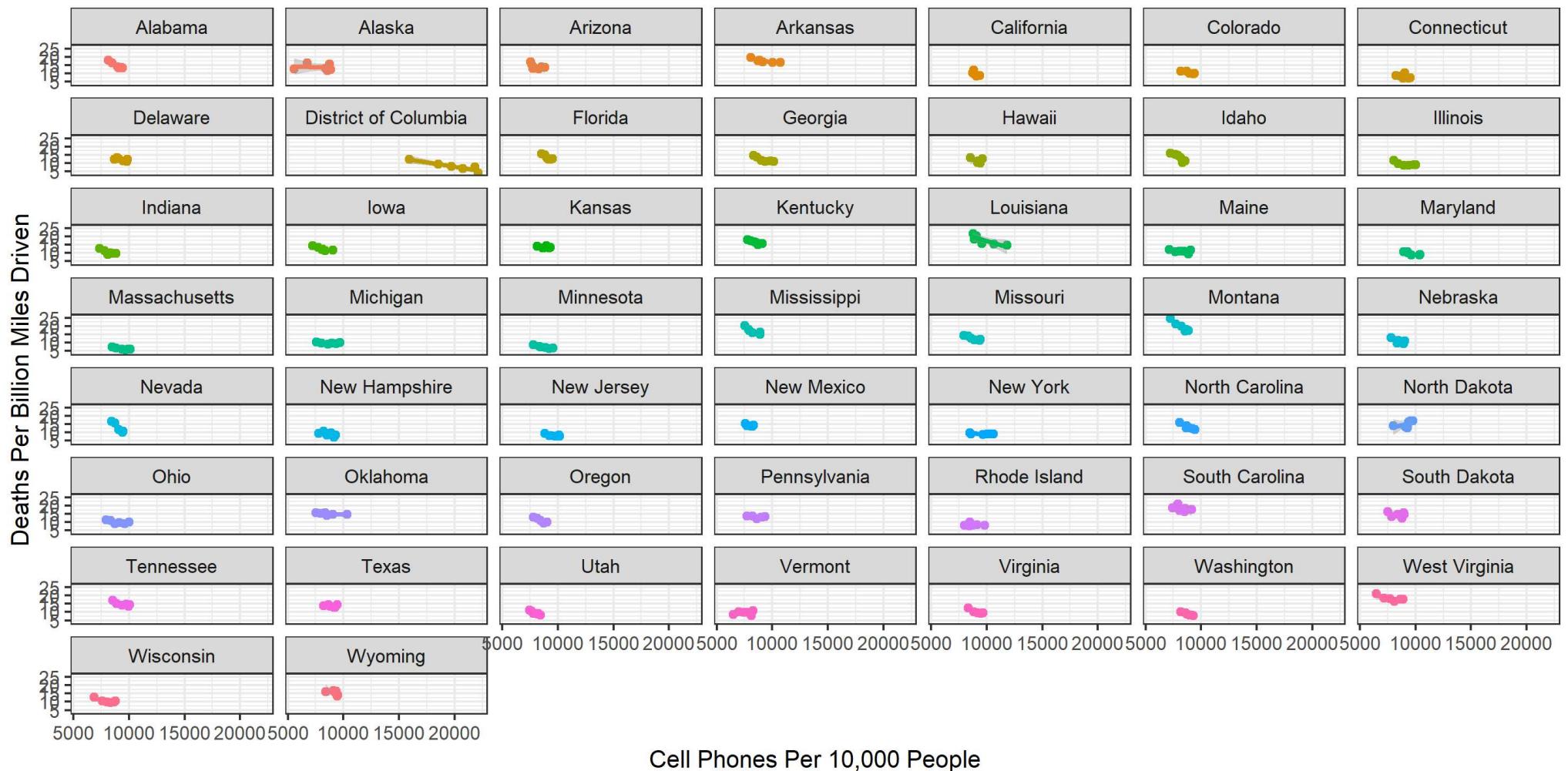
Example: Consider Just 5 States

```
1 phones %>%
2   filter(state %in% c("District of Columbia",
3                       "Maryland", "Texas",
4                       "California", "Kansas")) %>%
5   ggplot(data = .) +
6   aes(x = cell_plans,
7        y = deaths,
8        color = state) +
9   geom_point() +
10  geom_smooth(method = "lm") +
11  labs(x = "Cell Phones Per 10,000 People",
12        y = "Deaths Per Billion Miles Driven",
13        color = NULL) +
14  theme_bw(base_family = "Fira Sans Condensed",
15           base_size=14) +
16  theme(legend.position = "none") + #<<
17  facet_wrap(~state, ncol=3) #<<
```



Look at All States

```
1 ggplot(data = phones) + #<<
2   aes(x = cell_plans,
3       y = deaths,
4       color = state) +
5   geom_point() +
6   geom_smooth(method = "lm") +
7   labs(x = "Cell Phones Per 10,000 People",
8       y = "Deaths Per Billion Miles Driven",
9       color = NULL) +
10  theme_bw(base_family = "Fira Sans Condensed") +
11  theme(legend.position = "none") +
12  facet_wrap(~state, ncol=7) #<<
```



The Bias in our Pooled Regression

$$\widehat{\text{Deaths}}_{it} = \beta_0 + \beta_1 \text{Cell Phones}_{it} + u_{it}$$

- Cell Phones_{it} is **endogenous**:

