

ECON 255 — STATA Labs

Jihye Heo — Oberlin College

Contents

0	Introduction	3
0.1	Interface	3
0.2	Working Directory	3
0.3	File Types	4
0.4	Default Log Template	5
Week 1	— Introduction to STATA (Teaching)	6
1.1	Goals for Week 1	6
1.2	Setting for Week 1	6
1.3	Tasks for Week 1	6
Week 2	— Running OLS, Fitted Values, and Exporting Results (Teaching)	10
2.1	New Goals for Week 2	10
2.2	Setting for Week 2	10
2.3	Tasks for Week 2	11
Week 3	— Assignment 1: Housing Prices (Homework)	14
3.1	Setting	14
3.2	What to Submit	14
3.3	Assignment Tasks	14
Week 4	— Merging, Reshaping, and Inference (Teaching)	16
4.1	New Goals for Week 4	16
4.2	Setting for Week 4	16
4.3	Tasks for Week 4	17
Week 5	— Assignment 2 : Merging Movie Data (Homework)	23
5.1	Setting	23
5.2	What to Submit	23
5.3	Assignment Tasks	23
Week 6	— Nonlinearities in Regression (Teaching)	27
6.1	New Goals for Week 6	27
6.2	Setting for Week 6	27
6.3	Tasks for Week 6	27
Week 7	— Assignment 3: Nonlinearities in Regression (Homework)	31
7.1	Setting	31
7.2	What to Submit	31
7.3	Assignment Tasks	31
Week 8	— Introduction to Panel Data and Lagged Variables (Teaching)	34
8.1	New Goals for Week 8	34
8.2	Setting for Week 8	34
8.3	Tasks for Week 8	35
Week 9	— Assignment 4: Difference-in-Differences (DID) (Homework)	40

9.1	Setting	40
9.2	What to Submit	41
9.3	Assignment Tasks	41
Week 10	— Fixed Effects (Within Estimator) (Teaching)	46
10.1	New Goals for Week 10	46
10.2	Setting for Week 10	46
10.3	Tasks for Week 10	47
Week 11	— Assignment 5: Endogeneity & Instrumental Variables (IV) (Homework) 53	
11.1	Setting	53
11.2	What to Submit	54
11.3	Assignment Tasks	54
12	Command Directory	57
12.1	Basic	57
12.2	Data	58
12.3	Regression	61

0 Introduction

0.1 Interface

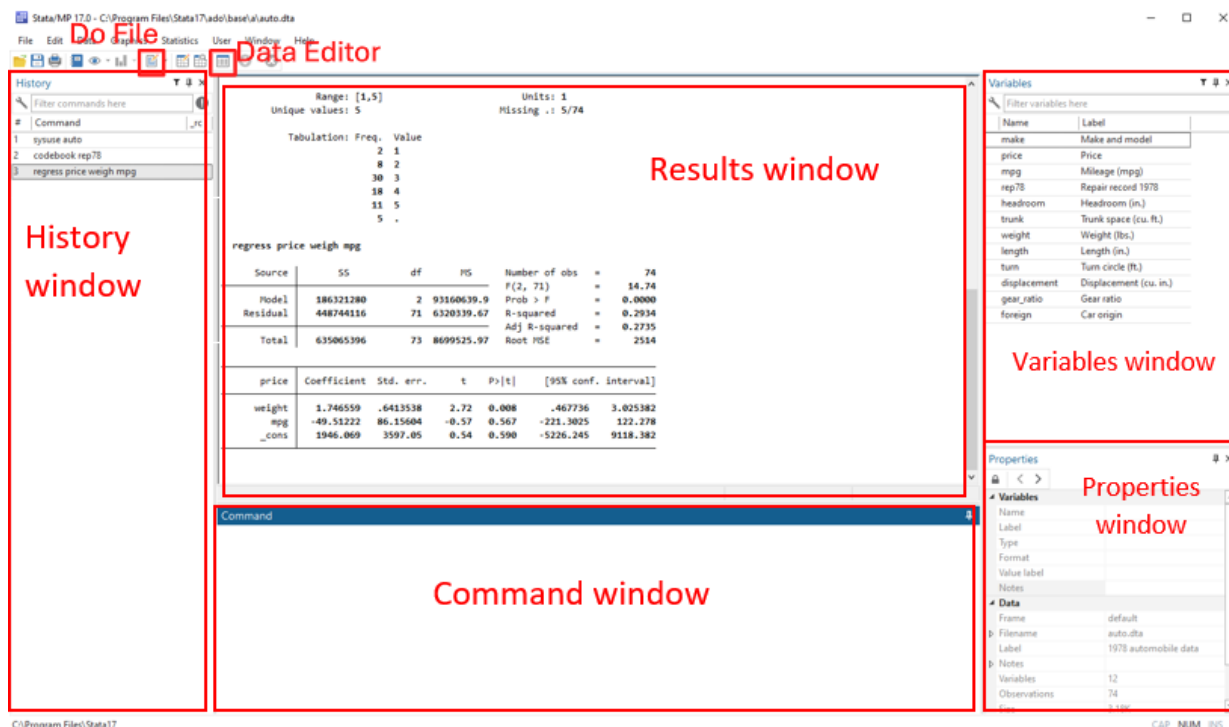


Figure 1: The STATA Interface

Key parts of the interface:

- **Command Window:** where you type Stata commands (not recommended for reproducibility).
- **Results Window:** displays command output.
- **Variables Window:** lists dataset variables currently in memory.
- **History Window:** shows your recent command history.

0.2 Working Directory

- **Default behavior:** Stata saves to and searches for files only in the current working directory, unless you specify a full file path.

Note

This means that if you do not set or change the working directory, Stata will not know where your dataset is located and will be unable to load it.

Tip: Working Directory for Our Lab Sessions

For lab sessions, set your working directory to the `Downloads` folder (where most datasets will be saved). To find the correct folder path:

1. Open your `Downloads` folder (or wherever you saved the dataset).
2. Click once on the dataset file to highlight it.
3. Right-click the file and select **Properties**.
4. In the Properties window, locate the field labeled **Location**.
5. Copy and paste this Location into Stata's `cd` command.

Set your working directory

```
1  * Display current directory
2  pwd
3
4  * Change working directory to Downloads (common in labs)
5  cd "C:\Users\YourUniquePath\Downloads"
6
7  * List files in the current directory
8  ls
```

0.3 File Types

- **Do-files (.do):** Contain the commands you write and execute. → We will always use Do-files instead of typing directly into the command window, because they let you:
 - Run multiple commands at once
 - Save your workflow for reproducibility
- **Data files (.dta):** Stata's native data format. Stata can also import `.csv`, `.xlsx`, and other formats.
- **Log files (.log):** Store both your commands and Stata's output. → Think of them as a record of what you did. Always start a log file at the beginning of each session.
 - You can add comments inside your do-file (which will also appear in the log) using an asterisk `*` at the start of a line.
Example: `* This is a comment and will not be executed.`
 - Alternatively, use `//` for end-of-line comments.
Example: `summarize mpg // quick summary of miles per gallon`

Note

When you have a question or encounter an error in Stata, always save both your **Do-file** and the corresponding **Log file**. Send them together to the teaching fellow or to me — this allows us to *reproduce exactly what you did* and quickly identify where the error occurred.

0.4 Default Log Template

Below is the default template you should use at the beginning of every lab session. It ensures your work is reproducible, well-documented, and easy to debug.

Default Do-file Template

```

1 *****
2 * [Your Name]
3 * Lab Class #1
4 *****
5
6 cls                // clear command window
7 clear all          // remove any existing dataset
8 set more off       // prevent pauses in output
9
10 capture log close // close any open log file
11 cd "C:\Users\YourName\Downloads" // set working directory
12
13 log using "lab1_initials.log", replace // starting log file named "
   lab1_initials"
14
15 * --- your work goes here ---
16 * use comments (*) to explain steps
17
18 log close          // close the log at the end
19 *****

```

Below is the same template again for you to **copy and paste directly**:

```

*****
* [Your Name]
* Lab Class #1
*****

cls                // clear command window
clear all          // remove any existing dataset
set more off       // prevent pauses in output

capture log close // close any open log file
cd "C:\Users\YourPath\Downloads" // set working directory

log using "lab1_initials.log", replace // starting log file named "lab1_initials"

* --- your work goes here ---
* use comments (*) to explain steps

log close          // close the log at the end
*****

```

Week 1 — Introduction to STATA (Teaching)

1.1 Goals for Week 1

1. Change the working directory (`cd`, `pwd`, `ls`).
2. Load datasets (`.dta`, `.csv`)
3. Construct new variables using `gen` and `replace`:
 - Apply conditions with `if`, logical operators (`&`, `|`)
 - Create dummy variables
 - Transform variables (squared, logarithm, interaction terms $A \times B$)
4. Summarize and describe data:
 - `describe`, `summarize`, `tabulate`
5. Keep or drop variables (`keep`, `drop`).
6. Create plots (histogram, scatter plots)
7. Save and export plots.
8. Save datasets for later use.

1.2 Setting for Week 1

In this week's lab, we will work with a *simulated* dataset that represents students who have taken **ECON255**. Each observation corresponds to one student in a given term, identified by a unique **StudentID**. The dataset includes demographic and academic information such as **Gender**, **Year** (Sophomore, Junior, Senior), and whether the student is an **EconMajor** or has a **DoubleMajor**. It also records measures of engagement, including the number of **Absences** and **OfficeHours** attended. Performance is captured by three midterm scores (**Mid1**, **Mid2**, **Mid3**) and a **FinalProject** score. This dataset is designed to illustrate how different academic behaviors (e.g., attending office hours, missing classes) and characteristics (e.g., major, year in school) are related to student outcomes.

1.3 Tasks for Week 1

0 Don't forget to start with your [default do-file template!](#)

1 Change the working directory to Downloads folder

```
1  pwd
2  cd "C:\...\Downloads"
```

2 Load the dataset (You can load either `.dta`, or `.csv` file)

```
1  * if using the DTA version:
2  use "simulated_students.dta", clear
3
4  * Or, if using the CSV version:
5  import delimited "simulated_students.csv", clear
```

3 Get descriptives of the data

- Use `describe(des)` to see an overview of variables.

```
1 des
2 des Mid1 Mid2 FinalProject
```

- Use `list` to inspect rows of data.

```
1 list StudentID Gender Year Mid1 Mid2 in 1/10
2 list StudentID Absences OfficeHours if Year=="Sophomore"
```

- Use `summarize(sum)` for summary statistics.

```
1 sum Mid1 Mid2 Mid3 FinalProject
2 sum Absences OfficeHours, detail
```

- Use `tabulate(tab)` for categorical variables.

```
1 tab Gender
2 tab Year, summarize(FinalProject)
```

- Use `count` to see how many observations meet conditions.

```
1 count if Absences > 5
2 count if OfficeHours > 10 & EconMajor==1
```

- Try logical conditions with `&`(and), `|`(or), `==(equal to)`, `!=(not equal to)`.

```
1 ls StudentID Mid1 Mid2 if Absences > 3 & EconMajor==1
```

4 Create new variables using `gen` and `egen`

- **Final score (0–100)**

Midterms are 25% each; FinalProject is out of 25 points and counts fully.

```
1 gen FinalScore = 0.25*Mid1 + 0.25*Mid2 + 0.25*Mid3 +
   FinalProject
```

- **Letter grade from FinalScore**

```
1 gen LetterGrade = ""
2 replace LetterGrade = "A" if FinalScore >= 90
3 replace LetterGrade = "B" if FinalScore >= 80 & FinalScore <
   90
4 replace LetterGrade = "C" if FinalScore >= 70 & FinalScore <
   80
5 replace LetterGrade = "D" if FinalScore >= 60 & FinalScore <
   70
6 replace LetterGrade = "F" if FinalScore < 60
```

- **Log of office hours (add 1 to handle zeros)**

```
1 gen LogOH = log(OfficeHours + 1)
```

– Squared absences

```
1 gen AbsencesSq = Absences^2
```

– Term average of FinalScore (using egen)

```
1 egen TermAvgScore = mean(FinalScore), by(Term)
```

Tip: Difference between gen and egen

- * **gen** (generate) is used to create a new variable based on a calculation for each observation.
 - Example: `gen LogOH = log(OfficeHours+1)` creates one value of LogOH per student.
- * **egen** (extensions to generate) is used when you need to create variables that rely on **groups of observations** or more complex functions.
 - Example: `egen TermAvgScore = mean(FinalScore), by(Term)` gives every student in the same term the same average score.
- * Rule of thumb: Use **gen** for simple, row-by-row calculations; use **egen** for grouped or advanced statistics.

5 Create and export Plots

– Scatter plot: Office Hours vs. Final Score

```
1 twoway scatter FinalScore OfficeHours, ///
2 title("Office Hours vs. Final Score") ///
3 xtitle("Office Hours Attended") ///
4 ytitle("Final Score") ///
5 name(sc_off, replace)
6 graph export "scatter_officehours_finalscore.png", replace
```

– Histogram of Letter Grades by Term (overlaid)

```
1 twoway ///
2 (histogram FinalScore if Term=="2024SP", width(5) start(50)
3 fcolor(none) lcolor(red)) ///
4 (histogram FinalScore if Term=="2024FA", width(5) start(50)
5 fcolor(blue%40) lcolor(blue)), ///
6 title("FinalScore by Term") ///
7 xtitle("FinalScore") ytitle("Frequency") ///
8 legend(order(1 "2024SP" 2 "2024FA") cols(2) position(11) ring
9 (0) region(lcolor(none))) ///
10 name(h_overlay, replace)
11 graph export "histogram_overlay_by_term.png", replace
```

– Joint Graph (1x2 layout)

```
1 graph combine sc_off h_overlay, cols(2) name(gcombo, replace)
2 graph export "combined_plots.png", replace
```

Tip: Basic Structure of twoway Graphs in Stata

The general syntax is:

```
twoway (plottype yvar xvar [if] [in] [weight], options) ///
      (plottype yvar xvar, options) ... , ///
      graph_options
```

- **plottype**: the kind of plot (e.g., `scatter`, `line`, `histogram`).
- **yvar xvar**: the variables to plot (y-axis first, then x-axis).
- **options**: specific settings for that plot type (e.g., colors, titles).
- Multiple plot types can be combined by stacking them in parentheses, separated by spaces or line breaks with `///`.
- After the plot definitions, you can add **graph_options** to control the overall look (titles, axis labels, legends, etc.).

