

# Problem Set 4

## Nonstationarity, Causality, Instrumental Variables

**EC 421:** Introduction to Econometrics

Due *before* midnight (11:59pm) on Wednesday, 05 June 2019

DUE Your solutions to this problem set are due *before* midnight on Wednesday, 05 June 2019. Your files must be uploaded to [Canvas](#).

IMPORTANT Your submission must include (1) **your responses/answers to the question in a PDF, Word, or similar file** and (2) the R script you used to generate your answers. **The R script is just for your code. To receive credit, your answers/figures/etc. must be in the PDF/Word document.** Each student must turn in her/his own answers.

OBJECTIVE This problem set has three purposes: (1) reinforce econometrics topics from class; (2) build your R toolset; (3) strengthen your intuition on causality and time series.

## Problem 1: Nonstationarity—the Basics

**1a.** Define stationarity.

Note: You can define it using math or words (or both).

**1b.** If our disturbance term  $u_t$  follows a [random walk](#), i.e.,

$$u_t = u_{t-1} + \varepsilon_t$$

then its variance is  $\text{Var}(u_t) = t\sigma_\varepsilon^2$ . Explain how this expression of its variance shows that the disturbance is [nonstationary](#) (i.e., it violates [stationarity](#)).

**1c.** We previously discussed autocorrelated disturbances, e.g., an AR(1) process such that

$$u_t = \rho u_{t-1} + \varepsilon_t$$

Under which circumstances would this AR(1) process become a random walk?

*Hint:* Consider the values of  $\rho$ .

## Problem 2: Nonstationarity—the Simulation

In this problem, we are going to create two independent, **nonstationary** time series. Specifically, we'll create two random walks. Then, we'll regress the first random walk on the second random walk.

*Hint:* Generating random walks is *nearly* identical to generating AR(1) processes, as you did in lab.

**2a.** Generate the first 50-period random walk. We will name it `v`.

$$v_t = v_{t-1} + \varepsilon_t$$

where  $\varepsilon_t$  comes from a normal distribution with mean 0 and standard deviation 1.

Here is some `R` to help.

```
# Set a seed (so your results stay the same)
set.seed(1234)
# Generate the initial number, (this will be v[1])
v <- rnorm(1, mean = 0, sd = 1)
# For loop to create the random walk
for (t in 2:50) {
  # Create the 'next' observation
  ...
}
```

while you're filling in the `for` loop, keep in mind **(1)** our equation for the random walk at the beginning of this question (meaning  $v_t$  depends upon  $v_{t-1}$  and  $\varepsilon_t$ ) and **(2)** the fact that you can reference different observations in `R`, e.g.,

- `v[t]` refers to the  $t^{\text{th}}$  observation
- `v[t-1]` refers to the  $(t - 1)^{\text{th}}$  observation
- `v[3]` refers to the  $3^{\text{rd}}$  observation

If you need more help on for loops, don't forget there are lab materials on Canvas and resources online (e.g., [datamentor.io](http://datamentor.io) and [datacamp.com](http://datacamp.com) have lots of resources).

**2b.** Generate a second 50-period random walk called `w`. This part is exactly the same as (2a), but you **use a different seed** (i.e., `set.seed(456)`) and **name the variable** `w`.

**2c.** We **independently** generated these two time series. Ideally (from a statistical point of view), should we find a statistically significant relationship between the two series? Explain.

**2d.** Regress `w` on `v`. Report the results from the  $t$  test. Do they match your expectations from (2c)? Explain.

**2e.** As we've mentioned, one (simple) way you can work with the nonstationary from random walks is to take differences, i.e.,  $v_t - v_{t-1}$ . The interpretation of the relationship does not change, whether we regress `w` on `v` or  $\Delta w$  on  $\Delta v$  (where  $\Delta w = w_t - w_{t-1}$ ). In `R`, you can use the `diff()` function to take difference, i.e., `diff(v)` will calculate the differences for the variable `v`.

Regress the differenced `w` on the differenced `v`. Does it change your results from **2d**?

## Problem 3: Causality

Following the Rubin causal model, imagine that we observe the following data (which would be impossible observe in real life):

**Table: Imaginary dataset**

$i$	Trt.	$y_1$	$y_0$
1	0	25	17
2	0	15	11
3	1	11	3
4	1	13	9

**3a.** Calculate the treatment effect **for each individual** (i.e.,  $\tau_i$ ).

**3b. [T/F]** The treatment effect is constant across individuals.

**3c.** Calculate the **average treatment effect**.

**3d Estimate the average treatment effect** by comparing the **mean of the treatment group** to the **mean of the control group**.

**3e.** Should we expect our estimator in (3d) to provide unbiased estimates? **Explain.**

**3f.** Why would it be impossible to actually observe all of the data in the table (in real life)?

**3g.** How does your answer in (3f) relate to *the fundamental problem of causal inference*?

## Problem 4: Instrumental Variables

**4a.** What are the two requirements for a valid instrument?

We're interested in estimating  $\beta_1$  in

$$\text{Wage}_i = \beta_0 + \beta_1 \text{Education}_i + u_i$$

but we have a problem with omitted-variable bias. Instrumental variables can potentially help.

**4b.** As we've discussed, we need an instrument for (endogenous) education. Do you think the number of children would be a valid instrument? Explain why it passes/fails each of the two requirements for a valid instrument.

**4c.** Which estimates would you trust more—OLS or IV, where number-of-children is your instrument? Explain.