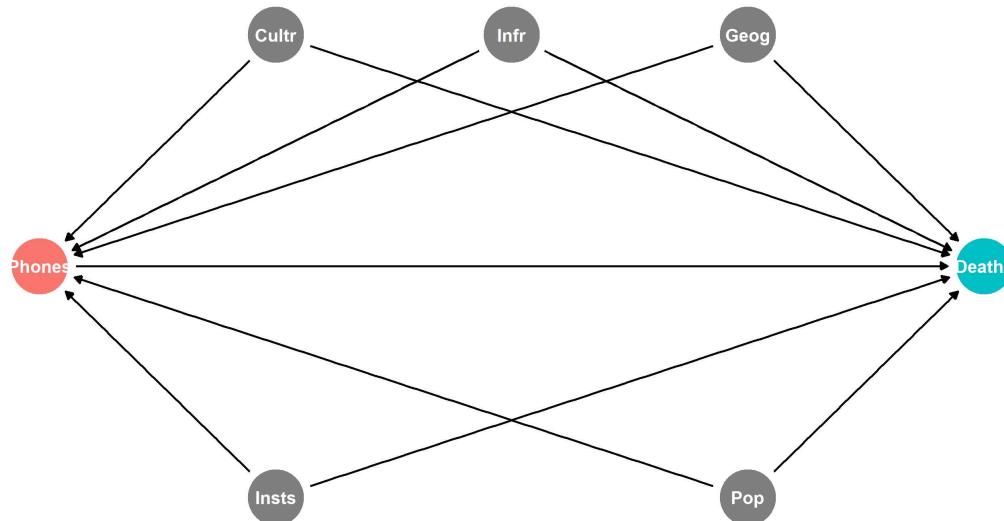
$$\text{cor}(u_{it}, \text{cell phones}_{it}) \neq 0 \quad E[u_{it} | \text{cell phones}_{it}] \neq 0$$

- Things in u_{it} correlated with Cell phones $_{it}$:
 - infrastructure spending, population, urban vs. rural, more/less cautious citizens, cultural attitudes towards driving, texting, etc
- A lot of these things vary systematically **by State!**
 - $\text{cor}(\mathbf{u}_{it_1}, \mathbf{u}_{it_2}) \neq 0$
 - Error in State i during t_1 correlates with error in State i during t_2
 - things in State that don't change over time

Fixed Effects Model

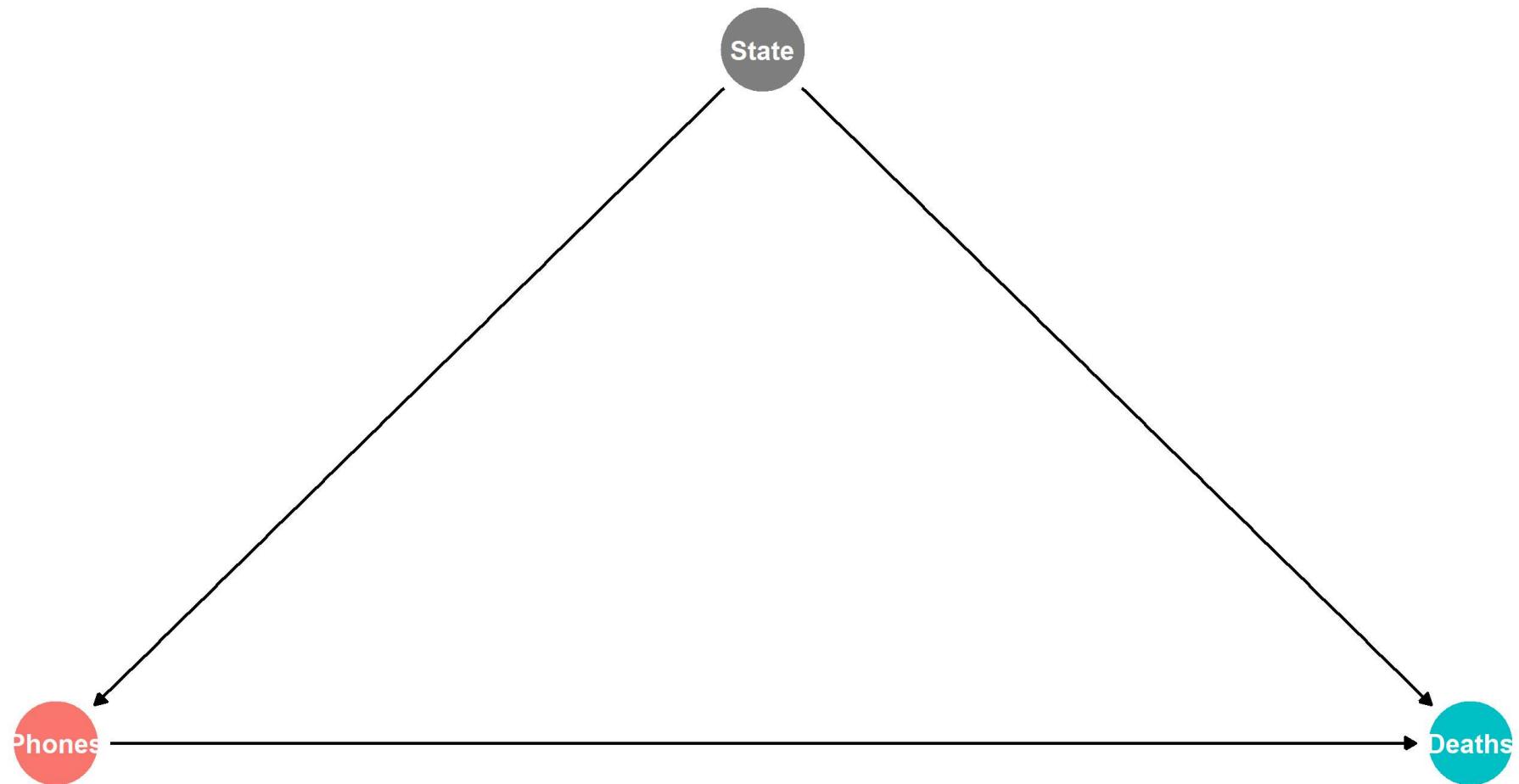
Fixed Effects: DAG

- A simple pooled model likely contains lots of omitted variable bias
- Many (often unobservable) factors that determine both Phones & Deaths
 - Culture, infrastructure, population, geography, institutions, etc



Fixed Effects: DAG

- A simple pooled model likely contains lots of omitted variable bias
- Many (often unobservable) factors that determine both Phones & Deaths
 - Culture, infrastructure, population, geography, institutions, etc
- But the beauty of this is that **most of these factors systematically vary by U.S. State and are stable over time!**
- We can simply **“control for State”** to safely remove the influence of all of these factors!