Example:

```
twoway (scatter y x, mcolor(blue)) ///
      (lfit y x, lcolor(red)), ///
      title("Y vs X") legend(order(1 "Data" 2 "Fit"))
```

This produces a scatter plot of `y` on `x` with a fitted regression line.

6 Save dataset

Save dataset in different formats

```
1  * Save as Stata .dta file
2  save "modified_students.dta", replace
3
4  * Save as CSV file
5  export delimited using "modified_students.csv", replace
```

Week 2 — Running OLS, Fitted Values, and Exporting Results

(Teaching)

2.1 New Goals for Week 2

1. Load and prepare datasets for regression analysis.
2. Run a basic Ordinary Least Squares (OLS) regression using `reg`.
3. Obtain and interpret key regression outputs:
 - Coefficient estimates($\hat{\beta}$) and standard errors($s.e(\hat{\beta})$)
 - R^2 and goodness-of-fit
4. Generate fitted values (`fitted_val - \hat{y}`) and residuals (`resid - \hat{u}`).
5. Summarize regression diagnostics with plots:
 - Residual plots
 - Fitted vs. actual scatter plots
6. Export regression results using `outreg2`.

2.2 Setting for Week 2

For this week we use the dataset `WAGE1.dta`, one of the standard examples from Wooldridge's textbook *Introductory Econometrics: A Modern Approach*. This dataset is a cross-sectional sample of working individuals in the United States. The main focus is on the determinants of hourly wages, with special attention to the role of education. The data contain information on wages, education, experience, demographic characteristics, and job-related variables.

Key variables we will work with include:

- `wage` : hourly wage in U.S. dollars.
- `educ` : years of education completed.
- `exper` : years of work experience.
- `tenure` : years with current employer.
- `female` : indicator variable for gender (=1 if female).
- `married` : indicator variable for marital status (=1 if married).

The dataset is widely used to illustrate the **Mincer wage equation**, which studies the returns to education. In this lab, however, we begin with a simple regression of `wage` on `educ` to practice running OLS, interpreting regression output, and generating fitted values and residuals in Stata.

2.3 Tasks for Week 2

0 Start with your [default do-file template](#).

1 Load the dataset:

```
1 use "WAGE1.dta", clear
```

2 Explore the data (quick diagnostics before OLS)

– Average years of education; min/max

```
1 sum educ
2 di "Mean educ = " %5.2f r(mean) " min = " r(min) " max = "
   r(max)
```

Tip: Displaying Results from summarize

- * `r(mean)`, `r(min)`, `r(max)` : stores results from the last `summarize` command.
- * `%5.2f` : numeric display format — width 5, with 2 decimal places (e.g., 12.3456 → 12.35).
- * You can change the format (e.g., `%8.3f`) to adjust spacing and decimals.

– Wage summary by gender (female vs. male)

```
1 bysort female: sum wage
```

– Distribution of education (histogram)

```
1 hist educ, width(1) start(0) freq ///
2 title("Distribution of Education") xtitle("Years")
```

– How many women vs. men?

```
1 tab female
```

– Education by gender (cross-tab)

```
1 tab educ female
```

– Overlay histograms of education by gender

```
1 twoway ///
2 (histogram educ if female==1, width(1) start(0) fcolor(none)
   lcolor(red)) ///
3 (histogram educ if female==0, width(1) start(0) fcolor(none)
   lcolor(blue)), ///
4 title("Education by Gender") xtitle("Years of Education")
   ytitle("Frequency") ///
5 legend(order(1 "Female" 2 "Male") cols(2) position(11) ring(0)
   region(lcolor(none))) ///
6 name(educ_gender_overlay, replace)
```

3 Run a simple regression

- Estimate the regression of `wage` on `educ`.
- Retrieve and display the estimated intercept and slope coefficient.
- Compute fitted values and residuals from the regression.

Tip: reg Syntax and Saved Results

- Syntax

```
regress depvar indepvars [if] [in] [weight], [options]
```

- Accessing estimates (after reg)

- * Coefficient for `x1`: `_b[x1]`
- * Std. error for `x1`: `_se[x1]`
- * Intercept: `_b[_cons]`
- * R^2 : `e(r2)` Adjusted R^2 : `e(r2_a)`
- * Sample size: `e(N)` Residual df: `e(df_r)`
- * RMSE: `e(rmse)` Model/Residual SS: `e(mss)`, `e(rss)`

- Predictions (after running reg)

- * Fitted values \hat{y} : `predict yhat, xb`
- * Residuals u : `predict uhat, residuals`
- * SE of \hat{y} : `predict se_yhat, stdp`

```
. reg wage educ
```

Source	SS	df	MS		
Model	1179.73204	1	1179.73204	Number of obs	= 526
Residual	5980.68225	524	11.4135158	F(1, 524)	= 103.36
Total	7160.41429	525	13.6388844	Prob > F	= 0.0000
				R-squared	= 0.1648
				Adj R-squared	= 0.1632
				Root MSE	= 3.3784

	Coefficient	Std. err.	t	P> t	[95% conf. interval]	
wage						
educ	.5413593	.053248	10.17	0.000	.4367534	.6459651
_cons	-.9048516	.6849678	-1.32	0.187	-2.250472	.4407687

Figure 2: Regression Results

Note: The top-right panel reports the number of observations, R^2 , and related fit statistics. The coefficient table displays the estimated coefficients, their standard errors, t -statistics (two-sided tests of $H_0 : \beta = 0$), corresponding two-sided p -values, and 95% confidence intervals.

[Question] What is the estimated increase in wage for a person with four more years of education?

4 Plot fitted vs actual, and residual plots:

```
1 twoway (scatter wage educ) (lfit wage educ)
2
3 twoway scatter resid educ, ///
4 title("Residuals vs. Education") ///
5 yline(0, lcolor(red))
```

5 Export regression results:

In this course, we will use the `outreg2` command to export regression tables to Word or other file formats. `outreg2` is not part of Stata's built-in commands, so you must first install the package (only once per computer).

```
1 // Install outreg2 if not already available
2 capture which outreg2 // checks if outreg2 is installed
3 if _rc ssc install outreg2 // if not found (_rc != 0), install
4 // from SSC
5
6 // Run regression and export results
7 reg wage educ
8 outreg2 using "Lab2_results.doc", ///
9 replace ctitle("OLS: wage on educ") dec(3) bdec(3) tdec(3) ///
addstat("R-squared", e(r2))
```

Week 3 — Assignment 1: Housing Prices (Homework)

3.1 Setting

Public real estate data, such as listing prices and housing characteristics, are widely accessible. Many providers even offer programmatic access (e.g., the [Zillow API](#)). In this assignment, you will use the dataset `HPRICE1.dta` from our textbook, which contains information on Boston-area home sales in 1990. Your tasks include: loading the data, generating descriptive statistics, estimating a simple OLS regression of housing price on key characteristics, and exporting both tables and figures.

3.2 What to Submit

Upload your **final and working** versions of the `.do` file, the corresponding `.log` file, and the exported `lab3_summary.doc` and `lab3_results.doc` to Blackboard. Grading will be based only on the files you submit, so please double-check that they are correct and complete before uploading.

3.3 Assignment Tasks

1. Set up your **default do-file template** and change the working directory.
2. Load `HPRICE1.dta` dataset.
3. **Summarize the main variables.** Generate descriptive statistics (mean, standard deviation, minimum, and maximum) for `price`, `sqrft`, and `bdrms`. Export these results into a Word document using the `outreg2` command. Save the file as `lab3_summary.doc`.

Tip: `outreg2` to save Summary Statistics

Type `help outreg2` in the Stata command window to view the official manual. In the section *Example 7. Summary tables*, the subsection *Regular summary table* shows how to save summary statistics (instead of regression results). Be sure to specify the file name with the `.doc` extension so that Stata saves it as a Word document rather than as a plain text file.

4. **Run an OLS regression.** Regress housing price on `sqrft` and `bdrms`.
5. **Save regression results.** Store the estimated coefficients as `beta0hat`, `beta1hat`, and `beta2hat`. Also generate fitted values (call them `pricehat`) and residuals (call them `resid`).
6. **Interpret the coefficient on bedrooms.** Use the `display` command and your saved `beta2hat` to calculate the estimated increase in house price from one additional bedroom, holding square footage constant. Report this in `lab3_summary.doc` using *dollar terms*, not thousands.
7. **Extend the calculation.** Again using `display`, compute the estimated increase in price for a house that has *one additional bedroom and 200 more square feet*. Report this in `lab3_summary.doc` and briefly explain in words why this increase is larger than in the previous step.
8. **Export regression results.** Export the OLS results as `lab3_results.doc` using `outreg2`. Make sure your table includes coefficients, standard errors, t-statistics, p-values, and R^2 .

9. **Plot the fitted regression line.** Create a graph with housing price on the vertical axis and `sqrft` on the horizontal axis. Overlay the fitted regression line using the `pricehat` values you generated earlier, and export the figure to a file. Then include the figure into the lab3_summary.doc

Note: Because the fitted values were estimated using both `sqrft` and `bdrms`, the line plotted against only `sqrft` will not be perfectly smooth—it will appear jagged.

Week 4 — Merging, Reshaping, and Inference (Teaching)

4.1 New Goals for Week 4

1. Acquire public country–year data from the World Bank API.
2. Tidy data by reshaping from *wide-by-year* to *long*, and pivoting to *wide-by-variable*.
3. Merge multiple country–year sources into one analysis dataset.
4. Create log and growth variables and run two macro-style regressions.
5. Export tables and (optionally) simple diagnostic plots.

4.2 Setting for Week 4

In practice, the data required for regression analysis are rarely clean or ready to use. They are often scattered across multiple datasets, recorded in different formats, or contain inconsistencies. To analyze them properly, researchers must first unify their structure and merge them using common identifiers.

In this week, we will practice these skills by using the World Bank API to download data in multiple formats and construct a **country–year panel dataset** with several key indicators.

With the completed dataset, our objective is to study two relationships at the country–year level:

- **Growth & Trade:** Does trade openness predict faster GDP per capita growth?

$$\Delta \log(\text{GDPpc}_{ct}) = \alpha + \beta \frac{\text{Exports}_{ct}}{\text{GDP}_{ct}} + \gamma \text{Education}_{ct} + u_{ct}. \quad (1)$$

- **Unemployment & Growth:** Are richer or high-inflation countries associated with different unemployment rates?

$$\text{Unemp}_{ct} = \alpha + \beta \log(\text{GDPpc}_{ct}) + \gamma \text{Inflation}_{ct} + u_{ct}. \quad (2)$$

We will use the `wbopendata` package in Stata to access World Bank indicators. For more information, type `help wbopendata`. The main variables of interest are:

- `NY.GDP.PCAP.KD` : real GDP per capita (constant 2015 US\$)
- `NE.EXP.GNFS.ZS` : exports of goods & services (% of GDP)
- `NE.IMP.GNFS.ZS` : imports of goods & services (% of GDP)
- `SL.UEM.TOTL.ZS` : unemployment rate (% of labor force)
- `FP.CPI.TOTL.ZG` : inflation, CPI (annual %)
- `SE.SEC.ENRR` : gross secondary enrollment (% , proxy for education)

Note

Recall the syntax of the `reg` command: for example, `reg y x1 x2`. Keeping this in mind, and referring back to equations (1) and (2), what should the structure of our dataset look like?

4.3 Tasks for Week 4

0 Start with your [default do-file template](#) and set your working directory.

1 Download country–year data directly in Stata (World Bank).

- Download GDP, EXP, IMP data as *wide format* and save it as dta file
- Download CPI, UEM, SEC data as *long format* and save it as dta file

```

1  * Install package for World Bank API
2  ssc install wbopendata, replace
3
4  * Download GDP, EXP, IMP as wide format
5  wbopendata, indicator(NY.GDP.PCAP.KD; NE.EXP.GNFS.ZS; NE.IMP.GNFS.
6  ZS) year(1990:2015) clear
7  save "week4_data1.dta", replace
8  * Check with Data Editor to see how this data looks like
9
10 * Download unemployment, inflation, education as long format
11 wbopendata, indicator(FP.CPI.TOTL.ZG; SL.UEM.TOTL.ZS; SE.SEC.ENRR)
12 long year(1990:2015) clear
    save "week4_data2.dta", replace
    * Check with Data Editor to see how this data looks like
    
```

Tip: Long vs Wide

- **Wide format:** Each row is a country–indicator, and each column is a year. Wide format is convenient for humans to read, but hard to merge across different indicators. Example: one row for “GDP per capita in Brazil” with columns yr1990, yr1991, ...
- **Long format:** Each row is a country–year–indicator observation. Long format is the standard for regression analysis (country–year variables stacked in one column). Example: one row for “Brazil, 1990, GDP per capita” and another row for “Brazil, 1991, GDP per capita”.

2 Understanding the structure of first data

For our regression equations (1) and (2), we need the dataset organized at the *country–year* level. That is, each row should represent a unique country–year observation. For each row, we want variables such as $\Delta \log(GDPpc_{ct})$, $\frac{Exports_{ct}}{GDP_{ct}}$, $Education_{ct}$, $Unemp_{ct}$, and $Inflation_{ct}$ as separate columns. The target regression dataset should look conceptually like this:

Country	Year	$\Delta \log(GDPpc)$	Exports/GDP	Education	Unemployment	Inflation
Brazil	1990	0.021	13.5	67.2	6.1	3.5
Brazil	1991	0.018	14.0	68.1	6.3	5.2
India	1990	0.025	9.2	52.7	4.8	8.4
India	1991	0.031	10.1	53.5	5.0	9.0
⋮	⋮	⋮	⋮	⋮	⋮	⋮

Table 1: Desired regression-ready structure: one row per country–year, variables as columns.

