# Tutorials
## *Week 3*

# Regression Analysis with time series data II

| Pdf file on Blackboard | Dataset on Blackboard | Papers related to the data | Description |
|---|---|---|---|
| C 18.2 | hseinv.dta | McFadden, D., 1994. Demographics, the housing market, and the welfare of the elderly. In: D.A. Wise, ed. 1994. Studies in the Economics of Aging. Chicago: University of Chicago Press, pp.225-285. | Use of lagged variables, test for unit root (Dickey Fuller test), use of ADF (augmented DF) test, consequences of unit root, meaning and consequences of co-integration. |
| C 18.3 | volat.dta | Hamilton, J. D., & Lin, G. (1996). Stock Market Volatility and the Business Cycle. Journal of Applied Econometrics, 11(5), 573–593. http://www.jstor.org/stable/2285217 | AR(3) model, Granger causality, VAR, Wald Test, logit, forecast, MAE. |
| C.18.3 | fertil3.dta | Whittington, L. A., Alm, J., & Peters, H. E. (1990). Fertility and the Personal Exemption: Implicit Pronatalist Policy in the United States. *The American Economic Review*, *80*(3), 545–556. http://www.jstor.org/stable/2006683 | Random walk with drift, AR(2), forecasting, MAE |

**C.18.2 Use the data HSEINV.RAW for this exercise.**

**Demographics, the Housing Market, and the Welfare of the Elderly**

**Variables:**
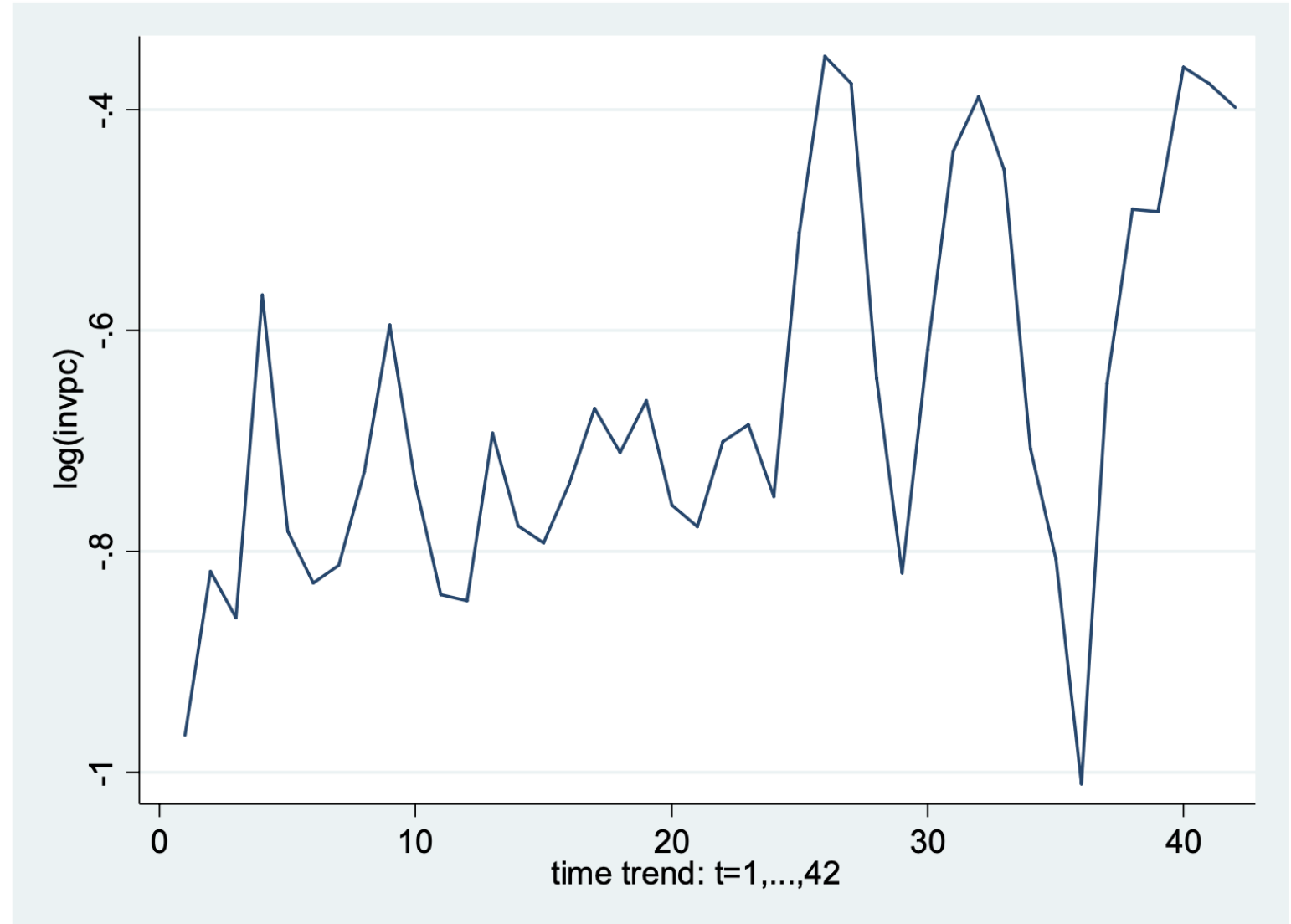invpc: per capita investment: inv/pop
price: housing price index; 1982=1
year:1947-1988
t: time trend: t=1,...,42

Before we check if log(*invpc)* has a
unit root, plot log(invpc) and time.
Discuss the graph.