Fixed Effects: Decomposing u_{it}

- Much of the endogeneity in X_{it} can be explained by systematic differences across i (groups)
- Exploit the systematic variation across groups with a **fixed effects model**
- *Decompose* the model error term into two parts:

$$u_{it} = \alpha_i + \epsilon_{it}$$

Fixed Effects: α_i

- Decompose the model error term into two parts:

$$\mathbf{u}_{it} = \alpha_i + \epsilon_{it}$$

- α_i are **group-specific fixed effects**

- group i tends to have higher or lower \hat{Y} than other groups given regressor(s) X_{it}
- estimate a separate α_i for each group i
- essentially, estimate a separate constant (intercept) *for each group*
- notice this is stable over time within each group (subscript only i , no t)
- **This includes all factors that do not change *within* group i over time**

Fixed Effects: ϵ_{it}

$$u_{it} = \alpha_i + \epsilon_{it}$$

- ϵ_{it} is the remaining random error
 - As usual in OLS, assume the 4 typical assumptions about this error:
 - $E[\epsilon_{it}] = 0, var[\epsilon_{it}] = \sigma_\epsilon^2, cor(\epsilon_{it}, \epsilon_{jt}) = 0, cor(\epsilon_{it}, X_{it}) = 0$
- ϵ_{it} includes all other factors affecting Y_{it} *not* contained in group effect α_i
 - i.e. differences *within* each group that *change* over time
 - Be careful: X_{it} can still be endogenous from other factors!

Fixed Effects: New Regression Equation

$$\hat{Y}_{it} = \beta_0 + \beta_1 X_{it} + \alpha_i + \epsilon_{it}$$

- We've pulled α_i out of the original error term into the regression
- Essentially we'll estimate an intercept for each .pink[group] (minus one, which is β_0)
 - avoiding the dummy variable trap
- Must have multiple observations (over time) for each group (i.e. panel data)

Fixed Effects: Our Example

$$\widehat{\text{Deaths}}_{it} = \beta_0 + \beta_1 \text{Cell phones}_{it} + \alpha_i + \epsilon_{it}$$

- α_i is the .hi-pink[State fixed effect]
 - Captures everything unique about each state i that *does not change over time*
 - culture, institutions, history, geography, climate, etc!
- There could *still* be factors in ϵ_{it} that are correlated with Cell phones_{it} !
 - things that do change over time within States
 - perhaps individual States have cell phone bans for *some* years in our data

Estimating Fixed Effects Models

$$\hat{Y}_{it} = \beta_0 + \beta_1 X_{it} + \alpha_i + \epsilon_{it}$$

- Two methods to estimate fixed effects models:
 1. Least Squares Dummy Variable (LSDV) approach
 2. De-meaned data approach

Least Squares Dummy Variable Approach

$$\widehat{Y}_{it} = \beta_0 + \beta_1 X_{it} + \beta_2 D_{1i} + \beta_3 D_{2i} + \cdots + \beta_N D_{(N-1)i} + \epsilon_{it}$$

##

- A dummy variable $D_i = \{0, 1\}$ for each possible group
 - = 1 if observation it is from group i , otherwise = 0
- If there are N groups:
 - Include $N - 1$ dummies (to avoid **dummy variable trap**) and β_0 is the reference category^{magenta[†]}
 - So we are estimating a different intercept for each group
- Sounds like a lot of work, automatic in R

If we do not estimate β_0 , we could include all N dummies. In either case, β_0 takes the place of one category-dummy.

Least Squares Dummy Variable Approach: Our Example

$$\widehat{\text{Deaths}}_{it} = \beta_0 + \beta_1 \text{Cell Phones}_{it} + \text{Alaska}_i + \cdots + \text{Wyoming}_i$$

- Let Alabama be the reference category (β_0), include all other States

Our Example in R I

$$\widehat{\text{Deaths}}_{it} = \beta_0 + \beta_1 \text{Cell Phones}_{it} + \text{Alaska}_i + \cdots + \text{Wyoming}_i$$

- If `state` is a `factor` variable, just include it in the regression
- R automatically creates $N - 1$ dummy variables and includes them in the regression
 - Keeps intercept and leaves out first group dummy

Our Example in R II

```

1 fe_reg_1 <- lm(deaths ~ cell_plans + state, data = phones)
2 fe_reg_1 %>% tidy()
3 ## # A tibble: 52 × 5
4 ##   term          estimate std.error statistic p.value
5 ##   <chr>        <dbl>     <dbl>      <dbl>    <dbl>
6 ## 1 (Intercept)  25.5      1.02      25.1  1.24e-70
7 ## 2 cell_plans   -0.00120  0.000101   -11.9  3.48e-26
8 ## 3 stateAlaska   -2.48     0.675     -3.68  2.82e- 4
9 ## 4 stateArizona  -1.51     0.670     -2.25  2.51e- 2
10 ## 5 stateArkansas  3.19     0.666      4.79  2.83e- 6
11 ## 6 stateCalifornia -4.98    0.666     -7.48  1.21e-12
12 ## 7 stateColorado  -4.34     0.665     -6.53  3.59e-10
13 ## 8 stateConnecticut -6.60     0.665     -9.91  8.70e-20
14 ## 9 stateDelaware   -2.10     0.667     -3.15  1.84e- 3
15 ## 10 stateDistrict of Columbia  6.36     1.29      4.93  1.50e- 6
16 ## # ... with 42 more rows
17 ## # i Use `print(n = ...)` to see more rows