Now let us first investigate the first dataset. If you reload the file and open the Data Editor, you will see that the data are stored in a *wide format*, where each indicator has values spread across multiple year columns:

```
1 use "week4_data1.dta", clear
2 browse
```

An excerpt of the data looks like this:

countrycode	indicatorname	yr1990	yr1991	yr1992	yr1993	...
BRA	NY.GDP.PCAP.KD	5211.3	5297.4	5402.1	5534.6	...
BRA	NE.EXP.GNFS.ZS	11.8	12.1	12.5	13.2	...
BRA	NE.IMP.GNFS.ZS	10.6	11.3	12.0	12.8	...
IND	NY.GDP.PCAP.KD	1347.2	1382.5	1425.6	1487.3	...
IND	NE.EXP.GNFS.ZS	7.1	7.5	8.0	8.6	...
⋮	⋮	⋮	⋮	⋮	⋮	...

Table 2: Example structure of `week4_data1.dta` (wide format).

To prepare the dataset for regression, we need to *reshape* the data. Currently, each year appears as a separate column, and indicators are stacked under the `indicatorname` column. Our goal is to convert the dataset so that:

- Each row corresponds to a unique `countrycode` – `year` pair.
- Each indicator (GDP per capita, exports, imports) becomes its own column.

This requires two steps in Stata: (i) reshape from wide to long for years, and (ii) reshape back from long to wide for indicators to become variables.

```
1 * Step 1. Reshape wide -> long (stack years into a single column)
2 * Wide format: separate columns for yr1990, yr1991, yr1992, ...
3 * After this step:
4 * - Those year columns are combined into one variable "year"
5 * - Their values go into a single variable "yr"
6 * - (countrycode, indicatorname) identifies each panel (i)
7 reshape long yr, i(countrycode indicatorname) j(year)
```

After first step, we would have data structure that looks like below.

countrycode	indicatorname	year	yr
BRA	NY.GDP.PCAP.KD	1990	5211.3
BRA	NY.GDP.PCAP.KD	1991	5297.4
BRA	NE.EXP.GNFS.ZS	1990	11.8
BRA	NE.EXP.GNFS.ZS	1991	12.1
⋮	⋮	⋮	⋮

Table 3: Example structure *after* Step 1 (`reshape long`): each row is a country–indicator–year observation, with values stored in `yr`.

Now in the second step, we need to transform the dataset so that each value of `indicatorcode` becomes its own separate variable. In other words, instead of stacking all indicators under one column, we spread them out so that for each country–year we directly observe GDP per capita, exports, imports, and so forth as separate columns.

However, Stata does not allow dots (.) or spaces in variable names. Therefore, before reshaping, we first create short, valid names (e.g. `NY.GDP.PCAP.KD` → `gdp_pc`, `NE.EXP.GNFS.ZS` → `exp`, `NE.IMP.GNFS.ZS` → `imp`). This step ensures that the reshape command works properly.

```

1      * Map raw codes to short, valid names (no spaces/dots)
2      gen str20 ind_short = "" // First create an empty string variable
3      replace ind_short = "gdp_pc" if indicatorcode == "NY.GDP.PCAP.KD"
4      replace ind_short = "exp"    if indicatorcode == "NE.EXP.GNFS.ZS"
5      replace ind_short = "imp"    if indicatorcode == "NE.IMP.GNFS.ZS"
6
7      * Step 2. Reshape long -> wide (make indicators into separate
8      columns)
9      reshape wide yr, i(countrycode year) j(ind_short) string
10
11     * Tidy variable names
12     rename yrgdp_pc gdp_pc
13     rename yrexp    exp
14     rename yrimp    imp

```

Our final dataset now has one observation per country–year, with clean variables for GDP per capita, exports, and imports.

countrycode	year	gdp_pc	exp	imp
BRA	1990	5211.3	11.8	10.6
BRA	1991	5297.4	12.1	11.3
IND	1990	1347.2	7.1	8.5
IND	1991	1382.5	7.5	9.0
⋮	⋮	⋮	⋮	⋮

Table 4: Example structure *after* Step 2 (`reshape wide`): each row is now a country–year observation with indicators in separate columns.

Finally, we save the final version of this data:

```

1      * Saving final version
2      save "week4_data1_final.dta"

```

3 Merging and Renaming

Now, we load the second dataset, which was downloaded in long format (containing unemployment, inflation, and education variables):

```

1      use "week4_data2.dta", clear
2      browse

```

We now merge this dataset with the reshaped file `week4_data1_final.dta`, using the common identifiers `countrycode` and `year`. In Stata, merging requires us to specify the relationship between the datasets:

Tip: Merging Datasets

- 1:1 merge: one observation in master matches exactly one in using dataset.
- m:1 merge: many observations in master match one in using dataset.
- 1:m merge: one observation in master matches many in using dataset.
- m:m merge: many-to-many merge (generally discouraged).

In this exercise, both datasets are organized at the country–year level, so we use a 1:1 merge.

```

1      * Merge the two datasets by countrycode and year
2      merge 1:1 countrycode year using "week4_data1_final.dta"
3
4      * Inspect merge result
5      tab _merge

```

Tip: Interpreting _merge

`_merge==3`: matched in both files.
`_merge==1`: master-only (house has no income match).
`_merge==2`: using-only (income row with no houses).

Finally, we rename the variables for readability and clarity:

```

1      * Clean up variable names for clarity
2      rename sl_uem_totl_zs unemp           // Unemployment rate
3      rename fp_cpi_totl_zg infl           // Inflation (CPI, %)
4      rename se_sec_enrr sec_enroll       // Secondary school enrollment

```

Our final dataset should look like the following:

countrycode	year	gdp_pc	exp	imp	unemp	infl	sec_enroll
BRA	1990	5211.3	11.8	10.6	6.1	3.5	67.2
BRA	1991	5297.4	12.1	11.3	6.3	5.2	68.1
IND	1990	1347.2	7.1	8.5	4.8	8.4	52.7
IND	1991	1382.5	7.5	9.0	5.0	9.0	53.5
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Table 5: Final merged dataset at the country–year level: each row is a country–year observation with indicators as separate variables.

4 Create log and growth variables.

We need the country–year growth rate of real GDP per capita:

$$\Delta \log(\text{GDPpc}_{ct}) \equiv \log(\text{GDPpc}_{ct}) - \log(\text{GDPpc}_{c,t-1}). \quad (3)$$

To compute this correctly, we must take differences *within each country* across adjacent years. In Stata, we declare the panel structure (country as the panel identifier, year as the time index) and then use the time-series difference operator `D.`, which applies differences by panel.

```

1   * 1) Log level
2   gen lgdppc = log(gdp_pc)
3   label var lgdppc "log(GDP per capita)"
4
5   * 2) Declare panel structure so time operators work within country
6   *xtset needs a numeric id, so we will encode countrycode to an id
7   egen cid = group(countrycode), label
8   xtset cid year
9
10  * 3) First difference of log GDP per capita, within each country
11  gen dlgdppc = D.lgdppc
12  label var dlgdppc "Delta log GDP per capita (within-country)"

```

Tip: Why xtset + D.?

- `xtset cid year` tells Stata that observations are ordered by **country** (panel) and **year** (time).
- After `xtset`, the operator `D.var` computes $var_t - var_{t-1}$ for each country separately. It never “bleeds” across countries.
- If you skip `xtset`, `D.` won’t know where one country ends and the next begins, and differences may be wrong.

5 Regression 1: Growth & Trade (export table).

$$\Delta \log(\text{GDP}_{pc_{ct}}) = \alpha + \beta \cdot \frac{\text{Exports}_{ct}}{\text{GDP}_{ct}} + \gamma \cdot \text{Educ}_{ct} + u_{ct}.$$

```

1   * Baseline growth regression
2   reg dlgdppc exports_gdp educ
3
4   * Export table
5   capture which outreg2
6   if _rc ssc install outreg2
7
8   outreg2 using "lab4_results.doc", replace ///
9       ctitle("Growth on Trade & Education") dec(3) bdec(3) tdec(3)
10  ///
11  addstat("R-squared", e(r2), "N", e(N))

```

Note

How should we interpret $\hat{\beta}$ when the explanatory variable `exports_gdp` is measured in percentage points, while the dependent variable is the change in log GDP per capita?

6 Regression 2: Unemployment & Growth (append table).

$$\text{Unemp}_{ct} = \alpha + \beta \cdot \log(\text{GDP}_{pc_{ct}}) + \gamma \cdot \text{Inflation}_{ct} + u_{ct}.$$

```

1      * Levels regression: unemployment on income and inflation
2      reg unemp lgdppc infl
3
4      * Append to the same document
5      outreg2 using "lab4_results.doc", append ///
6          ctitle("Unemployment on log(GDPpc) & Inflation") dec(3) bdec
7          (3) tdec(3) ///
          addstat("R-squared", e(r2), "N", e(N))

```

Note

How should we interpret $\hat{\beta}$ here?

7 Optional quick plots.

```

1      * Scatter + fitted line: growth vs openness
2      twoway (scatter dldgppc exp, msize(vsmall)) ///
3          (lfit dldgppc exp, lcolor(red)), ///
4          title("Growth vs Exports/GDP") ///
5          ytitle("Change in log(GDPpc)") xtitle("Exports/GDP (pct of
6          GDP)")
7
8      graph export "figures\growth_vs_exports.png", replace
9
10     twoway (scatter unemp lgdppc, msize(vsmall)) ///
11         (lfit unemp lgdppc, lcolor(red)), ///
         title("Unemployment vs log(GDPpc)")
         graph export "figures\unemp_vs_lgdppc.png", replace

```

Week 5 — Assignment 2 : Merging Movie Data (Homework)

5.1 Setting

You will work with two simulated datasets about movies. Each file contains different information, and your task is to merge them together for analysis:

- `movies_char.csv`: movie characteristics, including ID, title, genre, and budget.
- `movies_outcomes.csv`: movie outcomes, including ID, release year, box office revenue, and rating.

Your goals are to (i) import and clean each dataset, (ii) merge them by the movie identifier `movie_id`, and (iii) investigate how characteristics such as budget and genre are related to outcomes such as box office performance and ratings. Core variables will include: budget (`budget`), genre (`genre`), box office revenue (`box_office`), and rating (`rating`).

5.2 What to Submit

Upload your **final and working** versions of the `.do` file, the corresponding `.log` file, and the exported `lab5_results.doc` to Blackboard. Grading will be based only on the files you submit, so please double-check that they are correct and complete before uploading.

5.3 Assignment Tasks

1. **Set up** your [default do-file template](#) and change the working directory to your project folder.

2. **Load the movie characteristics data** (`movies_char.csv`)

This dataset contains one observation per movie and includes the following variables:

- `movie_id`: the unique identifier for each movie.
- `title`: the movie title.
- `genre`: the movie genre (e.g., Action, Comedy, Drama).
- `budget`: the production budget (in millions of dollars).

After importing, save the dataset as `char.dta`.

3. **Load the movie outcomes data** (`movies_outcomes.csv`)

This dataset contains one observation per movie and includes the following variables:

- `movie_id`: the unique identifier for each movie.
- `year`: the release year of the movie.
- `box_office`: worldwide box office revenue (in millions of dollars).
- `rating`: the average IMDb-style rating (0–10 scale).

After importing, save the dataset as `outcomes.dta`.

4. **Merge the movie characteristics and movie outcomes data.** Be sure to include the answers to below question in your `lab5_results.doc` once you generate the document file after completing the regression analysis in question 7.

- What was the common identifier for matching?

- Is this 1–1, m–1, or 1–m matching? Why?
- Report the match distribution using `_merge` and the match rate (share of `_merge==3`).

Tip: Hints for merging

- Load the movie characteristics dataset:

```
use "char.dta", clear
```

- Merge with the movie outcomes file. Fill in the merge type and common identifier (replace `merge_type` and `common_identifier` with the correct inputs):

```
merge merge_type common_identifier using "outcomes.dta"
```

- Check how many matches you got:

```
tab _merge
```

- Compute the match rate:

```
count if _merge == 3
count
display "Match rate = " r(N1)/r(N2)
```

- For later analysis, you may want to keep only matched observations:

```
keep if _merge == 3
drop _merge
```

5. **Explore the merged dataset (Part 1).** After merging, each row now represents an *individual movie*, including its characteristics (budget, genre) and outcomes (box office, rating). Next, classify movies into “blockbusters” or “flops” based on their box office relative to their budget:

- Define a **blockbuster** as a movie with box office revenue at least twice its budget.
- Define a **flop** as a movie with box office revenue less than its budget.
- Generate indicator variables for each case (`blockbuster`, `flop`) and tabulate their frequency across genres.
- **[Write in `lab5_results.doc` after your regression]:** Which genre has the highest share of blockbusters? Which genre has the highest share of flops?

Tip: Tabulating shares in Stata

To see the share of blockbusters within each genre, use:

```
tab genre blockbuster, row
```

This reports row percentages: for each genre, the proportion of movies that are blockbusters vs. non-blockbusters.

6. **Explore the merged dataset (Part 2).** Now that you have created indicators for blockbusters and flops, let's focus on how the relationship between budget and box office differs across genres.

- Create a **scatter plot** of `box_office` (y-axis) against `budget` (x-axis) for **Fantasy** movies.
- Overlay a second scatter plot for **Romance** movies on the same graph, using a different color.
- Add a legend and clear axis titles so the two genres can be compared directly.
- Export this graph and include it in `lab5_results.doc` later.
- **[Write in `lab5_results.doc` after your regression]:** How do Fantasy and Romance movies differ in terms of budget size and box office performance?

7. **Estimate the regression model and export the table as `lab5_results.doc`.**

$$\log(\text{box_office}) = \beta_0 + \beta_1 \log(\text{budget}) + \beta_2 \text{rating} + u. \quad (4)$$

(Generate new variables: `lbox_office = log(box_office)`, `lbudget = log(budget)`.)