. line linvpc t

**C.18.2 Use the data HSEINV.RAW for this exercise.**

**i) Test for unit root in log(invpc), including a linear time trend and two lags of $\Delta \log(invpc)$.  Use a 1% significance level.**

- Perform the augmented Dickey-Fuller test for unit root separately for all variables of the regression equation.

$$\Delta\log(invpc_{t)} = \alpha_0 + \alpha_1 \log(invpc_{t-1}) + \alpha_2\Delta\log(invpc_{t-1)} + \alpha_3\Delta\log(invpc_{t-2)} + \alpha_4 t + \varepsilon_t$$

- Where: $\Delta\log(invpc_{t)} = \log(invpc_t - \log(invpc_{t-1})$

- In Stata: declare the dataset as time-series:

```
. tsset t
        time variable:  t, 1 to 42
                delta:  1 unit
```

$$\Delta \log(invpc_t) = \alpha_0 + \alpha_1 \log(invpc_{t-1}) + \alpha_2 \Delta \log(invpc_{t-1}) + \alpha_3 \Delta \log(invpc_{t-2}) + \alpha_4 t + e_t$$

Then estimate the test regression.

```
. reg d.linvpc l.linvpc dl(1/2).linvpc t
```

| Source | SS | df | MS | | Number of obs = | 39 |
|---|---|---|---|---|---|---|
| | | | | | F( 4, 34) = | 6.59 |
| Model | .34943844 | 4 | .08735961 | | Prob > F = | 0.0005 |
| Residual | .450566441 | 34 | .013251954 | | R-squared = | 0.4368 |
| | | | | | Adj R-squared = | 0.3705 |
| Total | .800004881 | 38 | .02105276 | | Root MSE = | .11512 |

| D.linvpc | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| linvpc | | | | | | |
| L1. | -.9557867 | .1977787 | -4.83 | 0.000 | -1.357721 | -.5538521 |
| LD. | .531659 | .1615547 | 3.29 | 0.002 | .2033404 | .8599776 |
| L2D. | .2900152 | .1646455 | 1.76 | 0.087 | -.0445847 | .6246151 |
| | | | | | | |
| t | .00676 | .0021276 | 3.18 | 0.003 | .0024361 | .0110839 |
| _cons | -.7863716 | .1699981 | -4.63 | 0.000 | -1.131849 | -.4408939 |

- The test statistics of interest is the t-statistics of the lagged linvpc , $\alpha_1 \log(invpc_{t-1})$ which is -4.83.

- We will test:

$$H_0 : \theta = 0 \leftrightarrow \rho = 1 \text{ unit root (non stationary)}$$

$$H_A : \theta < 0 \leftrightarrow \rho < 1 \text{ no unit root (stationary)}$$

- If t-value < t critical value, → reject Ho                    Check table 18.3

- -4,83 < -3,96 → we reject Ho at a 1% significance level. That means, the variable log(invpc) does not have a unit root. It is trend stationary.

```
. dfuller linvpc, lag(2) trend reg

Augmented Dickey-Fuller test for unit root          Number of obs   =        39

                        ---------- Interpolated Dickey-Fuller ---------
              Test          1% Critical       5% Critical      10% Critical
           Statistic           Value             Value             Value
------------------------------------------------------------------------------
 Z(t)         -4.833           -4.251            -3.544            -3.206
------------------------------------------------------------------------------
MacKinnon approximate p-value for Z(t) = 0.0004

------------------------------------------------------------------------------
D.linvpc    |     Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
------------+-----------------------------------------------------------------
    linvpc  |
       L1.  |  -.9557867   .1977787    -4.83   0.000    -1.357721   -.5538521
       LD.  |   .531659    .1615547     3.29   0.002     .2033404    .8599776
      L2D.  |   .2900152   .1646455     1.76   0.087    -.0445847    .6246151
    _trend  |   .00676     .0021276     3.18   0.003     .0024361    .0110839
     _cons  |  -.7796116   .1683262    -4.63   0.000    -1.121692   -.4375316
------------------------------------------------------------------------------
```

- The only difference you can observe is that Stata will not use the same critical values.

- This test is called an Augmented DF test, as lags have been added, requiring different critical values. Use these critical values instead of the ones from Table 18.3.

- Test:

$$H_0 : \theta = 0 \leftrightarrow \rho = 1 \text{ unit root (non stationary)}$$

$$H_A : \theta < 0 \leftrightarrow \rho < 1 \text{ no unit root (stationary)}$$

- If t value < t critical value → reject Ho

- -4,83 < -4,251 → reject Ho at 1% sig. level, which means the variable log(invpc) does not have a unit root. And it is stationary.