```

De-meaned Approach

Zahid Asghar



De-meaned Approach I

- Alternatively, we can control our regression for group fixed effects without directly estimating them
- We simply **de-mean the data for each group**
- For each group i , find the means (over time, t):

$$\bar{Y}_i = \beta_0 + \beta_1 \bar{X}_i + \bar{\alpha}_i + \bar{\epsilon}_{it}$$

- Where:
 - \bar{Y}_i : average value of Y_{it} for group i
 - \bar{X}_i : average value of X_{it} for group i
 - $\bar{\alpha}_i$: average value of α_i for group i ($= \alpha_i$)
 - $\bar{\epsilon}_{it} = 0$, by assumption 1

De-meaned Approach II

\$\$

$$\begin{aligned}\widehat{Y}_{it} &= \beta_0 + \beta_1 X_{it} + u_{it} \\ \bar{Y}_i &= \beta_0 + \beta_1 \bar{X}_i + \bar{\alpha}_i + \bar{\epsilon}_i\end{aligned}$$

\$\$

- Subtract the means equation from the pooled equation to get:

$$\begin{aligned}Y_i - \bar{Y}_i &= \beta_1(X_{it} - \bar{X}_i) + \tilde{\epsilon}_{it} \\ \tilde{Y}_{it} &= \beta_1 \tilde{X}_{it} + \tilde{\epsilon}_{it}\end{aligned}$$

- Within each group i , the de-meaned variables \tilde{Y}_{it} and \tilde{X}_{it} 's all have a mean of 0.^{magenta[†]}
- Variables that don't change over time will drop out of analysis altogether
- Removes any source of variation **across** groups to only work with variation **within** each group

Recall : Summation Operator: $\sum(X_i - \bar{X}) = 0$

De-meaned Approach III

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

- Yields identical results to dummy variable approach
- More useful when we have many groups (would be many dummies)
- Demonstrates **intuition** behind fixed effects:
 - Converts all data to deviations from the mean of each group
 - All groups are “centered” at 0
 - Fixed effects are often called the .hi-pink[“within” estimators], they exploit variation *within* groups, not *across* groups

De-meaned Approach IV

- We are basically comparing groups *to themselves* over time
 - apples to apples comparison
 - e.g. Maryland in 2000 vs. Maryland in 2005
- Ignore all differences *between* groups, only look at differences *within* groups over time

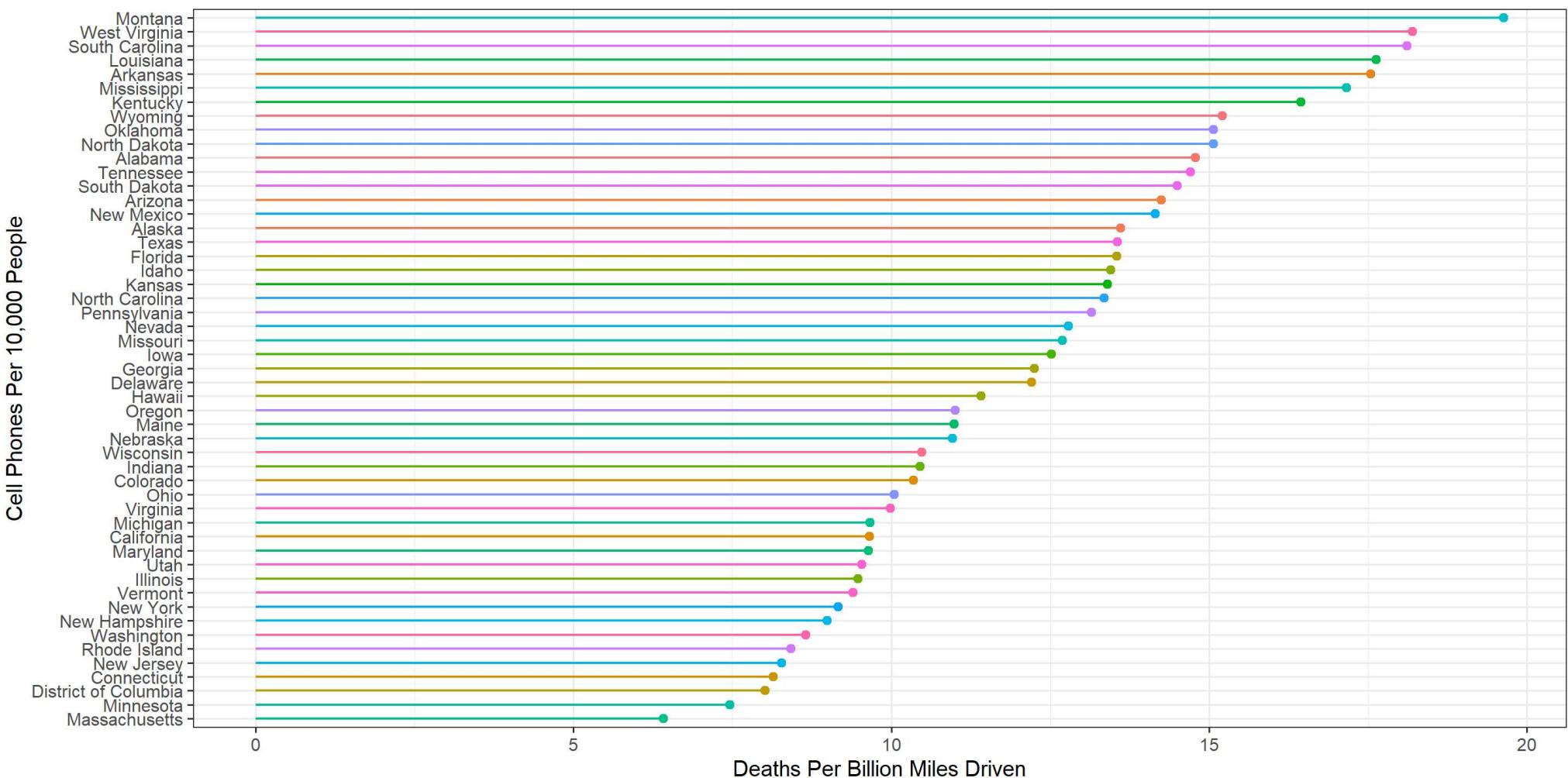
De-Meaning the Data in R I

```
1 # get means of Y and X by state
2 means_state<-phones %>%
3   group_by(state) %>%
4   summarize(avg_deaths = mean(deaths),
5             avg_phones = mean(cell_plans))
6
7 # look at it
8 means_state
```

```
## # A tibble: 51 × 3
##   state      avg_deaths avg_phones
##   <fct>        <dbl>       <dbl>
## 1 Alabama     14.8        8906.
## 2 Alaska      13.6        7818.
## 3 Arizona     14.2        8097.
## 4 Arkansas    17.5        9268.
## 5 California   9.66       9030.
## 6 Colorado     10.4       8982.
## 7 Connecticut  8.14       8948.
## 8 Delaware     12.2       9304.
## 9 District of Columbia  8.02      19811.
## 10 Florida     13.5       9079.
## # ... with 41 more rows
## # i Use `print(n = ...)` to see more rows
```