- **[Write in `lab5_results.doc`]:** Interpret $\hat{\beta}_1$: approximately how much higher is box office revenue if a movie's budget doubles?
- **[Write in `lab5_results.doc`]:** Interpret $\hat{\beta}_2$: how is box office revenue related to a one-point increase in rating (on the 0–10 scale)?

8. **Hypothesis testing** Use your regression model to formally test whether budget and rating matter for explaining box office revenues.

- Conduct a **t-test** of the null hypothesis that the coefficient on `log(budget)` is equal to zero.
(*Just reference t-stat column from regression results!*)
- Conduct an **F-test** of the null hypothesis that the coefficients on both `log(budget)` and `rating` are equal to zero jointly.
(*Use the tipbox below for F-test*)
- **[Write in `lab5_results.doc`]:** Report the test statistics, p-values, and whether you reject the null hypotheses at the 5% significance level. Briefly interpret what these results imply in the context of movies.

Tip: Hints for hypothesis testing in Stata

- When you run a regression, Stata already reports the **t-statistic** for each coefficient. For example:

```
reg lbox_office lbudget rating
```

The t-test for $H_0 : \beta_1 = 0$ (coefficient on `lbudget`) is shown in the regression output.

- To conduct a joint **F-test**, use the `test` command after regression:

```
test x1 x2
```

This tests the null hypothesis $H_0 : \beta_{x1} = \beta_{x2} = 0$ against the alternative that at least one of the coefficients is nonzero.

9. Wrap-up. Your `lab5_results.doc` must include:

- Answers to #4, #5, #6, #7, #8
- Figures from #6
- Regression table from #7
- One factor that might pose endogeneity problem in your regression and short explanation as to why.

Week 6 — Nonlinearities in Regression (Teaching)

6.1 New Goals for Week 6

1. Understand how to incorporate nonlinearities in regression:
 - Quadratic terms
 - Interaction terms
2. Learn how to calculate and interpret marginal effects.
3. Apply results to economic questions such as diminishing returns or heterogeneous effects.

6.2 Setting for Week 6

For this week we use the dataset `twoyear.dta`, which contains individual-level data on wages, education, experience, and demographics.

Key variables:

- `lwage` : log hourly wage (dependent variable).
- `educ` : years of education.
- `exper` : years of work experience.
- `female` : gender indicator (=1 if female).

We focus on the relationship between wages, experience, and gender, illustrating how quadratic and interaction terms are estimated and interpreted.

6.3 Tasks for Week 6

- 0 Start with your [default do-file template](#).
- 1 Load the dataset and explore the main variables

```
1 use "twoyear.dta", clear
2
3 sum lwage exper
4 tab female
```

- 2 Run the following baseline regression and export it using `outreg2` as `Lab3_results.doc`

$$\log(wage_i) = \beta_0 + \beta_1 \text{exper}_i + u_i \quad (5)$$

```
1 * Running simple regression
2 reg lwage exper
```

- 3 **Quadratics: does the return to experience differ by experience level?**

The effect of experience on log wages may be nonlinear, with returns diminishing as experience increases. To capture this, we now run a quadratic regression and append the results to `Lab3_results.doc` using `outreg2`:

$$\log(wage_i) = \beta_0 + \beta_1 \text{exper}_i + \beta_2 \text{exper}_i^2 + u_i \quad (6)$$

```

1 * Allowing decreasing/increasing returns
2 reg lwage c.exper##c.exper

```

Tip: Running Regression with Nonlinearities

There are two common ways to estimate quadratic regressions:

1. Manual approach:

```

1 gen expersq = exper^2
2 reg lwage exper expersq

```

This works for estimation, but when you later use `margins`, Stata does not automatically recognize that `expersq` is a function of `exper`. As a result, marginal effects of experience will be incorrect unless you compute them manually:

$$\frac{\partial \log(wage)}{\partial exper} = \hat{\beta}_1 + 2\hat{\beta}_2 \cdot exper$$

2. ## operator (recommended):

```

1 reg lwage c.exper##c.exper

```

The `##` notation tells Stata to include both the main effect (`c.exper`) and the squared term (`c.exper##c.exper`). Importantly, Stata then knows the functional form, so `margins` correctly computes partial derivatives.

Key syntax:

- * `c.var` : treat `var` as continuous.
- * `i.var` : treat `var` as categorical (indicator/dummy).
- * `c.var1##c.var2` : include both main effects and their interaction.
- * `i.var##c.var` : include main effect of the indicator, main effect of the continuous variable, and the interaction term.

- 4 Now, let's get marginal effect of an additional year of experience when experience is minimum, median, mean, and maximum values.

$$\frac{\partial \log(wage_i)}{\partial exper} = \hat{\beta}_1 + 2\hat{\beta}_2 \cdot exper \quad (7)$$

```

1 * Get key reference points for exper
2 quietly summarize exper, detail
3 local exp_min = r(min)
4 local exp_med = r(p50)
5 quietly summarize exper
6 local exp_mean = r(mean)
7 local exp_max = r(max)
8

```

```

9      * Marginal effects at min / median / mean / max
10     margins, dydx(exper) at( ///
11         exper = ('exp_min' 'exp_med' 'exp_mean' 'exp_max') )
12
13
14     * Plot (x-axis shows the exper values we chose)
15     marginsplot, xdimension(exper) ///
16         title("Marginal Effect of Experience at Key Points") ///
17         ytitle("d ln(wage) / d exper") xtitle("Experience (years)"
18         )
19     * Is the return to experience increasing or decreasing an
20         experience level increases?

```

5 Interactions: does the return to experience differ by gender?

If the payoff to an extra year of experience differs systematically for women and men (different slopes), a single coefficient on `exper` is too restrictive. An interaction allows the slope to vary by group:

$$\log(\text{wage}_i) = \beta_0 + \beta_1 \text{exper}_i + \beta_2 \text{female}_i + \beta_3(\text{female}_i \times \text{exper}_i) + u_i. \quad (8)$$

Here, β_1 is the experience effect for men (baseline), and $\beta_1 + \beta_3$ is the experience effect for women. Specifically,

- * β_1 : marginal effect of experience on log wage for men.
- * β_2 : female–male gap at zero experience.
- * β_3 : difference in the return to experience (women vs. men).

```

1      * Gender-specific returns to experience
2      reg lwage i.female##c.exper

```

Marginal effects with margins:

- * *Effect of one more year of experience, by gender, at chosen experience levels:*

```

1      * Average marginal effect of experience, by gender
2      margins female, dydx(exper)
3
4      * Marginal effect of experience at min/median/mean/max for
5      * each gender (we already saved these earlier)
6      margins female, dydx(exper) at(exper = ('exp_min' 'exp_med' '
7      * exp_mean' 'exp_max'))
8      marginsplot, title("ME of experience on log(wage) by gender")
9
10     * Or plot marginal effect at various experience level
11     margins female, dydx(exper) at(exper=(0(30) 150))
12     marginsplot
13     * Now at which dimension does the return to experience differ?
14     * Is it across experience levels or across genders?

```

6 [Optional] Interactions & Quadratics

Experience often shows diminishing returns, meaning the marginal effect of an extra year is smaller at higher levels of experience. Moreover, this curvature may differ by gender. To allow both the intercept, slope, and curvature to vary by gender, we estimate:

$$\log(\text{wage}_i) = \beta_0 + \beta_1 \text{exper}_i + \beta_2 \text{exper}_i^2 + \beta_3 \text{female}_i + \beta_4(\text{female}_i \times \text{exper}_i) + \beta_5(\text{female}_i \times \text{exper}_i^2) + u_i.$$

```
1 * Gender-specific nonlinear returns to experience
2 reg lwage i.female##c.exper##c.exper
```

Tip: Triple ##

Using `i.female##c.exper##c.exper` tells Stata to include:

- * The main effects: `female`, `exper`, and `exper2`.
- * All interactions: `female × exper`, `female × exper2`.
- * By construction, this setup lets the quadratic wage–experience profile differ by gender.

This is much safer than generating your own squared or interaction terms by hand, because `margins` automatically recognizes the full functional form.

Week 7 — Assignment 3: Nonlinearities in Regression (Homework)

7.1 Setting

You will use the dataset `LOANAPP.dta`, which contains information on mortgage loan applications. Your goal is to examine whether loan approval is influenced by race, gender, and other applicant characteristics, while allowing for nonlinearities such as quadratic terms and interaction effects.

7.2 What to Submit

Upload your **final and working** versions of the `.do` file, the corresponding `.log` file, and the exported `lab7_results.doc` (along with any figures) to Blackboard. Grading will be based only on the files you submit, so please double-check that they are correct and complete before uploading.

7.3 Assignment Tasks

1. **Set up your default do-file template and change the working directory.**
2. **Load the LOANAPP.dta dataset.** Explore the following variables using `des` and `sum`.
 - *white* (race indicator)
 - *male* (gender)
 - Housing expenditure/income ratio (*hrat*)
 - Other obligations/income ratio (*obrat*)
 - Industry unemployment rate (*unem*)
 - Marital status (*married*)
 - Number of dependents (*dep*)
 - Schooling > 12 years (*sch*)
 - Applicant income (*appinc*)
 - Late payment dummies (one, two, more than two) (*mortlat1*, *mortlat2*)

3. **Descriptive analysis.**

Be sure to include these results and figures in `lab7_results.doc` once you generate the document file after completing the regression analysis in question 4.

- Compute the proportion of applicants who are white vs. non-white.
- Calculate the approval rate for each group (four groups: white men, white women, non-white men, non-white women) and the difference in means.

Tip: Getting Group-Specific Averages and Figures

To compute averages by groups, use:

```
bysort group1 group2: summarize var
```

- *Figure exercise:* Above exercise is bit hard to see the full picture quickly. Let's create a bar chart showing approval rates by race and gender (four groups: white

men, white women, non-white men, non-white women). Save it so you can put it inside `lab7_results.doc` later.

Tip: Creating Visually Appealing Figures

Add readable category labels for 0/1 indicators:

```
1 label define race 0 "Non-white" 1 "White", replace
2 label values white race
3 label define sex 0 "Women" 1 "Men", replace
4 label values male sex
```

Keep the y-axis on 0–1 and label it clearly. Show bar labels with 2 decimals:

```
1 graph bar (mean) approve, over(male) over(white) ///
2 yscale(range(0 1)) ylabel(0(.1)1, format(%3.1f) angle
3 (0)) ///
4 blabel(bar, format(%4.2f) position(outside)) ///
5 bargap(20) bar(1, lcolor(black)) ///
6 title("Approval Rate by Race and Gender") ///
7 ytitle("Mean approval (share)") ///
8 note("Mean of 0/1 outcome")
```

4. **OLS regression without nonlinearities.** Estimate the probability of loan approval as a function of the applicant characteristics listed in question 2 and export the regression results to `lab7_results.doc`.

Tip: Wildcard Character ‘*’

In Stata, the asterisk `*` can be used as a *wildcard* in variable lists.

- `prefix*` will include all variables whose names begin with `prefix`.
- `*suffix` will include all variables whose names end with `suffix`.
- `*text*` will include all variables that contain `text` anywhere in their names.

This feature is especially useful when variables follow a consistent naming pattern (e.g., multiple dummies, ratios, or survey items). It helps avoid typing each variable name manually and reduces the chance of typos.

For example, if your dataset has columns named `yr_2001`, `yr_2002`, ..., `yr_2024`, instead of listing them one by one, you can simply type `yr_*`.

5. **Partial effects in linear models.**

- Compute the marginal effect of `appinc` at its minimum, average, and maximum values. (*Hint: This should be the same across different value because we didn't incorporate any non-linearities yet!*)
- Calculate the partial effect of `white` for both male and female applicants.
- Include interpretation of these marginal effects in your `lab7_results.doc`.
(*You are writing interpretation of above questions manually in your document before you submit!!!*)

6. **OLS regression with nonlinearities.** Re-estimate the probability of loan approval,

this time including the interaction between *white* and *male*, as well as a quadratic term for the applicant's income. Append these regression results to `lab7_results.doc`. (Use the `##` operator so that we can use `margins` command.)

7. **Nonlinearities and partial effects.**

- Compute the marginal effect of `appinc` at its minimum, average, and maximum values. (*Hint: Now that we have non-linearities, these should be different!!*)
- Calculate the partial effect of *white* for both male and female applicants.
- Include interpretation of these marginal effects in your `lab7_results.doc`.
(*You are writing interpretation of above questions manually in your document before you submit!!!*)

8. **Inference with linear combinations.** Using your regression with nonlinearities, test whether the effect of being white differs significantly between men and women.

- Recall: in the model with *white*, *male*, and *white* \times *male*, the effect of being white is:
 - * For women: coefficient on *white*.
 - * For men: coefficient on *white* + coefficient on (*white* \times *male*).
- Use the `lincom` command to test whether the white–nonwhite gap is the same for men and women.
- Interpret the result in `lab7_results.doc` : is there evidence that the impact of race on loan approval differs by gender?

9. **Wrap-up.** Your `lab7_results.doc` must include:

- Answers to #3, #5, #7, #8
- Figure from #3
- Regression result from #4, #6
- **Potential Endogeneity.** Where might endogeneity arise in our loan approval regression? Discuss possible sources and why they could bias the OLS estimates.

Week 8 — Introduction to Panel Data and Lagged Variables

(Teaching)

8.1 New Goals for Week 8

1. Learn practical tools for **model robustness**:
 - Compare alternative functional forms.
 - Use **heteroskedasticity-robust standard errors** and understand why they matter.
2. Begin working with data that have both a **cross-sectional** and a **time** dimension:
 - Recognize what panel data are and how they differ from a single cross-section.
 - Understand why many economic questions require tracking the same units over time.
 - Create and interpret a one-period lag of the dependent variable as a first step toward panel analysis.