Utrecht
University

**Variables:**

invpc: per capita investment:
inv/pop

price: housing price index;
1982=1

year:1947-1988

t: time trend: t=1,...,42

```
. dfuller lprice , lag(2) trend reg

Augmented Dickey-Fuller test for unit root          Number of obs    =          39

                                    ---------- Interpolated Dickey-Fuller ---------
                    Test           1% Critical          5% Critical         10% Critical
                  Statistic           Value                Value                Value
------------------------------------------------------------------------------------
 Z(t)              -2.409             -4.251               -3.544               -3.206
------------------------------------------------------------------------------------
MacKinnon approximate p-value for Z(t) = 0.3749


------------------------------------------------------------------------------------
D.lprice     |      Coef.    Std. Err.       t      P>|t|      [95% Conf. Interval]
-------------+----------------------------------------------------------------------
lprice       |
         L1. |  -.2216337    .092006      -2.41    0.022     -.4086124    -.0346549
         LD. |   .327572     .1551807      2.11    0.042      .012207      .642937
        L2D. |   .1300876    .1491206      0.87    0.389     -.172962     .4331372
      _trend |   .000971     .0004867      1.99    0.054     -.0000182    .0019602
       _cons |  -.039384     .0190149     -2.07    0.046     -.0780269    -.0007412
------------------------------------------------------------------------------------
```

- Test:

$$H_0 : \theta = 0 \leftrightarrow \rho = 1 \text{ unit root (non stationary)}$$

$$H_A : \theta < 0 \leftrightarrow \rho < 1 \text{ no unit root (stationary)}$$

- If t value < t critical value → reject Ho

- -2,41 > - 4,251 → we can not reject Ho. The variable log(price) has a unit root. It follows a non-stationary process, so not weakly dependent.

**iii) Given the outcomes in parts i) and ii), does it make sense to test for cointegration between log(invpc) and log(price)?**

- If one variable is non-stationary – unit root- (I(1)), and the other is stationary – no unit root - (I(0)), then it does not make sense to co-integrate them.
- Both variables need to have a unit root (be non-stationary) to be co-integrated.

- The answer is no. Cointegration makes sense between two non-stationary processes (unit roots) integrated in the same order.
- If we take any nontrivial linear combination of an I(0) process (which may have a trend) and an I(1) process, the result will be an I(1) process (possibly with drift).

**C.18.3 Use the data in VOLAT.RAW for this exercise.**
**Graph pcip against time. Does it contain a clear upward or downward trend over the entire sample period?**

gen time = _n
tsset time

We plot PCIP against time, in this case, date.

line pcip date ||lfit pcip date

Pcip: annualized percentage change in the Industrial Production index.
date: 1947.01 to 1993.06



- The graph shows significant fluctuations in the annual percentage change of industrial production over the years.
- There is plenty of volatility
- Clear upward and downward trends over time.
- However, the trend appears to stabilize around the zero line, especially in the later years. That means that the PCIP (annual % change in industrial production) does not show a strong upward or downward trend in the long term.

**C.18.3 Use the data in VOLAT.RAW for this exercise.**
**Check for unit root in pcip** : annualized percentage change in the Industrial Production index.

```
. dfuller pcip, lags(3) trend

Augmented Dickey-Fuller test for unit root

Variable: pcip                          Number of obs  = 553
                                        Number of lags =   3

H0: Random walk with or without drift

                                    Dickey-Fuller
                 Test          ──────── critical value ────────
              statistic         1%           5%           10%
          ─────────────────────────────────────────────────────
Z(t)          -9.053         -3.960       -3.410       -3.120
          ─────────────────────────────────────────────────────
MacKinnon approximate p-value for Z(t) = 0.0000.
```

- Test:

$$H_0 : \theta = 0 \leftrightarrow \rho = 1 \text{ unit root (non stationary)}$$

$$H_A : \theta < 0 \leftrightarrow \rho < 1 \text{ no unit root (stationary)}$$

- If |t value| < t critical value → reject Ho
- -9.053 < -3.960 → reject Ho at 1% sig. level (and at all other levels). pcip does not have a unit root. It is stationary.

**C.18.3 Use the data in VOLAT.RAW for this exercise.**
**i) Estimate an AR(3) model for pcip. Now, add a fourth lag and verify that it is very insignificant.**

Frist generate a time trend:

$$pcip_t = \beta_0 + \beta_1(pcip_{t-1}) + \beta_2(pcip_{t-2}) + \beta_3(pcip_{t-3})$$

. gen time = _n

Then declare the data as time-series:

. tsset time
       time variable:  time, 1 to 558
               delta:  1 unit

We now estimate the AR(3) specification:
. reg pcip l(1/3).pcip

|      Source |         SS |    df |         MS |
|------------:|-----------:|------:|-----------:|
|       Model | 16126.3579 |     3 | 5375.45264 |
|    Residual | 81224.8954 |   550 | 147.681628 |
|       Total | 97351.2533 |   553 | 176.042049 |

Number of obs =     554
F( 3,    550) =   36.40
Prob > F      =  0.0000
R-squared     =  0.1657
Adj R-squared =  0.1611
Root MSE      =  12.152

|    pcip |      Coef. |  Std. Err. |    t |  P>|t| | [95% Conf. | Interval] |
|--------:|-----------:|-----------:|-----:|-------:|-----------:|----------:|
| pcip    |            |            |      |        |            |           |
| L1. |  .3491232 |  .0425232 | 8.21 | 0.000 |  .2655954 | .4326509 |
| L2. |  .0707984 |  .0449501 | 1.58 | 0.116 | -.0174965 | .1590932 |
| L3. |  .0673713 |  .0425274 | 1.58 | 0.114 | -.0161647 | .1509073 |
|     |           |           |      |       |           |          |
| _cons |  1.804189 |  .5480442 | 3.29 | 0.001 |  .7276729 | 2.880704 |

# We test for heteroskedasticity and autocorrelation in the error term

## Test for heteroskedasticity (Breush-Pagan Test)

```
. hettest

Breusch-Pagan/Cook-Weisberg test for heteroskedasticity
Assumption: Normal error terms
Variable: Fitted values of pcip


H0: Constant variance

    chi2(1) =   62.06
Prob > chi2 = 0.0000
```

The error term is heteroscedastic at 1%. For the exam, you have to write all the statistical steps.