De-Meaning the Data in R II

```
1 ggplot(data = means_state)+  
2   aes(x = fct_reorder(state, avg_deaths),  
3       y = avg_deaths,  
4       color = state)+  
5   geom_point() +  
6   geom_segment(aes(y = 0,  
7                     yend = avg_deaths,  
8                     x = state,  
9                     xend = state)) +  
10  coord_flip() +  
11  labs(x = "Cell Phones Per 10,000 People",  
12        y = "Deaths Per Billion Miles Driven",  
13        color = NULL) +  
14  theme_bw(base_family = "Fira Sans Condensed",  
15            base_size=10) +  
16  theme(legend.position = "none")
```



Visualizing "Within Estimates" for the 5 States

Visualizing "Within Estimates" for All 51 States

De-meaned Approach in R I

- The `plm` package is designed for panel data
- `plm()` function is just like `lm()`, with some additional arguments:
 - `index="group_variable_name"` set equal to the name of your `factor` variable for the groups
 - `model=` set equal to `"within"` to use fixed-effects (within-estimator)

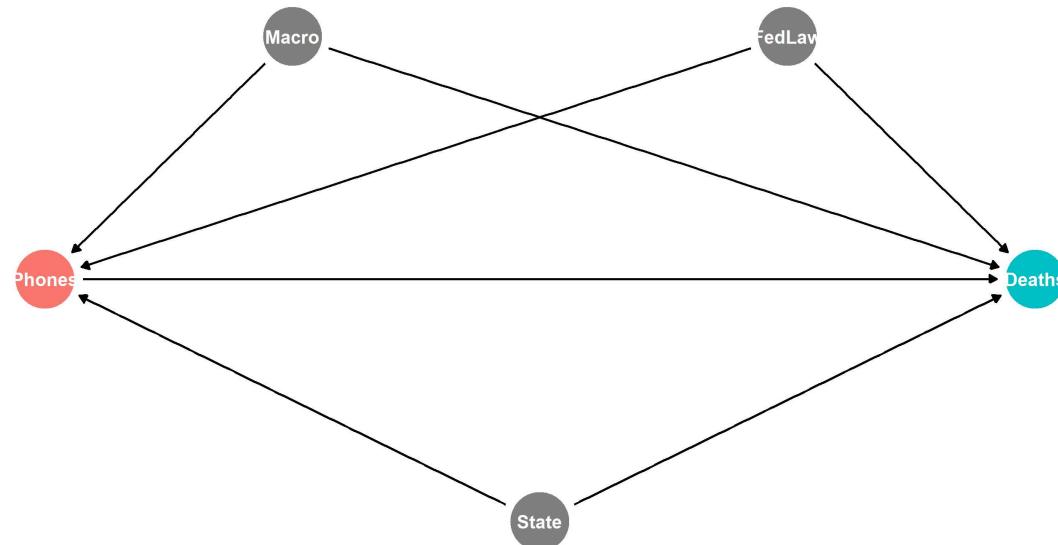
```
1 #install.packages("plm")
2 library(plm)
3 fe_reg_1_alt<-plm(deaths ~ cell_plans,
4                     data = phones,
5                     index = "state",
6                     model = "within")
```

De-meaned Approach in R II

```
1 fe_reg_1_alt %>% tidy()  
2 ## # A tibble: 1 × 5  
3 ##   term      estimate std.error statistic p.value  
4 ##   <chr>      <dbl>     <dbl>     <dbl>    <dbl>  
5 ## 1 cell_plans -0.00120  0.000101    -11.9  3.48e-26
```