8.2 Setting for Week 8

We will use the simulated student–year panel `lab8_students_panel.dta` (balanced panel with 600 students observed for 4 years). This is the first lab in which we explicitly use the **time dimension** of the data. While previous weeks analyzed only cross-sectional variation, many empirical studies in economics require examining how the same individuals or firms change over time. Such datasets—called *panel data*—allow researchers to investigate both *cross-sectional* and *time-series* variation simultaneously. Starting this week, we will begin to utilize this structure by creating and interpreting lagged variables, setting the stage for full panel techniques in later weeks.

Key dataset features

- Balanced panel: `student` (id), `year` (1–4).
- Major (categorical, 7 values): Economics, ComputerScience, Business, Biology, Education, Sociology, Psychology.
- Employment indicators:
 - * `employed` : 1 if employed in that year, 0 otherwise.
 - * `underemp` : 1 if employed but underemployed (major-specific probability).
- Wages:
 - * `logw` : log(wage) for employed observations (missing if unemployed).
 - * `wage` : level wage (0 if unemployed).
- Human-capital and demographics:
 - * `exper` : years of experience.
 - * `tenure` : an approximate tenure measure.
 - * `age` : age in each year.
 - * `female` : gender indicator.
 - * `race` : categorical race indicator.

8.3 Tasks for Week 8

0 Start with your [default do-file template](#).

1 **Load and explore the dataset.** The data are stored as a balanced panel, meaning each student (`id`) is observed for four years. Begin by loading the dataset and examining its main variables.

```
1 use "lab8_students_panel.dta", clear
2
3 describe
4 summarize logw wage exper tenure age
5 tab major
6 tab year
```

Discuss: What variables vary over time? Which are time-invariant?

2 **Check the panel structure.** In Stata, panel data must be *declared* using `xtset`, which tells Stata the unique identifier for each entity and the time variable.

```
1 * Declare panel structure: (id = student identifier, year =
   time)
2 xtset id year
3
4 * Describe the panel characteristics
5 xtdescribe
```

Stata now knows that observations are ordered by student and year. The output of `xtdescribe` confirms that this dataset is **strongly balanced**—each student has 4 yearly observations.

Tip: Panel terminology

- * **Cross section:** a single snapshot across many individuals (e.g., one year).
- * **Time series:** one individual observed over many periods.
- * **Panel:** combines both—many individuals observed over time.

3 **Visualize average trajectories by group (race).** To illustrate how outcomes move over time at the group level, compute and plot the mean log wages among employed by race across years.

```
1 * Save working data, compute group-year means, and restore
   original data
2 preserve
3
4 * 1) Mean log wage by race (condition on employed so logw exists)
5 collapse (mean) mean_logw = logw if !missing(logw), by(race year)
6
7 * Small-multiples: one panel per race
8 graph twoway line mean_logw year, ///
9     by(race, cols(3) compact title("Mean log(wage) by race (
   employed only)")) ///
10    ytitle("Mean log(wage)") xtitle("Year") note("Conditioned on
   employed (logw non-missing)")
```

```

11
12 graph export "mean_logw_by_race_year.png", replace
13
14 restore

```

Tip: Why preserve/restore when using collapse

Using `collapse` replaces the current dataset with a new aggregated dataset (one row per group). That *permanently* removes the original observation-level data in memory unless you save it first. In class we want to compute group-year summaries for plotting without destroying the full panel dataset we will continue to use.

* Quick and safe (recommended):

```

1 preserve
2 collapse (function) newvar = newvar_func , by(groupvar
   timevar)
3 restore

```

Here, `function` can be any summary operator such as `mean`, `sum`, `sd`, or `count`, and `by()` lists the grouping variables that define the aggregation. This temporarily swaps in the aggregated data and then restores the original dataset exactly as it was.

Short reminders:

- * Don't call `xtset` on the collapsed (group-level) data — it no longer has the original `id/year` structure.
- * `preserve/restore` is easy and safe for short in-class analyses; use `collapse + save` if you want to keep the aggregated table for reuse.

4 Wages among employed vs. overall average wage by major.

To see how employment status affects average earnings, create a variable that equals wage only for employed students. Then, compare the mean wage among employed students ($E[\text{wage} \mid \text{employed} = 1]$) with the overall mean wage (which includes zeros for unemployed students).

```

1 * Create a variable for wage among employed students
2 gen wage_emp = wage if employed==1
3
4 * Compare mean wage among employed vs. overall mean wage
5 tabstat wage_emp wage , by(major) statistics(mean) ///
6     columns(variables) format(%9.2f)

```

This table reports two columns for each major:

- * `wage_emp`: mean wage among employed students only (conditional mean)
- * `wage`: mean wage across all students, including zeros for unemployed (unconditional mean)

Discuss: Which majors have higher average wages among employed students? Are those rankings the same when you include non-employed students in the average?

5 Baseline regression and heteroskedasticity.

Now estimate a pooled wage regression using employed observations only. Because variance of residuals may increase with experience or differ by major, we use `, robust` to obtain heteroskedasticity-consistent standard errors. Then, plot residuals versus `exper` (experience) to visually check for heteroskedasticity by examining the distribution of residuals conditional on an independent variable.

```

1  * Baseline regression using employed observations
2  reg logw i.major exper tenure female i.race if employed==1, robust
   baselevels
3  * The option "baselevels" forces Stata to display the omitted (
   base) categories in the regression table.
4  * For example, the output should show "(base)" next to Economics (
   major==1) and White (race==1).
5  * All other category coefficients are interpreted relative to
   these base groups.
6
7  * Generate residuals
8  predict resid, resid
9
10 * Plot residuals versus experience (more direct for our DGP where
   sd grows with exper)
11 twoway (scatter resid exper, msymbol(o) msize(small) ///
12         (lfit resid exper), ///
13         ytitle("Residuals") xtitle("Experience (years)") ///
14         title("Residuals vs Experience") ///
15         note("Increasing spread in residuals across exper indicates
16             heteroskedasticity; robust SEs needed")
17 graph export "resid_vs_experience.png", replace

```

Tip: Interpreting residuals conditional on an independent variable

To detect heteroskedasticity visually, examine how the distribution (vertical spread) of residuals changes *conditional* on an independent variable such as experience or education. If the spread grows (a “fan” shape) as the conditioning variable increases, this indicates heteroskedasticity: $\text{Var}(u | X)$ depends on X . Using `, robust` produces consistent standard errors in this case.

6 Lagged variable and clustering intuition.

Next, we examine persistence by adding a lagged dependent variable. This helps us understand how outcomes today are related to past outcomes, and why repeated observations of the same student may have correlated errors.

```

1  * Create one-year lag of log wage within each student
2  bysort id (year): gen L1_logw = logw[_n-1]
3
4  * Check the correlation between current and lagged log wages
5  corr logw L1_logw if year>1 & employed==1
6
7  * Regression with lagged dependent variable
8  reg logw L1_logw i.major exper tenure female i.race ///

```

```

9     if employed==1 & year>1, robust
10
11 * Visualize persistence (strong correlation between L1_logw and
    logw)
12 twoway (scatter logw L1_logw if year>1 & employed==1, msymbol(o)
    msize(small)) ///
13     (lfit logw L1_logw if year>1 & employed==1), ///
14     ytitle("Current log(wage)") xtitle("Lagged log(wage)") ///
15     title("Persistence of Wages") ///
16     note("Strong correlation across years within the same
    student")
17 graph export "logw_vs_L1logw.png", replace

```

Tip: Understanding the `bysort id (year)` prefix

The prefix `bysort id (year)`: tells Stata to:

1. **Sort** the data by `id` and then by `year` within each individual.
2. **Execute** the command (after the colon) separately within each `id` group.

Generalized pattern for panel data operations:

```
bysort panel_id (time_variable): gen newvar = oldvar[_n - k]
```

Where:

- * `panel_id` = unit of observation (firm, person, region, etc.)
- * `time_variable` = time dimension (year, quarter, month, etc.)
- * `oldvar` = the variable you are lagging
- * `k` = number of periods back (1 for a one-period lag)

Examples:

```

1 bysort firmid (year): gen L1_sales = sales[_n-1]    // one-year
    lag of sales per firm
2 bysort country (year): gen L2_gdp = gdp[_n-2]      // two-year
    lag of GDP per country

```

Tip: From Heteroskedasticity to Clustering: Why We Need `, robust` and `, cluster(id)`

Heteroskedasticity means that the noise term (u_i) does not have a constant variance across observations—some data points are “louder” than others. This variation typically occurs *across individuals or firms*.

- * Example: In a wage regression, workers with more experience or education may have more variable earnings than entry-level workers.
- * We address this issue using the option `, robust`, which allows Stata to compute standard errors that are valid even when the error variance differs across observations.

Clustering (Serial Correlation) arises once we introduce a **time dimension**. In panel data, we observe the same individual (or firm) repeatedly, and their unobserved characteristics can persist over time. For instance, when we saw that wages are highly persistent, it indicates that some part of the residual variation—perhaps coming from ability, motivation, or background—remains similar from year to year.

- * Example 1: The same student’s unobserved ability affects wages in every year, so if their error is high in Year 1, it is likely high in Year 2 as well.
- * Example 2: A firm’s unobserved management quality influences profits across multiple quarters—errors move together within the firm.
- * Example 3: People living in the same city experience common shocks such as weather or local market conditions, generating correlated errors within that city.

Because these within-group observations are not independent, we need to adjust our inference to account for this *within-cluster correlation*. Clustering the standard errors by a group variable (here, the student identifier) tells Stata:

“Treat all years of the same student as one bundle: allow errors to be correlated within that bundle but independent across different students.”

Week 9 — Assignment 4: Difference-in-Differences (DID) (Homework)

9.1 Setting

Difference-in-differences (DID) compares how an outcome changes over time in a group that receives a treatment (the *treated* group) to how the outcome changes in a group that does not (the *control* group). The key identifying assumption is **parallel trends**: in the absence of treatment, the average outcomes for treated and control groups would have followed the same time path.

A simple motivating example. Imagine a state-level policy (e.g., a minimum-wage increase) implemented in a subset of counties in 2005. We observe county-level outcomes (e.g., average hourly wage) for several years before and after 2005. Some counties (“treated”) are exposed to the policy starting in 2005; others (“controls”) never receive the policy. DID asks: how much did wages in treated counties change after the policy, relative to the change in control counties?

Note

Notation. Let Y_{it} be the outcome for unit (county) i at time t . Define

$$D_i = \begin{cases} 1 & \text{if unit } i \text{ is in the treated group,} \\ 0 & \text{if unit } i \text{ is in the control group,} \end{cases} \quad \text{and} \quad \text{Post}_t = \begin{cases} 1 & \text{if } t \text{ is a post-treatment period,} \\ 0 & \text{if } t \text{ is a pre-treatment period.} \end{cases}$$

A simple DID regression (without fixed effects) can be written as:

$$Y_{it} = \beta_0 + \beta_1 D_i + \beta_2 \text{Post}_t + \beta_3 (D_i \times \text{Post}_t) + \varepsilon_{it}. \tag{9}$$

The DID estimate of the treatment effect is the interaction coefficient, $\hat{\beta}_3$.

2-by-2 table. The table below illustrates the logic of DID: take the change (After – Before) for each group, then difference those changes.

	Before	After	Difference (After–Before)
Treated ($D_i = 1$)	$\beta_0 + \beta_1$	$\beta_0 + \beta_1 + \beta_2 + \beta_3$	$\Delta_T = \beta_2 + \beta_3$
Control ($D_i = 0$)	β_0	$\beta_0 + \beta_2$	$\Delta_C = \beta_2$
Difference (Treated–Control)	$\Delta_B = \beta_1$	$\Delta_A = \beta_1 + \beta_3$	$\widehat{\text{DID}} = \beta_3$

Table 6: DID decomposition corresponding to regression equation (9).

As you can see from the table, in equation (9):

- β_1 captures any **pre-existing difference in levels** of Y between treated and control groups before the policy.
- β_2 represents the **common time trend** that affects both groups equally (for example, general economic growth).

- β_3 measures the **additional change in Y for the treated group** beyond the common time trend — the causal effect of treatment under the parallel trends assumption.

9.2 What to Submit

Upload your **final and working** versions of the `.do` file, the corresponding `.log` file, and the exported `lab9_results.doc` to Blackboard. Grading will be based only on the files you submit, so please double-check that they are correct and complete before uploading.

9.3 Assignment Tasks

We will use a simulated monthly panel dataset called `lab9_soda_tax_panel.dta`. The data represent 150 counties observed from 2010–2019. In 2015, a group of counties introduced a **soda tax** aimed at reducing sugary drink consumption. Other counties did not adopt the tax and will serve as comparison groups.

Each observation corresponds to a county–month pair and includes information on soda sales and several local characteristics. Key variables are:

- `group` – indicates the county group: there are three groups in total. One group implemented the soda tax (`treated`), while the other two did not (comparison groups).
- `soda_sales_pc` – monthly soda sales per capita (liters per person).
- `treated` – indicator for whether the county implemented the soda tax (1 = treated, 0 = otherwise).
- `post` – indicator for months in 2015 or later (post–policy period).
- `income_pc_month` – average monthly income per person (USD).
- `unemp_rate` – county unemployment rate (percent).
- `store_count` – number of retail outlets selling soda.
- `price_soda` – average local soda price (USD per liter).
- `diabetes_rate` – annual diabetes prevalence rate (percent).

Our goal is to use the **difference-in-differences (DID)** framework to estimate how the soda tax affected soda sales, test key assumptions such as parallel trends, and visualize the results. You will begin by exploring the dataset, checking pre-trends across groups, estimating the DID regression, and then extending the analysis to clustered inference and event-study plots.

1. **Set up your working environment.** Start with your [default do-file template](#). Change the working directory to the folder where you saved `lab9_soda_tax_panel.dta`. Before starting, ensure that the following Stata packages are installed:

```
1 ssc install outreg2, replace
2 ssc install coefplot, replace
```

2. **Explore and describe the data.** Load the dataset and examine its structure and main variables.

```

1 describe
2 summarize soda_sales_pc diabetes_rate income_pc_month unemp_rate
   store_count
3 tab group

```

Answer the following questions in your `lab9_results.doc` file:

- What are the three groups in the data, and how many counties belong to each?
(*Hint : We have monthly data for 10 years*)
- Which variables change over time within a county (i.e., vary within panel units)?
- Report the average monthly soda sales for each group (using `collapse` or `table`).