If chi2-stat > CVchi2, then reject Ho. 62.06 > 6.63, reject Ho. There is Heteroskedasticity in the error term at 1%. *Table G.4.*

## Test for serial correlation (Breusch Godfrey Test)

```
. bgodfrey, lags(1/3)

Breusch-Godfrey LM test for autocorrelation
---------------------------------------------------------------------------
   lags(p) |          chi2               df              Prob > chi2
-----------+---------------------------------------------------------------
     1     |          0.007              1                 0.9357
     2     |          0.890              2                 0.6408
     3     |          1.312              3                 0.7263
---------------------------------------------------------------------------
                 H0: no serial correlation
```
But we find no evidence for serial correlation in the error-term.

Ho: $\rho = 0$ (no 1st order serial correlation)
$H_1$: $\rho \neq 0$ (serial correlation)

If pvalue < 0.10 or < 0.05; then reject Ho. P values are not less than 0.10 or 0.05. Therefore, we can not reject the Ho. So, there is no serial correlation.

**We use the rob option**

```
. reg pcip l(1/3).pcip, rob

Linear regression                              Number of obs =       554
                                               F(  3,    550) =     16.24
                                               Prob > F       =    0.0000
                                               R-squared      =    0.1657
                                               Root MSE       =    12.152

------------------------------------------------------------------------------
             |              Robust
       pcip |     Coef.    Std. Err.      t     P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
       pcip |
        L1.  |   .3491232   .0623237     5.60   0.000     .2267015    .4715448
        L2.  |   .0707984   .0488291     1.45   0.148    -.025116     .1667127
        L3.  |   .0673713   .0414654     1.62   0.105    -.0140786    .1488212
             |
       _cons |   1.804189   .6400428     2.82   0.005     .5469612    3.061416
------------------------------------------------------------------------------
```

**Now, add a fourth lag and verify that it is very insignificant. And also use the rob option**

```
. reg pcip l(1/4).pcip, rob

Linear regression                              Number of obs =       553
                                               F(  4,    548) =     12.21
                                               Prob > F       =    0.0000
                                               R-squared      =    0.1659
                                               Root MSE       =    12.173

------------------------------------------------------------------------------
             |               Robust
       pcip  |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
       pcip  |
        L1.  |    .349382    .0632797     5.52   0.000     .2250816    .4736823
        L2.  |   .0702363    .0483599     1.45   0.147    -.0247571    .1652298
        L3.  |   .0657502    .0443265     1.48   0.139    -.0213205    .1528209
        L4.  |   .0043168    .0587231     0.07   0.941    -.1110331    .1196667
             |
       _cons |   1.787332    .6639847     2.69   0.007     .4830655    3.091599
------------------------------------------------------------------------------
```

- When $pcip_{t-4}$ is added, its coefficient is 0.0043 with a t-statistics of about 0.10.
- The t-statistics of the fourth lag is very small (0.07); the coefficient is close to zero.
- We conclude that the fourth lag is insignificant (Pvalue is greater than 0.10).

**ii) To the AR(3) model from part i), add three lags of pcsp to test whether pcsp Granger causes pcip. Carefully, state your conclusion.**

The Granger test specification:

$$pcip_t = \delta_0 + \alpha_1 pcip_{t-1} + \alpha_2 pcip_{t-2} + \alpha_3 pcip_{t-3} + \gamma_1 pcsp_{t-1} + \gamma_2 pcsp_{t-2} + \gamma_3 pcsp_{t-3} + u_t,$$

Variables:    **pcsp:** = %change, sp500, ann rate.        **pcip:** annualized %change in the industrial production index.

Test the Granger-causality in the following way:

If *pcsp* Granger causes *pcip* then, its past values should explain the current value of *pcip* in a statistically significant way.

$H_0$: $\gamma_1 = \gamma_2 = \gamma_3 = 0$     *pcsp* does not have Granger causality on *pcip*.

$H_1$:    At least one of the tested coefficients is not zero. In this case, *pcsp* Granger causes *pcip*.

Let us estimate the equation (since heteroskedasticity exists in the error term, we use robust estimation).

```
. reg pcip l(1/3).pcip l(1/3).pcsp, rob
```

Linear regression

```
                                              Number of obs =       554
                                              F(  6,    547) =     14.13
                                              Prob > F       =    0.0000
                                              R-squared      =    0.1895
                                              Root MSE       =     12.01
```

| pcip | Coef. | Robust Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
|---|---|---|---|---|---|---|
| pcip | | | | | | |
| L1. | .3258447 | .0621686 | 5.24 | 0.000 | .2037263 | .4479631 |
| L2. | .0691003 | .0475705 | 1.45 | 0.147 | -.0243429 | .1625434 |
| L3. | .0799492 | .0410298 | 1.95 | 0.052 | -.0006461 | .1605444 |
| | | | | | | |
| pcsp | | | | | | |
| L1. | .0234479 | .0134366 | 1.75 | 0.082 | -.0029458 | .0498416 |
| L2. | .0323316 | .012825 | 2.52 | 0.012 | .0071394 | .0575238 |
| L3. | .0195941 | .0137932 | 1.42 | 0.156 | -.0075 | .0466883 |
| | | | | | | |
| _cons | 1.245541 | .6232294 | 2.00 | 0.046 | .0213254 | 2.469757 |

We now test the exclusion restrictions for the lags of the *pcsp*:

```
. test 1.pcsp 12.pcsp 13.pcsp

( 1)   L.pcsp = 0
( 2)   L2.pcsp = 0
( 3)   L3.pcsp = 0

       F(  3,    547) =     5.71
             Prob > F =     0.0007
```

- $H_o$: $\beta_4 = \beta_5 = \beta_6 = 0$   *pcsp* does not have Granger causality on *pcip*.

