

Time series data

The Benefits and Problems of Persistence

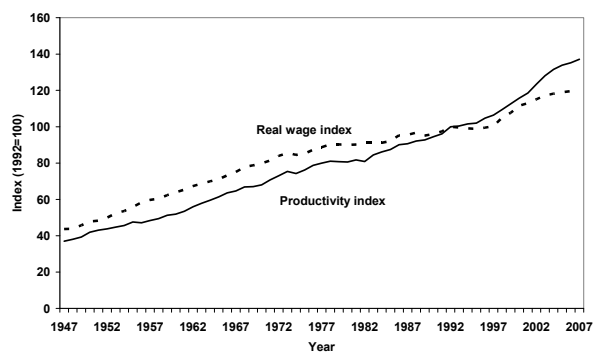
1

Key time series property: Persistence

- Persistence over time (trends)
- Persistence in cycle
- Both allow for excellent explanatory power of time series data
- Also produce persistence in errors
 - $\text{Cov}(\epsilon_{t-1}, \epsilon_t) \neq 0$
- Persistence can produce spurious correlation

2

Time Series: Productivity and Wage Index



3

```
* generate ln of outcomes
gen rwage=ln(rwage)

gen prodl=ln(productivity)

* run a model of productivity
* on just a time trend
* see how well it fits
reg prodl time
```

Look at the R-squared

Source	SS	df	MS	
Model	7.79439507	1	7.79439507	Number of obs = 61
Residual	.1113919	59	.001887998	F(1, 59) = 4128.39
Total	7.90578697	60	.131763116	Prob > F = 0.0000
				R-squared = 0.9859
				Adj R-squared = 0.9857
				Root MSE = .04345

prodl	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
time	.0203023	.000316	64.25	0.000	.01967 .0209346
_cons	3.684954	.0112649	327.12	0.000	3.662413 3.707495

Interpret coefficient

4

```

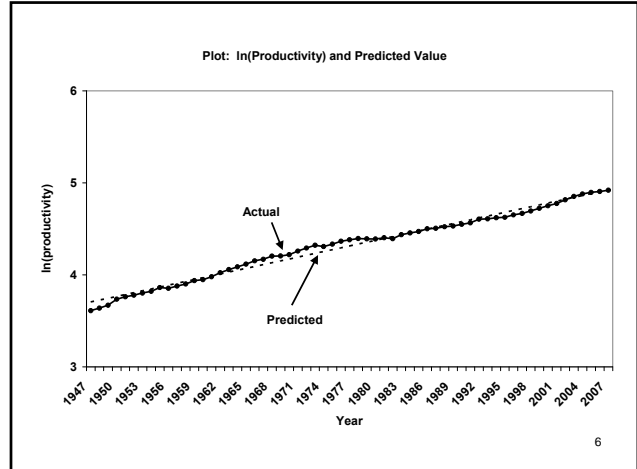
* output predicted value
predict prodl_pred
(option xb assumed, fitted values)

*output predicted values to csv
* file to graph in pretty graph
outsheet time year prodl prodl_pred using predict_value.csv, comma

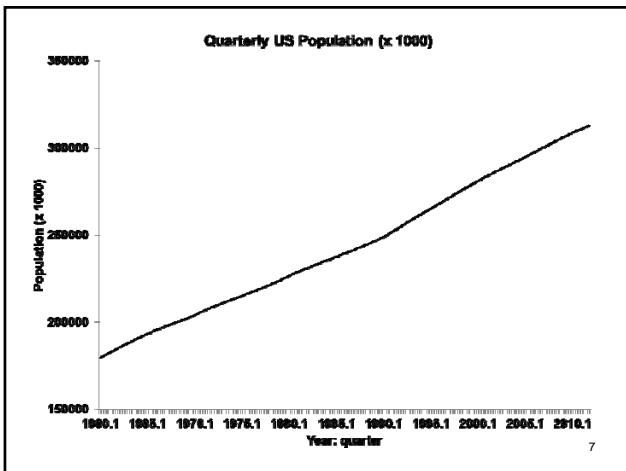
```

Output the predicted values from the regression, then
Output them to a csv file to be used in excel