3. **Visual check: parallel trends.** Before running the DID regression, we first need to check whether the **parallel-trends assumption** appears reasonable. Compute and plot the group means of `soda_sales_pc` over time. (This corresponds to $E[Y_{it} | D_i, Post_t]$). Compare the treated group to each control group using the code template below. Fill in the blanks where indicated.

```

1 preserve
2 collapse (mean) soda_sales_pc, by(year group)
3 twoway (line -- -- if group=="__", lcolor(red)) ///
4        (line -- -- if group=="__", lcolor(blue)) ///
5        (line -- -- if group=="__", lcolor(green)), ///
6        legend(order(1 "__" 2 "__" 3 "__")) ///
7        xline(2015, lpattern(dash) lcolor(black)) ///
8        title("Parallel Trends Check") ///
9        ytitle("Average soda sales per capita")
10 graph export "parallel_trends.png", replace
11 restore

```

In your `lab9_results.doc`, include the following:

- Which control group appears to satisfy the parallel-trends assumption? Briefly explain your reasoning.
 - Insert and label the exported graph (`parallel_trends.png`).
4. **Baseline DID regression.** We are going to modify equation (9) slightly to match our setting. Equation (9) was based on a simple two-group and two-period DID framework. However, our dataset contains multiple treated and control counties observed over many years. In this case, we generalize the model to the following two-way fixed effects DID specification:

$$\text{soda_sales_pc}_{it} = \alpha_i + \delta_t + \beta(\text{treated}_i \times \text{post}_t) + u_{it}. \quad (10)$$

Here, α_i represents a dummy variable for each county (**county fixed effect**), which captures all time-invariant differences across counties — such as local preferences for soda, climate, or population characteristics. The δ_t terms represent dummy variables for each year (**year fixed effects**), which capture common shocks that affect all counties in a given year — such as national health campaigns or changes in soda prices. Finally, the DID estimator β measures the **causal effect of the soda tax**: how much more (or less) soda sales changed in treated counties after 2015 relative to control counties, after removing county and year effects.

To estimate equation (10), follow the steps below.

1 Drop observations that belong to the non-chosen control group.

```
1 * keep treated and the chosen control group
2 keep if inlist(group,"treated","___")
```

2 Create a dummy variable called `post` equal to 1 if the year is 2015 or later.

```
1 gen post = (write condition here)
```

3 Declare the panel structure to Stata using the `xtset` command, and read the tipbox below before running your first fixed-effects regression.

```
1 * Generate monthly time variable for panel declaration
2 -> We have Year, and Month. We need Year-Month
3 gen ym = ym(year, month)
4
5 * This makes newly created ym be a "time" variable to STATA
6 format ym %tm
7
8 xtset __ __
```

Tip: Quick guide: `xtreg`, fixed effects DID

Generic syntax:

```
xtset panel_id time_var
xtreg y x1 x2 i.timevar, fe vce(cluster id)
```

What `xtreg, fe` does:

- * Estimates the model with α_i (unit fixed effects) absorbed. This removes all time-invariant unobserved heterogeneity across counties (e.g., persistent taste for soda or local climate differences).
- * Including `i.year` adds year fixed effects δ_t , capturing common time shocks that affect all counties (e.g., national soda trends).
- * Use `vce(cluster county_id)` to allow serial correlation within counties.

4 Run the baseline two-way fixed effects DID regression and export the table as `lab9_results.doc`.

```
1 * Run baseline two-way fixed effects DID regression
2 xtreg soda_sales_pc post##treated i.year, fe vce(cluster
3 county_id)
4 * i.year adds year dummies (delta_t) to control for common
5 time shocks
6 * vce(cluster county_id) gives clustered standard errors
7
8 outreg2 using "lab9_results.doc", replace ctitle("Original")
```

Interpretation task (write in `lab9_results.doc`): Interpret the coefficient on the interaction term `post##treated` as the causal effect of the soda tax. Explain whether the sign (positive or negative) makes sense.

Note

Why didn't we include other covariates in the baseline? We can, but the two-way FE DID (county FE + year FE) identifies the causal effect under the parallel-trends assumption using within-county changes over time. Time-varying covariates (e.g., `income_pc_month`, `unemp_rate`, `price_soda`, `store_count`) can be added to improve precision or adjust for pre-determined confounders, but do not “fix” a violated parallel-trends assumption. If you include covariates, prefer pre-treatment or lagged versions and cluster SEs at the county level.

5. **Inference: why clustering matters.** Let's first visually understand what clustering means.

```

1 * Visualize average residuals by county
2 * For xtreg , e gives residuals
3 predict resid, e
4
5 bysort county_id: egen mean_resid = mean(resid)
6 hist mean_resid, title("Distribution of county mean residuals")
7
8 graph export ....

```

If the error terms (u_{it}) were *independent across counties and over time*, the mean residuals by county should:

- fluctuate randomly around zero, and
- have a distribution tightly centered around zero.

If the histogram instead shows that:

- many counties have **consistently positive or negative mean residuals**, or
- the distribution is **wide or skewed**,

then **residuals are systematically correlated within counties**. This means certain counties tend to have residuals that are mostly positive (consuming more soda than predicted) and others mostly negative (consuming less than predicted). In other words, residuals are not independent within clusters. If a county has unobserved features that make soda consumption persistently high, these features will cause its residuals to remain correlated over time.

When this happens, the standard errors from the usual heteroskedasticity-robust specification are **too small**, and the resulting t -statistics are **too large**. To correct this, you should cluster standard errors at the county level. `vce(cluster county_id)` allows arbitrary correlation of errors within each county while assuming independence across counties.

Question (write in lab9_results.doc): The histogram shows that several counties have mean residuals noticeably different from zero (i.e., some density at the extremes). What does this imply about the independence of residuals across time within a county? Based on this evidence, should we be concerned about clustering?

6. **Placebo test.** To further check the credibility of your DID estimate, run a **placebo test** by pretending that the soda tax was implemented two years earlier (in 2013 instead of 2015). If the DID framework is valid and the parallel-trends assumption holds, **you**

should not find a statistically significant treatment effect when the “fake” policy date is used.

```

1  * Create a fake post-treatment dummy (pretend policy started in
   2013)
2  gen fake_post = (condition here)
3
4  * Re-estimate the DID regression using the fake treatment timing
5  xtreg ...
6
7  * Export results for comparison
8  outreg2 using "lab9_results.doc", append ctitle("Placebo")

```

Note

If the placebo DID estimate (interaction `fake_post#treated`) is small and statistically insignificant, it suggests that there were no pre-existing differences in soda sales trends between treated and control counties before the actual 2015 policy. A significant placebo effect, however, would cast doubt on the parallel-trends assumption.

7. **Extension: link to health outcome.** Aggregate data to the annual level and test whether lower soda sales correspond to lower diabetes rates.

```

1  preserve
2  collapse (mean) soda_sales_pc diabetes_rate income_pc_month
   unemp_rate, by(county_id year group)
3  regress diabetes_rate soda_sales_pc income_pc_month unemp_rate,
   vce(cluster county_id)
4  restore

```

Interpret the coefficient on `soda_sales_pc`. Does it have the expected sign?

8. **Wrap-up:** Include in `lab9_results.doc`:
- Answers to #2, #3, #4, #5, #7
 - Graphs from #2, #5
 - Regression result tables from #4

Week 10 — Fixed Effects (Within Estimator) (Teaching)

10.1 New Goals for Week 10

1. Understand why and when we use **fixed effects** in panel data:
 - Account for unobserved time-invariant heterogeneity correlated with regressors.
 - Interpret the fixed-effects (within) estimator as using only *within-unit* variation.
 - Recognize that including a full set of entity dummies (the *Least Squares Dummy Variable* or LSDV approach) produces the same coefficient estimates as the within estimator.
2. Learn the practical Stata workflow for fixed-effects estimation:
 - Declare a panel (`xtset`) and convert string IDs (`encode`).
 - Run `reg` with entity dummies, then `xtreg`, `fe`, and compare results.
 - Include time fixed effects and add lags to capture delayed responses.
 - Perform diagnostics: joint test for time FE, groupwise heteroskedasticity, serial correlation, and cross-sectional dependence.
3. Know alternatives and computational shortcuts:
 - Use `reghdfe` to absorb many fixed effects efficiently.
 - When and why to cluster standard errors by panel id.

10.2 Setting for Week 10

We will use the panel dataset `lab10_fixed_panel.dta`, which contains data for multiple countries observed over several years. The goal is to estimate how within-country changes in trade, labor, and related macroeconomic variables are associated with changes in GDP per capita. This lab focuses on applying fixed-effects estimators and conducting diagnostic tests in Stata.

Key dataset features

- Panel identifiers:
 - * `country` – country name.
 - * `country_code` – ISO country code.
 - * `year` – calendar year.
- Main economic variables:
 - * `gdppc` – GDP per capita (current US\$).
 - * `export` – export value index (2015 = 100).
 - * `import` – import value index (2015 = 100).
 - * `trade` – total trade index (sum of exports and imports).
 - * `labor` – labor force participation rate (% of population ages 15+).
 - * `capital` – gross capital formation (current US\$).
 - * `pop_growth` – population growth (annual %).
 - * `inflation` – inflation rate based on GDP deflator (annual %).
 - * `school_enroll` – tertiary school enrollment rate (% gross).
- Data source: World Development Indicators (World Bank), years 2000–2024.
- The panel may be unbalanced due to missing observations in some variables or years.

10.3 Tasks for Week 10

0 Start with your [default do-file template](#).

1 Prepare the panel and declare the panel structure.

Load the dataset `Week10_long_panel.csv`, convert the string identifier `country_code` into a numeric id, and declare the panel in Stata.

```

1 * Load the dataset (CSV format)
2 import delimited "Week10_long_panel.csv", clear
3
4 * Encode string country code to numeric id
5 encode country_code, gen(country_id)
6
7 * Declare panel structure: panel id = country_id, time = year
8 xtset country_id year
9
10 * Quick panel summary
11 xtdescribe

```

Tip: Encoding the panel identifier

The variable `country_code` is a string (e.g., “USA”, “KOR”). Stata’s panel commands require a numeric identifier, so we use:

```
encode country_code, gen(country_id)
```

This creates a numeric id while preserving the original string labels, allowing `xtset` to recognize the panel structure correctly.

2 Construct key variables and explore the data.

```

1 * Construct total trade (export + import)
2 gen trade = your code here
3
4 * Create log-transformed variables
5 gen ln_gdppc = your code here
6 gen ln_trade = your code here
7 gen ln_labor = your code here
8 gen ln_capital = your code here
9
10 * Basic pooled summary
11 sum gdppc trade labor
12
13 * Panel summary: overall / between / within variation
14 xtsum gdppc trade labor
15
16 * Histograms: GDP per capita (level vs log)
17 hist gdppc, normal name(h1, replace) ///
18     title("GDP per capita (level)")
19
20 hist ln_gdppc, normal name(h2, replace) ///
21     title("log(GDP per capita)")

```

```

22
23 graph combine h1 h2, col(2) ///
24     title("Distribution of GDP per capita (level vs log)") ///
25     note("Visual comparison shows that logging reduces skewness")

```

Tip: Interpreting xtsum

The command `xtsum` reports three components of variation for each variable:

- * **Overall:** total variation across all observations.
- * **Between:** variation across units (countries).
- * **Within:** variation within each unit over time.

The fixed-effects estimator relies on the *within* variation. If within variation is small, FE estimates may be imprecise or unstable.

3 Visualize trends in GDP per capita across time.

Before running regressions, it is helpful to **visualize how GDP per capita evolves over time**. This allows us to see whether growth patterns differ across countries and whether global income levels move together. For example, if most countries experience a similar upward trend, this suggests that **time effects (such as global business cycles or common shocks)** may be important to control for later with year fixed effects.

```

1  * Visualize cross-year variation in GDP per capita
2  preserve
3
4  * Compute yearly mean GDP per capita across countries
5  bysort year: egen mean_gdppc = mean(gdppc)
6
7  * Scatterplot of GDP per capita by year, with yearly mean line
8  twoway (scatter gdppc year, msymbol(circle_hollow) color(gs14))
9         ///
10         (connected mean_gdppc year, sort msymbol(diamond)), ///
11         title("GDP per capita over time") ///
12         xtitle("Year") ytitle("GDP per capita (current US$)") ///
13         legend(label(1 "Country observations") label(2 "Yearly mean
14                ")) ///
15         note("The upward trend in the mean line indicates global
16                income growth.")
17
18 graph export "gdppc_mean_by_year.png", replace
19
20 restore

```

4 Fixed Effects Estimation (stepwise)

In this section, we explore how within-country changes in key macroeconomic variables relate to within-country changes in GDP per capita. We begin by estimating a simple pooled OLS model that **ignores** country-specific effects, and then progressively introduce fixed effects to account for unobserved heterogeneity. For now, we ignore non-linearities

and lag effects so that we can focus on understanding the mechanics and interpretation of fixed effects.

$$\ln(\text{gdppc}_{it}) = \alpha_i + \beta_1 \ln(\text{trade}_{it}) + \beta_2 \ln(\text{labor}_{it}) + \beta_3 \ln(\text{capital}_{it}) \\ + \beta_4 \text{pop_growth}_{it} + \beta_5 \text{inflation}_{it} + \beta_6 \text{school_enroll}_{it} + u_{it}.$$

The term α_i represents all unobserved, time-invariant characteristics of each country that may influence its level of GDP per capita—factors such as *geography, colonial history, institutional quality, or long-run technological capacity*. By including α_i , we allow each country to have its own baseline intercept, acknowledging that nations differ in ways that do not change over time but still affect economic performance. If these unobserved characteristics are correlated with any of the explanatory variables, omitting α_i would bias the estimated coefficients. Fixed-effects estimation effectively removes α_i by focusing on how changes within a country over time relate to changes in its explanatory variables, thereby isolating the *within-country* variation that identifies β_1, \dots, β_6 .