- $H_1$   At least one of the tested coefficients is not zero. In this case, *pcsp* Granger causes *pcip*.

- $If\ Fstat > Fcv, reject\ Ho$

- 5.71 > 3.78, reject Ho, at 1% significance level.                                        Table G.3.c

- Pcsp (pct chg, sp500, ann rate) does Granger cause pcip (pct chg, IP, ann rate). That means that past values of the change of the Standard and Poor's 500 index (pcsp) can predict changes in the current value of industrial production growth rate (pcip).

**iii) To the model in part ii), add three lags of the change in i3, the three-month T-bill rate. Does pcsp Granger cause pcip conditional on past $\Delta i3$?**

Model from part ii)

$$pcip_t = \delta_0 + \alpha_1 pcip_{t-1} + \alpha_2 pcip_{t-2} + \alpha_3 pcip_{t-3} + \gamma_1 pcsp_{t-1} + \gamma_2 pcsp_{t-2} + \gamma_3 pcsp_{t-3} + u_t,$$

**Solution:**

- i3: 3 mo. T-bill annualized rate

- Difference of the interest rate for 3 months = $\Delta i_3$

- 3 lags of the change of i3: $\Delta i_{t-1} + \Delta i_{t-2} + \Delta i_{t-3}$

- Granger Causality: if the difference of the interest rate on the past 3 months will still be considered on the Granger causality of pCSP on pCIP

- We need to include three lags of the change of i3 (3 month Treasury Bill annualized interest rate) and retest Granger causality from pcsp to pcip. That means, test joint significance of $pcsp_{t-1}; pcsp_{t-2}; pcsp_{t-3}$.

. reg pcip l(1/3).pcip l(1/3).pcsp dl(1/3).i3, rob

Linear regression

```
Number of obs =      554
F( 9,    544) =    11.93
Prob > F      =   0.0000
R-squared     =   0.1959
Root MSE      =   11.995
```

| | | Robust | | | | |
|---|---|---|---|---|---|---|
| pcip | Coef. | Std. Err. | t | P>\|t\| | [95% Conf. Interval] | |
| pcip | | | | | | |
| L1. | .3145074 | .0649439 | 4.84 | 0.000 | .1869358 | .442079 |
| L2. | .0621721 | .0487805 | 1.27 | 0.203 | -.0336492 | .1579934 |
| L3. | .0789091 | .0411188 | 1.92 | 0.056 | -.001862 | .1596801 |
| pcsp | | | | | | |
| L1. | .028815 | .0137396 | 2.10 | 0.036 | .0018258 | .0558042 |
| L2. | .0314511 | .0128079 | 2.46 | 0.014 | .0062921 | .0566102 |
| L3. | .0141627 | .0139711 | 1.01 | 0.311 | -.0132812 | .0416067 |
| i3 | | | | | | |
| LD. | 1.519901 | 1.237951 | 1.23 | 0.220 | -.9118499 | 3.951651 |
| L2D. | 1.268064 | 1.151717 | 1.10 | 0.271 | -.9942935 | 3.530422 |
| L3D. | -.7773987 | 1.14427 | -0.68 | 0.497 | -3.025127 | 1.470329 |
| _cons | 1.311842 | .6380201 | 2.06 | 0.040 | .0585574 | 2.565127 |

. test l.pcsp l2.pcsp l3.pcsp

```
( 1)  L.pcsp = 0
( 2)  L2.pcsp = 0
( 3)  L3.pcsp = 0

    F( 3,   544) =     5.18
         Prob > F =    0.0016
```

**Solution:**

This was added: $\Delta i_{t-1} + \Delta i_{t-2} + \Delta i_{t-3}$

Retest Granger causality from pcsp to pcip. That means, test joint significance of $pcsp_{t-1}; pcsp_{t-2}; pcsp_{t-3}.$

$H_o: \gamma_1 = \gamma_2 = \gamma_3 = 0$

$H_1:$ At least one of the tested coefficients is not zero. In this case, *pcsp* Granger causes *pcip*.

$If\ Fstat > Fcv, reject\ Ho$

5.18 > 3.78, reject Ho at level 1%

Conclusion: Past $\Delta i3$ is considered on the Granger causality of pCSP on pCIP.

## Additional material

It is customary that Granger tests are carried out in all directions. That is, in this particular example, we should not only test id pcsp Granger causes pcip, but also vice versa. We can do this either in a single equation framework, just as we did before, or we can use a VAR (Vector Autoregression) model instead.