5



6



7

```

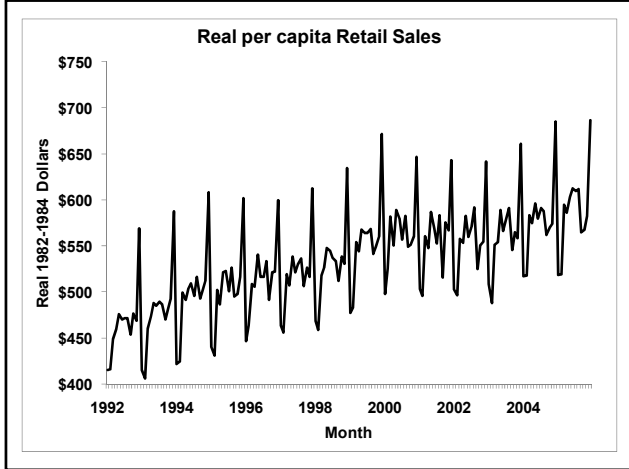
. reg population time

```

Source	SS	df	MS	
Model	3.0705e+11	1	3.0705e+11	Number of obs = 208
Residual	1.5447e+09	206	7498534.91	F(1, 206) = 40948.07
Total	3.0860e+11	207	1.4908e+09	Prob > F = 0.0000
				R-squared = 0.9950
				Adj R-squared = 0.9950
				Root MSE = 2738.3

population	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
time	639.8904	3.162196	202.36	0.000	633.656 646.1248
_cons	176904.6	301.1136	453.86	0.000	176033.4 177536.2

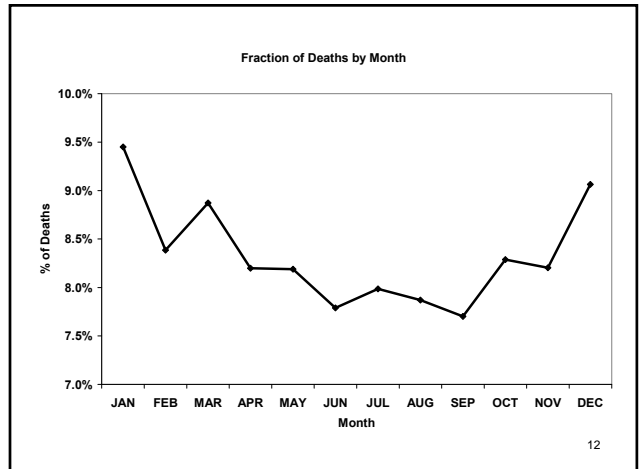
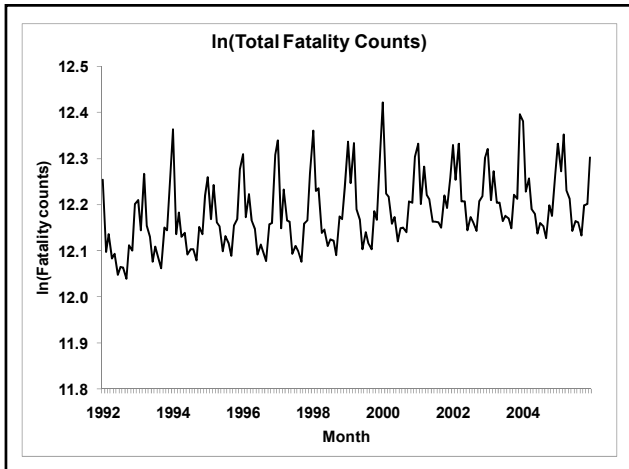
8

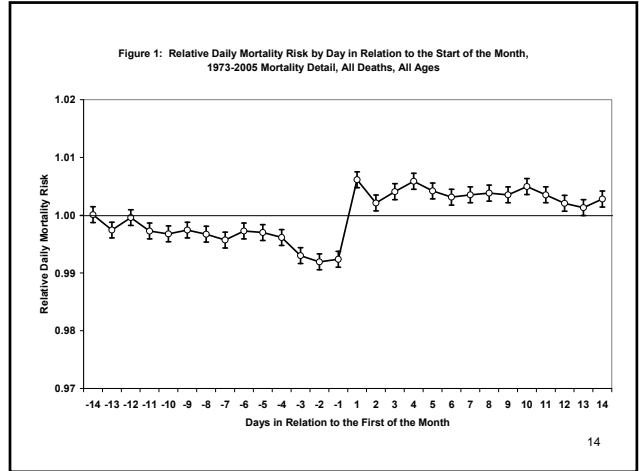
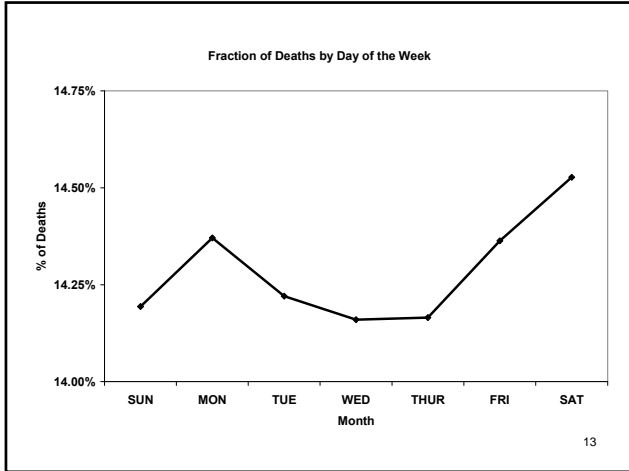


```
reg lr_retail_capita_I* time unrate
```