1. Pooled OLS (no fixed effects).

Begin by estimating a pooled OLS model, treating all observations as independent. This ignores any country-specific effects (α_i) and assumes all heterogeneity is captured in the error term.

```
1 * 4.1 Pooled OLS regression (ignores country-specific effects)
2 reg ln_gdppc ln_trade ln_labor ln_capital pop_growth inflation
   school_enroll, robust
3 outreg2 using week10_results.doc, replace ctitle("Pooled OLS")
```

Tip: Interpretation

In this pooled regression, unobserved country-specific factors such as geography, institutions, or natural resources are absorbed into the error term. If these factors are correlated with the regressors, OLS estimates will be biased.

2. LSDV: OLS with country dummies (explicit fixed effects).

Add a full set of country dummies (`i.country_id`) to estimate the Least Squares Dummy Variable (LSDV) model. This explicitly controls for time-invariant country effects.

```
1 * 4.2 LSDV regression (OLS with country dummies)
2 reg ln_gdppc ln_trade ln_labor ln_capital pop_growth inflation
   school_enroll i.country_id, ///
3 vce(cluster country_id) hascons
4 outreg2 using week10_results.doc, append ctitle("LSDV") ///
5 keep(ln_trade ln_labor ln_capital pop_growth inflation
   school_enroll)
```

Tip: LSDV vs. Within Estimator

The LSDV model estimates the same slope coefficients (β) as the fixed-effects (within) estimator but explicitly includes one dummy variable for each country. This can become computationally intensive when the number of groups is large, which is why we typically use `xtreg, fe`.

3. Within estimator using `xtreg, fe`.

Next, estimate the fixed-effects model using the within transformation.

```

1 * 4.3 Fixed effects (within estimator)
2 xtreg ln_gdppc ln_trade ln_labor ln_capital pop_growth
   inflation school_enroll, fe robust
3 outreg2 using week10_results.doc, append ctitle("xtreg") ///
4   keep(ln_trade ln_labor ln_capital pop_growth inflation
   school_enroll)

```

Tip: Interpreting the FE coefficients

The fixed-effects estimator removes all time-invariant country characteristics by demeaning the data within each country. Coefficients capture the effect of *within-country changes* in the explanatory variables on *within-country changes* in GDP per capita.

Discussion: If we included time-invariant variables such as `latitude` or `area`, they would be dropped automatically. This is because these variables do not vary over time within a country: once we include country fixed effects (α_i), knowing which country an observation belongs to already determines its latitude or area exactly. In other words, such variables are perfectly collinear with the country dummies and provide no additional within-country variation for estimation.

4. Diagnostics for the FE model.

After estimating the fixed-effects model, perform diagnostic tests to check key assumptions.

```

1 * 4.4 Diagnostic tests after FE estimation
2 * Groupwise heteroskedasticity
3 ssc install xttest3, replace
4 xttest3
5
6 * Serial correlation (Wooldridge test)
7 capture which xtserial
8 if _rc {
9   net from http://www.stata-journal.com/software/sj3-2/
10  net describe st0039
11  net install st0039
12 }
13
14 * Run Wooldridge test for first-order serial correlation
15 xtserial ln_gdppc ln_trade ln_labor ln_capital pop_growth
   inflation school_enroll

```

```

16
17 * Re-run with clustered SEs if issues are detected
18 xtreg ln_gdppc ln_trade ln_labor ln_capital pop_growth
      inflation school_enroll, fe vce(cluster country_id)
19 outreg2 using week10_results.doc, append ctitle("xtreg (
      clustered)") ///
20      keep(ln_trade ln_labor ln_capital pop_growth inflation
      school_enroll)

```

Tip: Why cluster?

If heteroskedasticity or serial correlation is detected, clustered standard errors allow arbitrary correlation of residuals within each country while assuming independence across countries.

5. Adding time fixed effects.

From Task 3, we saw that global GDP per capita trends upward over time, suggesting that time effects may be important. Add year dummies (`i.year`) to capture global shocks or common trends.

$$\ln(\text{gdppc}_{it}) = \alpha_i + \gamma_t + \beta X_{it} + u_{it}$$

Here, X_{it} represents the same set of explanatory variables used in the earlier specifications— $\ln(\text{trade}_{it})$, $\ln(\text{labor}_{it})$, $\ln(\text{capital}_{it})$, pop_growth_{it} , inflation_{it} , $\text{school_enroll}_{it}$. Including the year effects γ_t allows the model to account for macroeconomic shocks or global trends that influence all countries simultaneously.

```

1 * 4.5 Two-way fixed effects: country and year FE
2 xtreg ln_gdppc ln_trade ln_labor ln_capital pop_growth
      inflation school_enroll i.year, fe robust
3 outreg2 using week10_results.doc, append ctitle("xtreg + year
      FE") ///
4      keep(ln_trade ln_labor ln_capital pop_growth inflation
      school_enroll)
5
6 * Test whether year fixed effects are jointly significant
7 testparm i.year

```

Tip: What do year effects capture?

Time fixed effects (γ_t) absorb common shocks affecting all countries in a given year, such as global recessions, commodity price swings, or pandemics. The test `testparm i.year` checks whether these year effects are jointly significant.

6. Optional: Using `reghdfe` for many fixed effects.

For large datasets with many groups or additional dimensions of fixed effects, use the `reghdfe` command (after installing it). It is computationally efficient and produces identical coefficient estimates to `xtreg, fe` for the same specification.

```

1 * Install if needed

```

```
2 ssc install reghdfe, replace
3 ssc install ftools, replace
4
5 * Country fixed effects
6 reghdfe ln_gdppc ln_trade ln_labor ln_capital pop_growth
   inflation school_enroll, ///
7     absorb(country_id) vce(cluster country_id)
8
9 * Country and year fixed effects
10 reghdfe ln_gdppc ln_trade ln_labor ln_capital pop_growth
   inflation school_enroll, ///
11     absorb(country_id year) vce(cluster country_id)
```

Tip: When to use reghdfe

reghdfe is ideal when you need to absorb multiple high-dimensional fixed effects, such as country, industry, and year simultaneously. It is faster and more memory-efficient than the standard xtreg command.

Week 11 — Assignment 5: Endogeneity & Instrumental Variables (IV) (Homework)

11.1 Setting

In many empirical applications, we are interested in the causal effect of a firm's decision or strategic choice on an economic outcome. However, these decisions are often **not randomly assigned**. This leads to an **endogeneity** problem: the explanatory variable of interest is correlated with unobserved factors that also influence the outcome, causing simple OLS regressions to be biased.

A motivating example: advertising and product sales. Consider a panel of retail stores $i = 1, \dots, N$ observed over several years $t = 1, \dots, T$. Each store chooses how much to spend on advertising in each period, and we observe:

- `salesit` – product sales at store i in year t ,
- `ad_spendingit` – advertising expenditures in that year,
- `foot_traffici` – average pedestrian traffic near the store (time-invariant; reflects baseline local demand),
- `budget_shockit` – a firm-wide advertising budget shock (e.g., corporate marketing push; varies by year and is unrelated to individual store characteristics),
- `ad_cost_indexit` – a local media advertising cost index (e.g., fluctuations in radio/TV slot prices; varies over time and across markets).

Our goal is to estimate whether increasing advertising expenditures actually raises sales. A key difficulty in this setting is that stores differ in many unobserved ways. Some stores are located in busy commercial areas, have a loyal customer base, or enjoy strong brand recognition. These features are not recorded in our dataset but clearly affect both advertising decisions and sales.

To make this concrete, suppose the true (but unobserved) relationship between advertising and sales is:

$$\log \text{sales}_{it} = \beta_0 + \beta_1 \log \text{ad_spending}_{it} + \beta_2 \text{popularity}_i + u_{it}, \quad (11)$$

where `popularityi` is a time-invariant store characteristic that captures location quality, demand potential, and other persistent advantages.

Because more popular stores tend to spend more on advertising *and* have higher sales for reasons unrelated to advertising, `log ad_spendingit` is correlated with `popularityi`. If we estimate the simplified OLS regression

$$\log \text{sales}_{it} = \beta_0 + \beta_1 \log \text{ad_spending}_{it} + \varepsilon_{it}, \quad (12)$$

the omitted variable `popularityi` becomes part of the error term, leading to a biased estimate of β_1 . In other words, OLS confuses the effect of advertising with the underlying differences in store popularity.

To address endogeneity, we use **instrumental variables (IVs)**: factors that shift advertising spending but are otherwise unrelated to unobserved determinants of sales. In this lab, you will analyze a simulated panel dataset in which each store faces multiple potential sources of variation in advertising:

- a time-varying *firm-wide advertising budget shock* (`budget_shock`), which raises or lowers advertising resources across all stores in a given year and serves as a plausibly **valid and strong** instrument;
- a time-varying *local advertising cost index* (`ad_cost_index`), which captures fluctuations in media slot prices (e.g., radio, online ads) and acts as a plausibly **valid but weak** instrument;
- a time-invariant measure of pedestrian activity near each store (`foot_traffic`), which affects advertising decisions but is also directly related to local demand and store popularity, making it a **potentially invalid** instrument.

The goals of this assignment are to:

- show how endogeneity leads OLS to produce biased advertising–sales estimates,
- use IV methods to recover the causal effect of advertising under the IV assumptions, and
- compare two candidate instruments based on their **relevance** and **exogeneity**.

11.2 What to Submit

Upload your **final and working** versions of the `.do` file, the corresponding `.log` file, and the exported `lab11_results.doc` to Blackboard. Grading will be based only on the files you submit, so please double-check that they are correct and complete before uploading.

11.3 Assignment Tasks

1. **Set up your default do-file template.** Change your working directory to the location where you saved `ads_iv_full.csv`. Install the following packages if needed:

```
1 ssc install outreg2, replace
```

2. **Load and explore the dataset.** Import the dataset and examine its structure:

```
1 import delimited "ads_iv_full.csv", clear
2 describe
3
4 * Summarize main variables
5 summarize sales ad_spending foot_traffic budget_shock ///
6           ad_cost_index popularity
```

Create log values for `sales` and `ad_spending`. Then, answer the following questions in `lab11_results.doc`:

- How many stores and how many years are in the dataset?
- Which variables vary over time? Which are time-invariant?
- Report the average values of `sales`, `ad_spending`, and `foot_traffic`.

3. **Correct OLS regression (with popularity).** Suppose the magic cat Coco did wonders and, for this particular assignment, you now have access to each store's **popularity**. Estimate the “correct” regression from equation (11):

$$\log \text{sales}_{it} = \beta_0 + \beta_1 \log \text{ad_spending}_{it} + \beta_2 \text{popularity}_i + u_{it}.$$

Run this regression in Stata and export the results to `lab11_results.doc` using `outreg2` with an appropriate column title (e.g., “Correct OLS”).

4. **Naive OLS regression (biased).** Now the whimsical cat Coco vanishes, taking the **popularity** variable with her. From this point on, you no longer observe **popularity**, just as in real empirical work. Estimate the naive log–log OLS model from equation (12):

$$\log \text{sales}_{it} = \beta_0 + \beta_1 \log \text{ad_spending}_{it} + \varepsilon_{it}.$$

Export the results to `lab11_results.doc` (e.g., column title “OLS (Biased)”).

In `lab11_results.doc`, answer:

- What is the estimated effect of advertising on sales according to the **correct** model?
- What is the estimated effect according to the naive model?
- By how much does the naive OLS estimate differ from the correct one when **popularity** is omitted? (Report the difference in estimated elasticities.)

5. **Visual exploration: relationships among variables.** You now decide to use an **IV** approach. Before running any regressions, it is helpful to check visually whether each potential instrument is correlated with advertising.

Produce scatter plots of:

- `ln_ad_spending` vs. `budget_shock` (valid, strong IV),
- `ln_ad_spending` vs. `ad_cost_index` (valid, weak IV),
- `ln_ad_spending` vs. `foot_traffic` (potentially invalid IV).

Export all three plots and include them in `lab11_results.doc`. Briefly describe what each graph suggests about the strength and plausibility of each instrument.

6. **First-stage regressions: instrument relevance.** Now formally test whether each candidate IV is **relevant** by running separate first-stage regressions where the endogenous variable `ln_ad_spending` is regressed on each instrument:

- First stage 1: `ln_ad_spending` on `budget_shock` (valid, strong IV),
- First stage 2: `ln_ad_spending` on `ad_cost_index` (valid, weak IV),
- First stage 3: `ln_ad_spending` on `foot_traffic` (potentially invalid IV).

In `lab11_results.doc`, answer the following:

- Is each instrument **statistically relevant**? In other words, does it significantly predict advertising in the first-stage regression?
- Which instrument is **strongest**? When evaluating instrument strength, look not only at the size of the coefficient but also at the *first-stage F-statistic*. (A larger *F*-statistic means the instrument explains more variation in advertising and is therefore “stronger.”)
- Which instrument appears **weak**? An instrument is considered weak if it has only a small effect on advertising *and* produces a low first-stage *F*-statistic. Explain your reasoning using both pieces of evidence.

- Even if an instrument is strong (for example, `foot_traffic`), does that guarantee it is *valid*? Briefly explain why strength alone does **not** ensure exogeneity.

7. **IV estimation: recover the causal effect of advertising.** We now want to estimate the causal effect of advertising using each instrument separately. You will run three 2SLS regressions for the model:

$$\log \text{sales}_{it} = \beta_0 + \beta_1 \log \text{ad_spending}_{it} + \varepsilon_{it},$$

where `ln_ad_spending` is treated as endogenous.