If you use the following command:

```
. var pcip pcsp d.i3, lag(1/3) small

Vector autoregression

Sample:  5 - 558                                   No. of obs      =         554
Log likelihood = -5298.088                         AIC             =    19.23498
FPE            =    45313.31                        HQIC            =     19.3263
Det(Sigma_ml)  =    40661.66                        SBIC            =    19.46876


Equation           Parms       RMSE      R-sq          F          P > F
-----------------------------------------------------------------------------
pcip                  10     11.9954    0.1959    15.00105    0.0000
pcsp                  10     38.5903    0.0966     6.579677    0.0000
D_i3                  10     .455928    0.1582    11.57126    0.0000
-----------------------------------------------------------------------------
```

```
           --------------------------------------------------------------------
                    |      Coef.    Std. Err.       t     P>|t|    [95% Conf. Interval]
           ---------+----------------------------------------------------------
           pcip     |
              pcip  |
               L1.  |    .3145074    .0427499     7.36    0.000    .2305323     .3984825
               L2.  |    .0621721    .0444846     1.40    0.163   -.0252105     .1495547
               L3.  |    .0789091    .0423072     1.87    0.063   -.0041963     .1620145
                    |
              pcsp  |
               L1.  |     .028815     .013294     2.17    0.031    .0027011     .0549289
               L2.  |    .0314511    .0136768     2.30    0.022    .0045853      .058317
               L3.  |    .0141627      .01336     1.06    0.290   -.0120808     .0404063
                    |
                i3  |
               LD.  |    1.519901    1.135269     1.34    0.181   -.7101465     3.749948
              L2D.  |    1.268064    1.160591     1.09    0.275   -1.011725     3.547854
              L3D.  |   -.7773987     1.13927    -0.68    0.495   -3.015306     1.460509
                    |
             _cons  |    1.311842    .5563367     2.36    0.019    .2190109     2.404673
           ---------+----------------------------------------------------------
           pcsp     |
              pcip  |
               L1.  |     .021078    .1375305     0.15    0.878   -.2490778     .2912338
               L2.  |   -.0968251    .1431111    -0.68    0.499   -.3779432      .184293
               L3.  |   -.0490966    .1361062    -0.36    0.718   -.3164547     .2182615
                    |
              pcsp  |
               L1.  |    .2345174    .0427681     5.48    0.000    .1505065     .3185284
               L2.  |   -.0691706    .0439997    -1.57    0.117   -.1556007     .0172595
               L3.  |    .0608425    .0429805     1.42    0.157   -.0235856     .1452705
                    |
                i3  |
               LD.  |   -14.60858    3.652268    -4.00    0.000   -21.78285    -7.434302
              L2D.  |    1.383237    3.733734     0.37    0.711   -5.951065     8.717538
              L3D.  |    .0436589    3.665141     0.01    0.991   -7.155904     7.243222
                    |
             _cons  |    6.796643    1.789789     3.80    0.000      3.2809     10.31239
           ---------+----------------------------------------------------------
```

```
                   |      Coef.    Std. Err.       t     P>|t|      [95% Conf. Interval]
-------------------+---------------------------------------------------------------------- -
D_i3               |
          pcip     |
            L1.    |     .0030131    .0016249     1.85    0.064     -.0001787     .0062049
            L2.    |     .0032639    .0016908     1.93    0.054     -.0000574     .0065852
            L3.    |    -.0003259     .001608    -0.20    0.839     -.0034846     .0028328
                   |
          pcsp     |
            L1.    |     .0007001    .0005053     1.39    0.166     -.0002924     .0016927
            L2.    |     .0017094    .0005198     3.29    0.001      .0006882     .0027305
            L3.    |    -.0006568    .0005078    -1.29    0.196     -.0016543     .0003407
                   |
            i3     |
            LD.    |     .3029408      .04315     7.02    0.000      .2181797     .3877019
           L2D.    |    -.1863638    .0441125    -4.22    0.000     -.2730155    -.0997121
           L3D.    |    -.0050184    .0433021    -0.12    0.908     -.0900782     .0800415
                   |
          _cons    |    -.0304459    .0211456    -1.44    0.150     -.0719829     .0110911
-------------------+---------------------------------------------------------------------- '
```

- We estimate all relevant equations in one step. Observe that the VAR command has no robust option.
- VAR assumes that the error terms are stationary.
- Now, we only need to ask STATA to perform the exclusion tests. This allows us to determine whether certain variables should be excluded from the VAR model because they do not contribute significantly.
- To perform the exclusion test for **VAR analysis: Wald Test (which is an F-test)**

```
. vargranger

  Granger causality Wald tests
+--------------------------------------------------------------------------+
|          Equation          Excluded |     F       df    df_r   Prob > F  |
|-------------------------------------+------------------------------------|
|              pcip              pcsp  |  5.1703      3     544    0.0016   |
|              pcip              D.i3  |  1.4741      3     544    0.2206   |
|              pcip               ALL  |   3.479      6     544    0.0022   |
|-------------------------------------+------------------------------------|
|              pcsp              pcip  |  .28184      3     544    0.8385   |
|              pcsp              D.i3  |  5.6388      3     544    0.0008   |
|              pcsp               ALL  |  3.2168      6     544    0.0041   |
|-------------------------------------+------------------------------------|
|              D_i3              pcip  |  3.6886      3     544    0.0119   |
|              D_i3              pcsp  |  5.3425      3     544    0.0012   |
|              D_i3               ALL  |  4.9576      6     544    0.0001   |
+--------------------------------------------------------------------------+
```

Now, we can interpret the above equations. The following line, for example:

```
|          Equation          Excluded |     F       df    df_r   Prob > F  |
|-------------------------------------+------------------------------------|
```

is the result of the F-test of the joint significance of the three lags of the variable pcsp in the equation with pcip as dependent variable. We can see that the p-value is less than 0.01, hence we can say that at 1% pcsp Granger causes pcip, with the past values of Δi3 fixed. We find that the change of interest rates does not Granger-cause pcip.