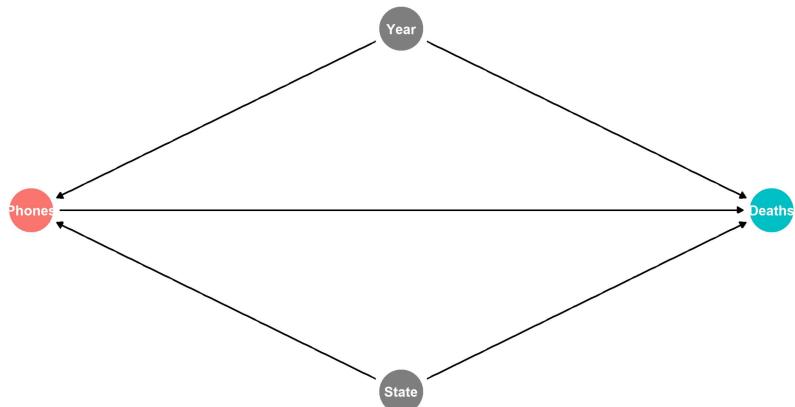
Two-Way Fixed Effects

- State fixed effect controls for all factors that vary by state but are stable over time
- But there are still other (often unobservable) factors that affect both Phones and Deaths, that *don't* vary by State
 - The country's macroeconomic performance, federal laws, etc



Two-Way Fixed Effects

- State fixed effect controls for all factors that vary by state but are stable over time
- But there are still other (often unobservable) factors that affect both Phones and Deaths, that *don't* vary by State
 - The country's macroeconomic performance, federal laws, etc
- If these factors systematically vary over time, but are the same by State, then we can **"control for Year"** to safely remove the influence of all of these factors!



Two-Way Fixed Effects

- A .hi[one-way fixed effects model] estimates a fixed effect for **groups**
- .hi[Two-way fixed effects model] estimates fixed effects for *both groups and time periods*

$$\hat{Y}_{it} = \beta_0 + \beta_1 X_{it} + \alpha_i + \theta_t + \nu_{it}$$

- α_i : group fixed effects
 - accounts for **time-invariant differences across groups**
- θ_t : time fixed effects
 - accounts for **group-invariant differences over time**
- ν_{it} remaining random error
 - all remaining factors that affect Y_{it} that vary by state *and* change over time

Two-Way Fixed Effects: Our Example

$$\widehat{\text{Deaths}}_{it} = \beta_0 + \beta_1 \text{Cell phones}_{it} + \alpha_i + \theta_t + \nu_{it}$$

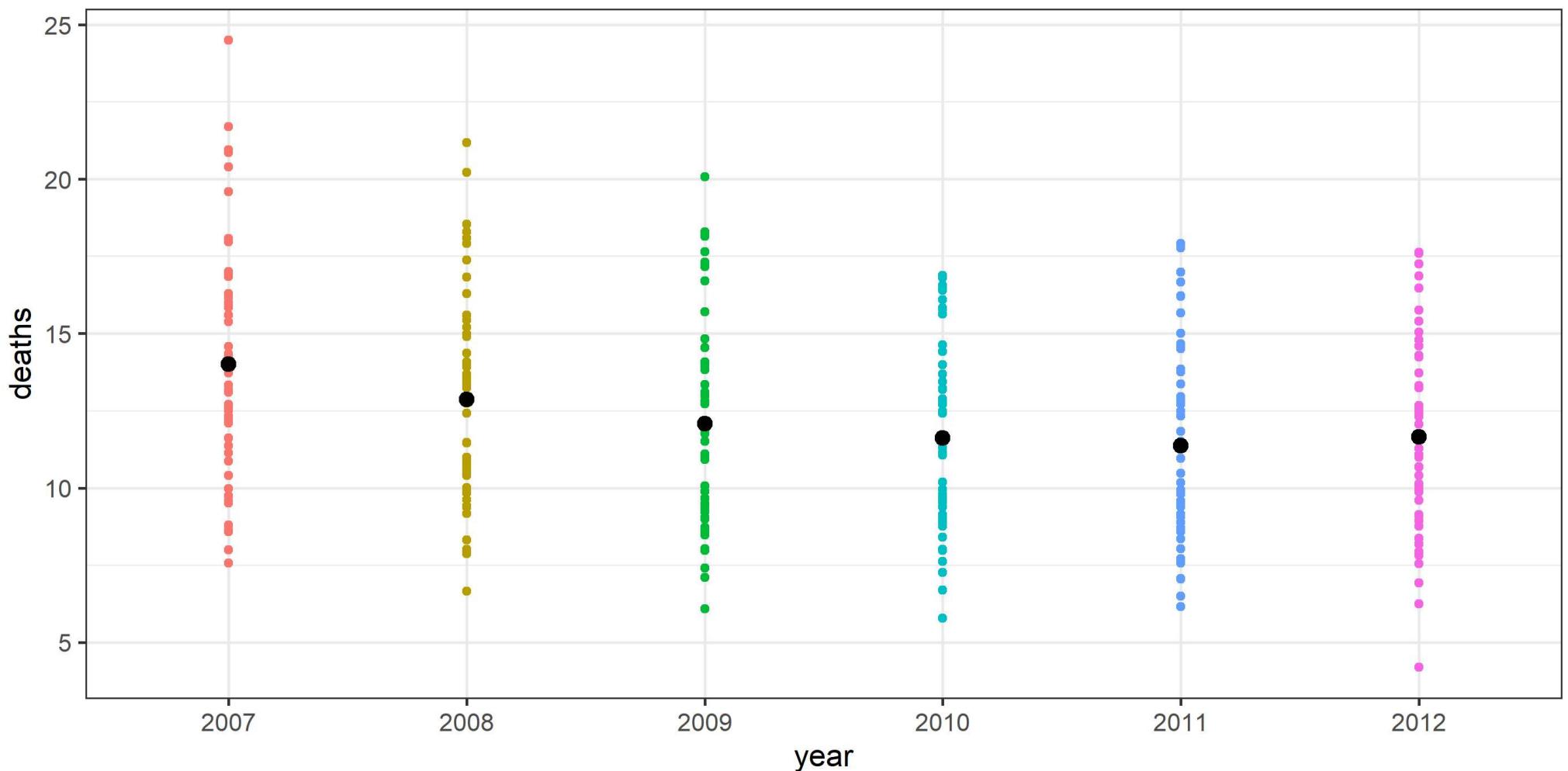
- α_i : **State fixed effects**
 - differences **across states** that are **stable over time** (note subscript i only)
 - e.g. geography, culture, (unchanging) state laws
- θ_t : Year fixed effects
 - differences **over time** that are **stable across states** (note subscript t only)
 - e.g. economy-wide macroeconomic changes, *federal* laws passed

