# Day 6: Panel Data Models

# Panel Data

Panels are a hybrid data structure that lives between the traditional data structures of microeconomics and forecasting.

- Contains observations of multiple individuals
  - Similar to standard cross-sectional data
- Contains multiple observations of each individual
  - Makes the data a collection of [possibly multivariate] time series data

# Panel Data

Forecasting algorithms like ARIMA models and GAMs cannot cope with this kind of data structure

- How do we difference out a time series when we have multiple observations (of different individuals) in any given period?

- How do we control for unobservable or unmeasurable differences between individuals?

# Panel Data

Panel data allows us to generalize much of what we can learn through time series analysis

- We can generalize the effect of covariates to more than one individual

- We can make forecasts for different groups simultaneously from the same model

# Working with Panel Data

$$y_{it} = \alpha_{it} + X_{it}\beta + \epsilon_{it}$$

$i$: individual index, $t$: time index

We might start with the model above, but we wouldn't get far.

- We have insufficient information to calculate the model!
    - $K + NT > NT$

# Working with Panel Data

$$y_{it} = \alpha + X_{it}\beta + \epsilon_{it}$$

If we remove the individual-level intercepts, we can remedy our information problem.

- Now, so long as we choose a reasonable number of covariates, $K < N$

# Working with Panel Data

$$y_{it} = \alpha + X_{it}\beta + \epsilon_{it}$$

Unfortunately, panel data means that we have correlated error terms within individuals.

- There is no good reason to believe

$$corr(y_{it}, y_{it+1}) = 0$$

# Working with Panel Data

$$y_{it} = \alpha + X_{it}\beta + \epsilon_{it}$$

We need to decompose our error terms so that

$$\epsilon_{it} = \mu_i + \nu_{it}$$

where $\mu_i$ is an individual **fixed effect**, and $\nu_{it}$ is the noise term.

# Working with Panel Data

$$y_{it} = \alpha + X_{it}\beta + \mu_i + \nu_{it}$$

Our model now has $K + N$ parameters, and $NT$ degrees of freedom.

- We can now solve our model!

# Working with Panel Data

$$y_{it} = \alpha + X_{it}\beta + \mu_i + \nu_{it}$$

The model can actually be solved using a modified form of OLS.

# Working with Panel Data

$$y_{it} = \alpha + X_{it}\beta + \mu_i + \nu_{it}$$

$$\downarrow$$

$$y_{it} - \bar{y}_i = (X_{it} - \bar{X}_i)\beta + \nu_{it} - \bar{\nu}_i$$

$$\downarrow$$

$$\ddot{y}_{it} = \ddot{X}_{it}\beta + \ddot{\nu}_{it}$$

# Working with Panel Data

$$\ddot{y}_{it} = \ddot{X}_{it}\beta + \ddot{\nu}_{it}$$

We difference each observation by subtracting the average values for a given individual over time, causing the intercept terms and individual fixed effects to be differenced out of the model.

$$\bar{X}_i = \frac{1}{T}\sum_{t=1}^{T} X_{it}$$

# Implementing A Fixed Effects Model

```python
# Import Libraries
import pandas as import pd
import numpy as np
import statsmodels.formula.api as sm

# Import Data
data = pd.read_csv(
        '/home/dusty/DatasetsDA/firmInvestmentPanel.csv')
```

If we import the formula module from `statsmodels`, then we are able to implement R-style formulas in our regressions.

# EXERCISE TIME!

With your classmates, find a way to difference out the firm-level means from the panel data in `firmInvestmentPanel.csv`.

Hint: These kinds of procedures always perform better if you can find pre-built functions in either Pandas or Numpy that do what you need.

# EXERCISE ANSWER

```python
# We now need to de-mean our data
vars = ['I_','F_','C_']
for i in data.FIRM.unique():
    data.loc[data.FIRM==1, vars] =\ # '\' indicates a line
        data.loc[data.FIRM==1, vars] -\ #            break
        data.loc[data.FIRM==1, vars].mean()
```

We only want to difference out means for numeric data on the firm-level, not on the year or firm columns. We need to select our columns for differencing accordingly.

# Implementing A Fixed Effects Model

```python
# Specify regression
reg = sm.ols("I_ ~ F_ + C_ + C(FIRM) + YEAR + I(YEAR**2)",
        data=data[data.YEAR<1954]) # Last year saved as
                                   # forecast
# Fit regression with clustered robust standard errors
# Note: We need to have same year restriction on groups
fit = reg.fit().get_robustcov_results(cov_type='cluster',
        groups=data.loc[data.YEAR<1954, 'FIRM'])
# Print results
print(fit.summary())
```

We can now explore our results, the effects of included variables, and what our forecasts might look like.

# Implementing A Fixed Effects Model

```python
# Store predictions and truth
pred = fit.predict(data[data.YEAR==1954])
truth = data.loc[data.YEAR==1954, "I_"]
# Store errors
errors = pred - truth
# Calculate Absolute Percentage Error
pce = np.abs(errors/truth)*100
```

We need to perform the calculations that will provide us with information on how well we do at predicting out of sample with our current panel.

# Implementing A Fixed Effects Model

```python
# Print MSE, Mean Absolute Error,
#    and Mean Abs Percentage Error
print("Mean Squared Error: %s" %
        str(np.mean(errors**2)))
print("Mean Absolute Error: %s" %
        str(np.mean(np.abs(errors))))
print("Mean Absolute Percentage Error: %s"
        % str(np.mean(pce)))
```

Mean Squared Error: 13288.423957448418

Mean Absolute Error: 77.27884184438867

Mean Absolute Percentage Error: 58.253213431705774

Yikes! It looks like we need more information...

# For Lab Today:

We will look at how well we can forecast [student's grades](#) based on information about their study habits, social patterns, and family situation.

In your teams, develop a model for the data contained in `continuousTrain.csv` that will allow you to forecast a student's final grade (G3).

Then, use the model that you have built to forecast the grades for the student data contained in `continuousTest.csv`.