We find, however, that the change in S&P 500 index is not Granger caused by the changes in the industrial production, but is Granger caused by the change of interest rate.

Finally, the change of the 3 month Treasury Bill interest rates is Granger caused by both pcip and pcsp at 5%.

## C.**18.8 Use the data in FERTIL3.RAW.**
i)Graph gfr against time. Does it contain a clear upward or downward trend over the entire sample period?

(i)      We plot gfr (gross fertility rate) against time.

. line gfr year



. line gfr year||lfit gfr year



There is a negative long-term in the data. We can even make this more apparent by introducing a linear time trend in the plot

But the negative trend is not clear since there is a period of increasing fertility after circa 1940 until the late 1950s.

(ii) Using the data through 1979, estimate a cubic time trend model for *gfr* (that is, regress *gfr* on $t$, $t^2{}_2$, and $t^3$ along with an intercept). Comment on the *R*-squared of the regression.

```
. reg gfr t tsq tcu if year<1980

      Source |       SS          df       MS                Number of obs =      67
-------------+------------------------------              F(  3,     63) =   59.47
       Model |  17288.927         3   5762.97566           Prob > F       =  0.0000
    Residual |  6104.84329        63   96.9022745           R-squared      =  0.7390
-------------+------------------------------              Adj R-squared =  0.7266
       Total |  23393.7703        66   354.451065           Root MSE       =  9.8439


-------------------------------------------------------------------------------
         gfr |      Coef.   Std. Err.      t     P>|t|     [95% Conf. Interval]
-------------+-----------------------------------------------------------------
           t |  -6.904217   .6438123    -10.72   0.000    -8.190773   -5.617661
         tsq |   .2426157   .0219125     11.07   0.000     .1988271    .2864042
         tcu |  -.0024194   .0002119    -11.42   0.000    -.0028429   -.0019959
       _cons |   148.7082   5.092812     29.20   0.000     138.5311    158.8854
-------------------------------------------------------------------------------
```

If we use the usual *t* critical values, all terms are very statistically significant, and the *R*-squared indicates that this curve-fitting exercise tracks $gfr_t$ pretty well, at least up through 1979.

```
. hettest
```

Breusch-Pagan / Cook-Weisberg test for heteroskedasticity
       Ho: Constant variance
       Variables: fitted values of gfr

      chi2(1)    =    0.38
      Prob > chi2  =   0.5401

```
. bgodfrey, lags(1/3)
```

Breusch-Godfrey LM test for autocorrelation

| lags(p) | chi2   | df | Prob > chi2 |
|---------|--------|----|-------------|
| 1       | 50.801 | 1  | 0.0000      |
| 2       | 51.222 | 2  | 0.0000      |
| 3       | 51.496 | 3  | 0.0000      |

H0: no serial correlation

**Heteroscedasticity is no problem.**

But there is autoregression in the error-term signifying specification problems.

**iii) Using this model $gfr_t = \beta_0 + \beta_1 t + \beta_2 t^2 + \beta_3 t^3$ compute the mean absolute error of the one step ahead forecast error for the years 1980 through 1984**

<mark>Solution</mark>

**Manually calculation**

| Year | GFR actual $\overset{...}{y}$ | Predicted GFR ("y-hat") | Absolute forecast error |
|------|------|------|------|
| 1979 | 67.2 | 47.6 | |
| 1980 | 68.4 | 40.3 | 28.1 |
| 1981 | 67.4 | 32.6 | 34.8 |
| 1982 | 67.3 | 24.4 | 42.9 |
| 1983 | 65.8 | 15.6 | 50.2 |
| 1984 | 65.4 | 6.3 | 59.1 |

**STATA**

```
. predict forecast if year>1979, xb
(67 missing values generated)

. gen abserror=abs(gfr-forecast)
(67 missing values generated)
```

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|------|------|------|------|------|------|
| abserror_iii | 5 | 43.01686 | 12.2686 | 28.06189 | 59.11368 |

The MAE is 43.02 (see STATA outcome), and the model performs badly in forecasting.
A high R2 is no guarantee of a good forecasting capability.

**iv) Using the data through 1979, regress $\Delta gfr_i$ on a constant only. Is the constat statistically sgnificant different from zero? Does it make sense to assume that any drift term is zero, if we assume that $gfr_i$ follows a random walk?**

We should estimate the $\Delta \boldsymbol{gfr_i}$ against a consant , using data up through 1and and decide if a drift is needed if we assume that gfr follows a random walk process.

```
. reg d.gfr if year<1980

      Source |       SS           df       MS                  Number of obs =        66
-------------+------------------------------                   F(  0,      65) =      0.00
       Model |          0           0          .               Prob > F        =          .
    Residual |  1264.15569          65   19.448549             R-squared       =     0.0000
-------------+------------------------------                   Adj R-squared   =     0.0000
       Total |  1264.15569          65   19.448549             Root MSE        =     4.4101


------------------------------------------------------------------------------
       D.gfr |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+----------------------------------------------------------------
       _cons |  -.8712121    .5428397    -1.60   0.113    -1.955338    .2129137
------------------------------------------------------------------------------
```

The R-squared is identically zero since there are no explanatory variables. But $\hat{\sigma}$ which estimates the standard deviation of the error, is comparable to that in part (ii), and we see that it is much smaller here.) The t statistic for the intercept is about −1.60, which is not significant at the 10% level against a two-sided alternative. Therefore, it is legitimate to treat gfrt as having no drift, if it is indeed a random walk. (That is, if gfrt = $\alpha_0$ + $gfrt_{t-1}$ + $e_t$, where {et} is zero-mean, serially uncorrelated process, then we cannot reject H0: $\alpha_0$ = 0.)