Visualizing Year Effects I

```
1 # find averages for years
2 means_year<-phones %>%
3   group_by(year) %>%
4   summarize(avg_deaths = mean(deaths),
5             avg_phones = mean(cell_plans))
6 means_year
7 ## # A tibble: 6 × 3
8 ##   year  avg_deaths avg_phones
9 ##   <dbl>     <dbl>      <dbl>
10## 1 2007     14.0     8065.
11## 2 2008     12.9     8483.
12## 3 2009     12.1     8860.
13## 4 2010     11.6     9135.
14## 5 2011     11.4     9485.
15## 6 2012     11.7     9660.
```

Visualizing Year Effects II

```
1 ggplot(data = phones) +  
2   aes(x = year,  
3       y = deaths) +  
4   geom_point(aes(color = year)) +  
5  
6   # Add the yearly means as black points  
7   geom_point(data = means_year,  
8             aes(x = year,  
9                 y = avg_deaths),  
10            size = 3,  
11            color = "black") +  
12  
13  geom_path(data = means_year,  
14             aes(x = year,  
15                 y = avg_deaths),  
16             size = 1) +  
17  theme_bw(base_family = "Fira Sans Condensed",  
18            base_size = 14) +  
19  theme(panel.grid.major = element_line()  
20        , panel.grid.minor = element_line())
```



Estimating Two-Way Fixed Effects

$$\hat{Y}_{it} = \beta_0 + \beta_1 X_{it} + \alpha_i + \theta_t + \nu_{it}$$

- As before, several equivalent ways to estimate two-way fixed effects models:
1. **Least Squares Dummy Variable (LSDV) Approach:** add dummies for both groups and time periods (separate intercepts for groups and times)
 2. **Fully De-meaned data:**

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

where for each variable: $\tilde{var}_{it} = var_{it} - \overline{var}_t - \overline{var}_i$

3. **Hybrid:** de-mean for one effect (groups or years) and add dummies for the other effect (years or groups)

LSDV Method

```

1 fe2_reg_1 <- lm(deaths ~ cell_plans + state + year,
2                     data = phones)
3 fe2_reg_1 %>% tidy()
4 ## # A tibble: 57 × 5
5 ##   term          estimate std.error statistic p.value
6 ##   <chr>        <dbl>     <dbl>      <dbl>    <dbl>
7 ## 1 (Intercept) 18.9       1.45      13.0  5.43e-30
8 ## 2 cell_plans -0.000300  0.000172   -1.74 8.34e- 2
9 ## 3 stateAlaska -1.50      0.624     -2.40  1.70e- 2
10 ## 4 stateArizona -0.779     0.611     -1.27  2.04e- 1
11 ## 5 stateArkansas 2.87      0.599      4.79  2.90e- 6
12 ## 6 stateCalifornia -5.09     0.596     -8.55 1.30e-15
13 ## 7 stateColorado -4.41      0.595     -7.41  1.95e-12
14 ## 8 stateConnecticut -6.63     0.595    -11.1  1.17e-23
15 ## 9 stateDelaware -2.46      0.599     -4.10  5.55e- 5
16 ## 10 stateDistrict of Columbia -3.50     1.97    -1.78  7.66e- 2
17 ## # ... with 47 more rows
18 ## # i Use `print(n = ...)` to see more rows