Tip: Quick guide: `ivregress 2sls`

The basic syntax for a single endogenous regressor and a single instrument is:

```
ivregress 2sls depvar (endogvar = instrument), vce(robust)
```

For example:

```
ivregress 2sls y (x = z), vce(robust)
```

Here:

- `y` is the dependent variable,
- `x` is the endogenous regressor,
- `z` is an instrumental variable that shifts `x` but is otherwise unrelated to the error term.

Important: If an instrument is weak (i.e., has a low first-stage F -statistic), the IV estimate may be very imprecise or unstable. If an instrument is invalid (i.e., correlated with omitted variables), the IV estimate will be biased even if the first stage is strong.

Using this template, run:

- an IV regression using `budget_shock` as the instrument for `ln_ad_spending` (valid, strong IV),
- an IV regression using `ad_cost_index` as the instrument for `ln_ad_spending` (valid, *weak* IV),
- an IV regression using `foot_traffic` as the instrument for `ln_ad_spending` (potentially *invalid* IV).

Export all IV results to `lab11_results.doc` with clear column titles (e.g., “IV: budget_shock”, “IV: ad_cost_index”, and “IV: foot_traffic”).

In your write-up:

- Which IV produces an advertising elasticity closest to the **correct** OLS estimate from equation (11)?
- How does the estimate from the `ad_cost_index` IV illustrate the issues that arise with a **weak** instrument?
- Based on the economic story, why is `foot_traffic` likely invalid, even if it strongly predicts advertising?

(Optional: You may try a fixed-effects regression for comparison. In this particular setting, FE performs remarkably well and provides a useful benchmark!)

12 Command Directory

12.1 Basic

12.1.1 Comments & Line Control

```
1 * single-line comment
2 /* multi-line comment */
3 /// line continuation (use at end of line)
```

12.1.2 Working Directory

- `pwd` : show current working directory.
- `cd "path"` : change working directory.
- `dir` or `list` : list files in current directory.

12.1.3 Using help to check syntax

```
1 help regress
2 help summarize
3 search outreg2
4 viewsource regress.ihlp // view help file source (read-only)
```

12.1.4 Loops (forvalues and foreach)

```
1 // Numeric loops
2 forvalues y = 2018/2024 {
3     display "Year = 'y'"
4     // use file 'y'.dta, clear
5 }
6
7 // Over lists (variables, strings, macros)
8 foreach v in price sqrft bdrms {
9     summarize 'v'
10 }
11 foreach ext in png pdf eps {
12     graph export "myplot.'ext'", replace
13 }
```

12.1.5 Session Control

- `set more off` : prevent “–more–” pauses when running do-files.
- `clear all` : remove data and value labels from memory.
- `cls` : clear the Results window for a fresh output screen.

12.1.6 Log Files

- `capture log close` : safely close any open log (avoids an error if none is open).
- `log using "path\logfile.smcl", replace` : start a new log in SMCL format.
- `log using "path\logfile.log", replace` : start a new plain-text log.
- `log close` : end the current log.

```
1 capture log close
2 log using "output\session_log.smcl", replace
3 // ... your commands ...
4 log close
```

12.2 Data

12.2.1 Load, Browse, Export

- `use filename [, clear]` : load .dta.
- `import delimited filename [, options]` : import CSV/TXT.
- `browse` : open data browser (read-only).
- `save filename [, replace]` : save .dta.
- `export delimited using filename [, options]` : export CSV/TXT.

```
1 use "housing.dta", clear
2 browse
3 save "housing_clean.dta", replace
4 export delimited using "housing_clean.csv", replace
```

12.2.2 Descriptive tools

Variable overview and listing

- `des [varlist]` : describe variables.
- `list varlist [if] [in]` : display rows with conditions/ranges.

```
1 des
2 des price sqrft bdrms
3
4 list price sqrft in 1/10
5 list price sqrft bdrms if bdrms >= 3
```

Summaries and tabulations

- `sum [varlist] [, detail]` : summary stats; `detail` adds percentiles.
- `tab var` : one-way tabulation; `tab var1 var2` two-way.
- `count [if]` : count obs meeting a condition.

```

1 sum price sqrft bdrms
2 sum price, detail
3
4 tab bdrms
5 tab bdrms region
6 count if price > 500000

```

12.2.3 Visual descriptive tools (graphs & export)

Flexible twoway overlays

- General pattern:

```

twoway (plottype1 y1 x1 [, opts1]) (plottype2 y2 x2 [, opts2]) ..., graph_opts
    * Common plottypes: scatter, line, connected, lfit, qfit, kdensity, function.
    * Useful graph_opts: title("..."), xtitle("..."), ytitle("..."), legend(...)
      , name(gname, replace).

```

```

1 // Example A: scatter + fitted line
2 twoway (scatter price sqrft) (lfit price sqrft), ///
3     title("Price vs. Square Feet") xtitle("Sq Ft") ytitle("Price")
4
5 // Example B: overlay two kernel densities
6 twoway (kdensity wage if female==0) (kdensity wage if female==1), ///
7     legend(order(1 "Male" 2 "Female")) title("Wage Distributions")
8
9 // Example C: add a theoretical function
10 twoway (scatter y x) (function y = 2 + .5*x, range(x)), ///
11     title("Data with Theoretical Line")

```

Histograms (with overlays)

- histogram var [, freq | percent | density bin(#) start(#) width(#) ...]
 To *overlay* other plots (e.g., kdensity), use addplot().

```

1 // Basic histogram
2 hist price, bin(30) percent title("Price Distribution")
3
4 // Histogram + kernel density overlay (natural pairing)
5 hist price, bin(30) percent addplot(kdensity price) ///
6     title("Price: Histogram with Kernel Density")
7
8 // Two histograms by group (save and combine)
9 hist price if bdrms==2, percent name(h2, replace) title("Price: 2
10     Bedrooms")
11 hist price if bdrms==3, percent name(h3, replace) title("Price: 3
12     Bedrooms")
13 graph combine h2 h3, cols(2) name(h23, replace)

```

Exporting graphs

- `graph export "filename.png", replace : also .pdf/.eps.`

```
1 graph export "outputs/price_scatter.png", replace
```

12.2.4 Generating and transforming variables

gen and egen Stata provides two main commands to create new variables:

- `gen newvar = expression`

Used for **row-wise transformations**. Each new value is computed directly from variables in the same observation. Example: taking logs, squares, or creating indicator dummies.

- `egen newvar = function(varlist) [, by(group)]`

Used for **group-wise or more complex operations**. It can compute summary statistics (mean, max, count, rank, etc.) within groups, or create unique group identifiers. These are functions that `gen` cannot do directly.

```
1 // Row-wise transforms
2 gen lprice = log(price)
3 gen sqrft2 = sqrft^2
4 gen large = (sqrft >= 2000)
5
6 // Group stats
7 egen mean_price_by_city = mean(price), by(city)
8 egen gid = group(city neighborhood)
```

Difference and lag variables Many analyses require creating *lagged*, *lead*, or *differenced* versions of variables. These transformations allow us to capture dynamic relationships (e.g. how last year's trade affects current growth) or to compute growth rates from levels.

In Stata, you must first declare the data's panel/time structure with `xtset` (for panel data) or `tsset` (for pure time series). After this, you can use Stata's built-in time operators:

- `L.var` : lag of `var` (e.g. `L1.price` = last year's price).

- `F.var` : lead of `var` (e.g. `F1.price` = next year's price).

- `D.var` : difference of `var` (e.g. `D.price` = $price_t - price_{t-1}$).

```
1 // Example: city-by-year panel
2 xtset city year // declare panel: city = panelvar, year = timevar
3
4 gen d_price = D.price // first difference of price
5 gen L1_price = L.price // lag-1 price
6 gen g_price = 100*(D.lprice) // approx % growth (log * 100)
```

12.2.5 Reshaping data

Stata's `reshape` command reorganizes the layout of your dataset:

- `reshape long stub, i(id) j(time)`
Converts from **wide** format (one row per id, many year columns) to **long** format (multiple rows per id, with year stacked in one column).
- `reshape wide stub, i(id) j(time)`
Converts from **long** format (one row per id–time) back to **wide** format (one row per id, variables spread across year columns).

```
1 // Wide -> Long: price2019 price2020 ...
2 reshape long price, i(houseid) j(year)
3
4 // Long -> Wide
5 reshape wide price, i(houseid) j(year)
```

For a more detailed walk-through with World Bank panel data, see the Week 4 teaching materials.

12.2.6 Merging data

Stata's `merge` command combines datasets by matching on key variables:

- `merge 1:1 keyvars using file.dta`
One observation in master matches one observation in using (1-to-1).
- `merge m:1 keyvars using file.dta`
Many observations in master match one in using (many-to-1).
- `merge 1:m keyvars using file.dta`
One observation in master matches many in using (1-to-many).

```
1 use "houses.dta", clear
2 merge m:1 zipcode using "zip_income.dta"
3
4 // Inspect merge result via _merge
5 tab _merge
6 drop if _merge==2 // drop using-only if appropriate
7 drop _merge
```

For more on interpreting merge diagnostics (such as `_merge` codes), see the Week 4 materials.

12.3 Regression

12.3.1 Basic regress syntax

- `regress depvar indepvars [if] [in] [weight], [options]`

```
1 reg price sqrft bdrms
2 reg lprice sqrft bdrms i.city, robust
```

12.3.2 Accessing results & robust SE

- `_b[var]` : coefficient; `_se[var]` : standard error.
- `e(r2)`, `e(N)` : stored results.
- `predict newvar, xb (fitted)` `predict newvar, resid (residuals)`.
- `, robust` : heteroskedasticity-robust SE.

```

1 reg lprice sqrft bdrms, robust
2 di "Beta_sqrft = " _b[sqrft]
3 di "R2 = " e(r2)
4 predict yhat, xb
5 predict uhat, resid

```

12.3.3 Factor variables: ##, interactions, and margins

Stata's **factor variable** notation (`i.`, `c.`, `#`, `##`) lets you specify interactions and polynomial terms directly inside regression commands. This avoids the need to manually generate new variables and ensures that Stata correctly handles dummies, interactions, and post-estimation tools such as `margins`.

- `c.x1##c.x2` : interaction of two continuous variables (includes x_1 , x_2 , and $x_1 \times x_2$).
- `i.group##c.x` : interaction of a discrete factor and a continuous variable (e.g. fixed effects \times slope).
- `margins` : compute predicted values or contrasts at specific values of regressors.

```

1 // Interaction of square feet and bedrooms
2 reg lprice c.sqrft##c.bdrms, robust
3
4 // City fixed effects with interaction
5 reg lprice c.sqrft##i.city bdrms, robust
6
7 // Marginal effects at specified values
8 margins, at(sqrft=(1000(500)3000) bdrms=(2 3 4)) nose
9 marginsplot, xdimension(sqrft)

```

Tip: Why use ##?

- If you manually create $x_1 \times x_2$ with `gen`, Stata will not automatically include the main effects or treat categorical dummies correctly.
- Using `##` tells Stata to include the main effects *and* the interaction, and it also integrates seamlessly with `margins` and `marginsplot`.

Tip: Getting coefficient labels with factor variables

When you use ##, Stata generates expanded variable names internally (e.g. 1.city#c.sqrft). To see the exact names you can use in `lincom` or `test`, run:

```
1 reg lprice c.sqrft##i.city, robust
2 reg ..., coeflegend
```

The option `coeflegend` displays the internal coefficient labels you can copy directly into post-estimation commands.

12.3.4 Inference : test and lincom

Stata's `test` and `lincom` commands allow you to form and test linear combinations of coefficients.

- `test varlist` Wald test of linear or joint hypotheses.

```
1 // Example regression with interaction
2 reg lprice c.sqrft##c.bdrms, robust
3
4 // ---- Using test ----
5 // (1) Does the marginal effect of sqrft depend on bdrms?
6 // H0: beta3 = 0
7 test c.sqrft#c.bdrms
8
9 // (2) Does sqrft have no effect on lprice at all?
10 // H0: beta1 = 0 and beta3 = 0
11 test c.sqrft c.sqrft#c.bdrms
```

- `lincom expression` Estimate a specific linear combination of coefficients.

```
1 // ---- Using lincom ----
2 // (3) Effect of sqrft when bdrms = 3
3 // H0: beta1 + 3*beta3 = 0
4 lincom c.sqrft + 3*(c.sqrft#c.bdrms)
5
6 // (4) Difference in sqrft effect between 2 vs 4 bedrooms
7 // H0: (beta1 + 4*beta3) - (beta1 + 2*beta3) = 0
8 // i.e., the marginal effect of sqrft is the same at 2 vs 4
9 // bedrooms
10 lincom (c.sqrft + 4*(c.sqrft#c.bdrms)) - (c.sqrft + 2*(c.sqrft#c.
11 bdrms))
```

Tip: Interpreting interaction effects

- `test` is useful for joint significance of coefficients: e.g. $\beta_3 = 0$ (interaction irrelevant) or $\beta_1 = \beta_3 = 0$ (no effect at all).
- `lincom` is best for computing the effect of x_1 at a particular value of x_2 , or differences across values.
- Always check coefficient labels using `reg ...`, `coeflegend` so you know the correct names to use in `lincom`.

12.3.5 Exporting using outreg2

Most exports use user-written commands. In this class we standardize on `outreg2` for Word/Excel.

- `outreg2 using filename, [options]` : export summary or regression tables.

```

1  * One-time installation
2  capture which outreg2
3  if _rc ssc install outreg2
4
5  * [A] Summary statistics to Word
6  outreg2 using "output\lab3_summary.doc", replace ///
7      sum(log) keep(price sqrft bdrms) dec(3)
8
9  * [B] Regressions to Word (append models)
10 reg lprice sqrft bdrms
11 outreg2 using "output\reg_results.doc", replace ///
12     ctitle("OLS: lprice on sqrft, bdrms") dec(3) bdec(3) tdec(3) ///
13     addstat("R-squared", e(r2))
14
15 reg lprice c.sqrft#c.bdrms i.city, robust
16 outreg2 using "output\reg_results.doc", append ///
17     ctitle("OLS: + interactions + city FE, robust") dec(3) bdec(3)
18     tdec(3)

```