**v) Now, forecast, $gfr$ for 1980 through 1984, using a random walk model: the forecast of $gfr_{n+1}$ is simply $gfr_n$. Find the MAE. How does it compare with the MAE from part iii). Which method of forecasting do you prefer?**

Now we should do the forecasting for 1980-84, but this time under the assumption that GFR is a random walk process without drift. In other words, our forecast for period t+1 is the value in t.

| Year | GFR actual (1) | GFR forecast (2) | Diff (1-2) | Absolute forecast error |
|---|---|---|---|---|
| 1979 | 67.2 | | | |
| 1980 | 68.4 | 67.2 | 1.2 | 1.2 |
| 1981 | 67.4 | 68.4 | -1 | 1.0 |
| 1982 | 67.3 | 67.4 | -0.1 | 0.1 |
| 1983 | 65.8 | 67.3 | -1.5 | 1.5 |
| 1984 | 65.4 | 65.8 | -0.4 | 0.4 |

The MAE is 0.84, that is, the random-walk model outperforms the deterministic trend model in terms of forecasting accuracy. With stata, we can calculate the absolute forecasting errors as the absolute value of the differenced gfr series, since the gfr(t) is the actuals series and gfr(t-1) is the forecast for period t.

That is the forecast error for period t is: $\hat{e}_t = y_t - y_{t-1} = \Delta y_t$

Hence the mean absolute error of the forecast is:

$$\frac{1}{5} \sum_{i=1980}^{1984} \hat{e}_i = \frac{1}{5} \sum_{i=1980}^{1984} \Delta y_i$$

We can calculate this in stata:

```
. gen absfe=abs(d.gfr) if year>1979
(67 missing values generated)

. sum absfe
```

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|---|---|---|---|---|---|
| absfe | 5 | .8400009 | .5770624 | .0999985 | 1.5 |

**vi) Now, estimate an AR(2) model for gfr, again using the data only through 1979. Is the second lag significant?**

<mark>Solution:</mark>  We need to estimate an AR(2) model for gfr until 1979.

```
. reg gfr l(1/2).gfr if year<1980

      Source |       SS           df       MS                  Number of obs =      65
-------------+------------------------------------            F(  2,     62) =  571.67
       Model |  20669.7236         2   10334.8618             Prob > F       =  0.0000
    Residual |  1120.86601        62   18.078484             R-squared      =  0.9486
-------------+------------------------------------            Adj R-squared =  0.9469
       Total |  21790.5896        64   340.477963             Root MSE       =  4.2519


-------------------------------------------------------------------------------------
         gfr |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-------------+-----------------------------------------------------------------
         gfr |
         L1. |   1.272076   .1203391    10.57   0.000     1.031522    1.512631
         L2. |  -.3113864   .1213988    -2.56   0.013    -.5540592   -.0687136
             |
        _cons |   3.215658   2.924166     1.10   0.276    -2.629667    9.060983
-------------------------------------------------------------------------------------
```

The second lag is significant. (Recall that its $t$ statistic is valid even though *gfrt* apparently contains a unit root: the coefficients on the two lags sum to .961.) The standard error of the regression is slightly below that of the random walk model.

**vii) Obtain the MAE for 1980 through 1984, using the AR(2) model. Does this more general model work better out-of-sample than the random walk model?**

**Solution:**   Now we are asked to calculate the MAE for the forecasts from the AR(2) model.

| Year | GFR actual (1) | GFR forecast (2) | Diff (1-2) | Absolute forecast error |
|------|------|------|------|------|
| 1978 | 65.5 | | | |
| 1979 | 67.2 | | | |
| 1980 | 68.4 | 68.3 | 0.1 | 0.1 |
| 1981 | 67.4 | 69.3 | -1.9 | 1.9 |
| 1982 | 67.3 | 67.7 | -0.4 | 0.4 |
| 1983 | 65.8 | 67.8 | -2.0 | 2.0 |
| 1984 | 65.4 | 66.0 | -0.6 | 0.6 |

The out-of-sample forecasting performance of the AR(2) model is worse than the random walk without drift: the MAE for 1980 through 1984 is about .991 for the AR(2) model.

And the MAE is 1, which is slightly higher than the MAE from our forecast under the random walk assumption. Even though the AR(2) model seems more sophisticated than the random-walk, it cannot outperform the random-walk in terms of forecasting.

With stata:

```
. predict forecastar2 if year>1979, xb
(67 missing values generated)

. gen absfe2=abs(forecastar2-gfr) if year>1979
(67 missing values generated)

. mean absfe2

Mean estimation                           Number of obs    =         5

-----------------------------------------------------------------
             |       Mean    Std. Err.      [95% Conf. Interval]
-------------+---------------------------------------------------
      absfe2 |   .9905746    .4070746       -.1396457    2.120795
-----------------------------------------------------------------
```