```

With plm

```

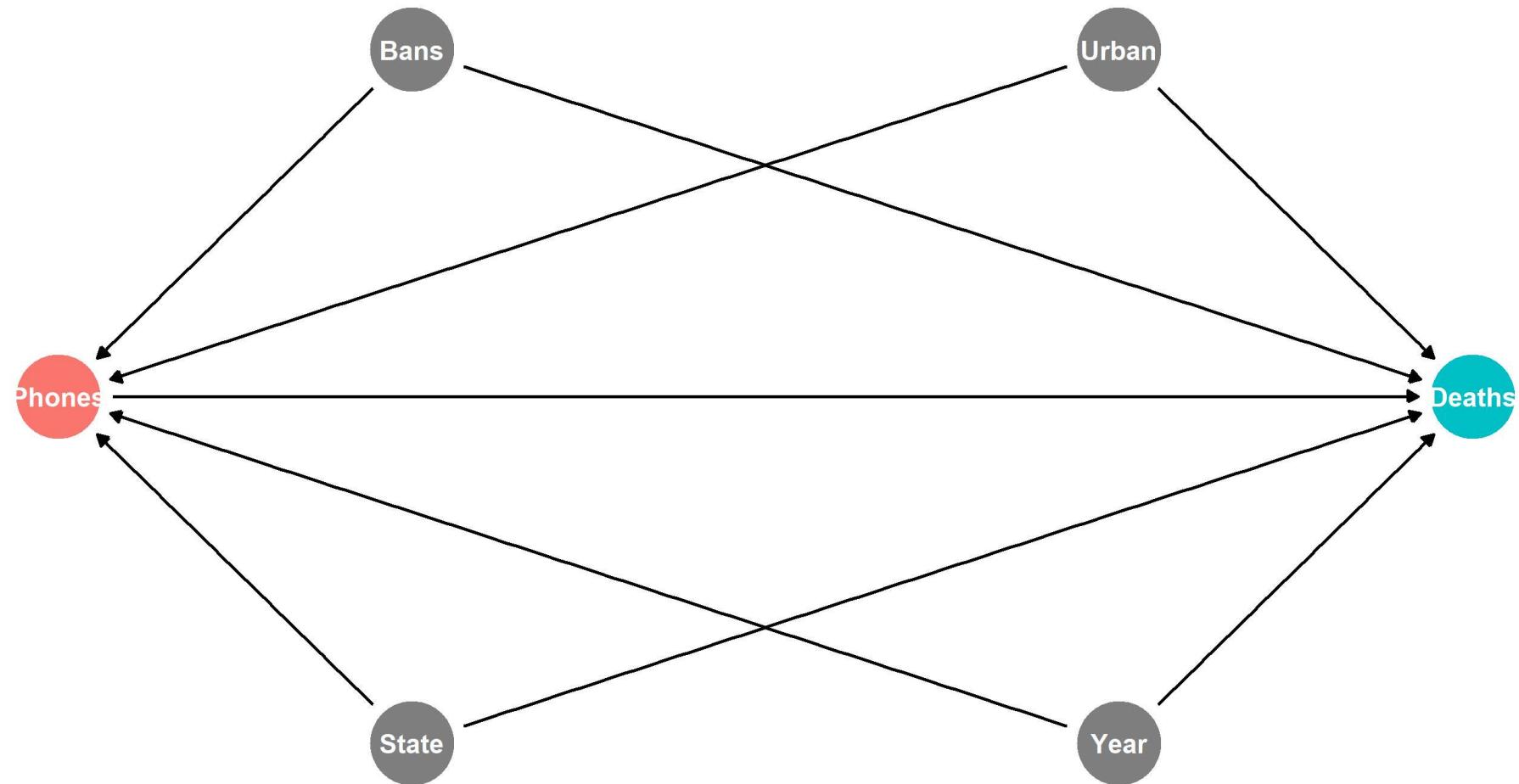
1 fe2_reg_2 <- plm(deaths ~ cell_plans,
2                         index = c("state", "year"),
3                         model = "within",
4                         data = phones)
5 fe2_reg_2 %>% tidy()
6 ## # A tibble: 1 × 5
7 ##   term      estimate std.error statistic p.value
8 ##   <chr>      <dbl>     <dbl>     <dbl>    <dbl>
9 ## 1 cell_plans -0.00120  0.000101    -11.9  3.48e-26

```

- `plm()` command allows for multiple effects to be fit inside `index=c("group", "time")`

Adding Covariates

- State fixed effect absorbs all unobserved factors that vary by state, but are constant over time
- Year fixed effect absorbs all unobserved factors that vary by year, but are constant over States
- But there are still other (often unobservable) factors that affect both Phones and Deaths, that *vary* by State *and* change over time!
 - Some States *change* their laws during the time period
 - State *urbanization* rates *change* over the time period
- We will also need to .hi-pink[control for these variables] (*not* picked up by fixed effects!)
 - Add them to the regression



Adding Covariates I

$$\widehat{\text{Deaths}}_{it} = \beta_1 \text{Cell Phones}_{it} + \alpha_i + \theta_t + \text{urban pct}_{it} + \text{cell ban}_{it} + \text{text ban}_{it}$$

- Can still add covariates to remove endogeneity not soaked up by fixed effects
 - factors that change within groups over time
 - e.g. some states pass bans over the time period in data (some years before, some years after)

Adding Covariates II

```

1 fe2_controls_reg <- plm(deaths ~ cell_plans + text_ban + urban_percent + cell_ban,
2                           data = phones,
3                           index = c("state", "year"),
4                           model = "within",
5                           effect = "twoways")
6
7 fe2_controls_reg %>% tidy()
8 ## # A tibble: 4 × 5
9 ##   term      estimate std.error statistic p.value
10 ##   <chr>     <dbl>    <dbl>     <dbl>    <dbl>
11 ## 1 cell_plans -0.000340  0.000173    -1.97   0.0502
12 ## 2 text_ban1    0.256     0.222     1.15    0.251
13 ## 3 urban_percent  0.0131    0.0112     1.17    0.242
14 ## 4 cell_ban1    -0.680     0.403     -1.69   0.0929

```

Comparing Models

```
1 library(huxtable)
2 huxreg("Pooled" = pooled,
3        "State Effects" = fe_reg_1,
4        "State & Year Effects" = fe2_reg_1,
5        "With Controls" = fe2_controls_reg,
6        coefs = c("Intercept" = "(Intercept)",
7                  "Cell phones" = "cell_plans",
8                  "Cell Ban" = "cell_ban1",
9                  "Texting Ban" = "text_ban1",
10                 "Urbanization Rate" = "urban_percent"),
11        statistics = c("N" = "nobs",
12                      "R-Squared" = "r.squared",
13                      "SER" = "sigma"),
14        number_format = 4)
```

	Pooled	State Effects	State & Year Effects	With Controls
Intercept	17.3371 *** (0.9754)	25.5077 *** (1.0176)		18.9305 *** (1.4511)
Cell phones	-0.0006 *** (0.0001)	-0.0012 *** (0.0001)		-0.0003 (0.0002)
Cell Ban				-0.6798 (0.4029)
Texting Ban				0.2559 (0.2222)
Urbanization Rate				0.0131 (0.0112)
N	306	306	306	306

R-Squared	0.0845	0.9055	0.9259	0.0329
SER	3.2791	1.1526	1.0310	

*** p < 0.001; ** p < 0.01; * p < 0.05.