Source	SS	df	MS	Number of obs = 168
Model	2.06236449	13	.158643422	F(13, 154) = 514.11
Residual	.047521099	154	.000308579	Prob > F = 0.0000
Total	2.10988559	167	.012634045	R-squared = 0.9775
				Adj R-squared = 0.9756
				Root MSE = .01757

lr_retail_a	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
_Imonth_2	-.0202347	.0066395	-3.05	0.003	-.033351 - .0071183
_Imonth_3	.0899394	.0066398	13.55	0.000	.0768226 .1030562
_Imonth_4	.0861822	.00664	12.98	0.000	.073065 .0992995
_Imonth_5	.1351768	.0066405	20.36	0.000	.1220585 .1482952
_Imonth_6	.1094872	.006641	16.49	0.000	.0963681 .1226064
_Imonth_7	.1074848	.0066416	16.18	0.000	.0943643 .1206053
_Imonth_8	.1316515	.0066427	19.82	0.000	.118529 .144774
_Imonth_9	.0720061	.0066434	10.84	0.000	.0588821 .0851301
_Imonth_10	.1109814	.0066449	16.70	0.000	.0978546 .1241082
_Imonth_11	.1427717	.0066459	21.48	0.000	.1296427 .1559006
_Imonth_12	.3417103	.0066471	51.41	0.000	.328579 .3548416
time	.0013152	.0000319	41.24	0.000	.0012522 .0013782
unrate	-.0049376	.001584	-3.12	0.002	-.0080668 -.0018083
_cons	8.11711	.0112221	723.31	0.000	8.094941 8.139279





```

* get log of counts
gen lcounts=ln(counts)

* construct month and weekday dummy variables
xi i.month i.weekday
i.month      _Imonth_1-12      (naturally coded; _Imonth_1 omitted)
i.weekday    _Iweekday_1-7     (naturally coded; _Iweekday_1 omitted)

* construct variable for the 1st week of the month
gen firstweek=(day>=1&day<=7)

gen sept911=(month==9&day==11&year==2001)

```

Given a variable x that has k unique values, the xi command in STATA will generate a series of k-1 dummy variables, one for each unique value With the smallest value set as the reference group

```

reg lcounts trend sept911 firstweek _I*

```

Source	SS	df	MS	Number of obs = 12053		
Model	127.396712	20	6.36983559	F(20, 12032) = 4084.99		
Residual	18.7618363	12032	.001559328	Prob > F = 0.0000		
Total	146.158548	12052	.012127327	R-squared = 0.8716		
				Adj. R-squared = 0.8714		
				Root MSE = .03949		

	lcounts	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
trend		.0007336	3.15e-06	233.04	0.000	.0007275 .0007398
sept911		.3857986	.0395213	9.76	0.000	.3083306 .4632667
firstweek		.0052759	.0008548	6.17	0.000	.0036003 .0069515
_Imonth_2		-.0260462	.0017882	-14.57	0.000	-.0295514 -.022541
_Imonth_3		-.0638372	.001746	-36.56	0.000	-.0672597 -.0604147

DELETE SOME RESULTS

_Imonth_11		-.11388	.0017608	-64.68	0.000	-.1173314 -.1104286
_Imonth_12		-.0493627	.0017464	-28.27	0.000	-.0527859 -.0459396
_Iweekday_2		.0113788	.001346	8.45	0.000	.0087405 .0140171

DELETE SOME RESULTS

Iweekday_6		.010537	.001346	7.83	0.000	.0078987 .0131753
_Iweekday_7		.0225822	.001346	16.78	0.000	.0199439 .0252205
_c_oms		8.641988	.001646	5250.26	0.000	8.638761 8.645214

Great Britain --1947

- Raised compulsory education age from 14-15
- Within a year – education levels for the country increased
- Did this generate better outcomes?

17

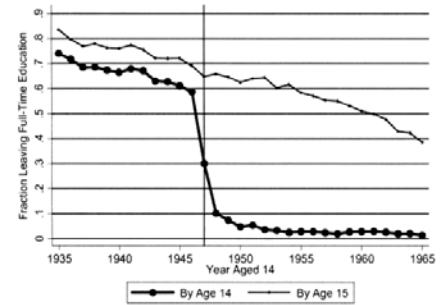


FIGURE 1. FRACTION LEFT FULL-TIME EDUCATION BY YEAR AGED 14 AND 15 (Great Britain)

Note: The lower line shows the proportion of British-born adults aged 32 to 64 from the 1983 to 1998 General Household Surveys who report leaving full-time education at or before age 14 from 1935 to 1965. The upper line shows the same, but for age 15. The minimum school-leaving age in Great Britain changed in 1947 from 14 to 15.

18

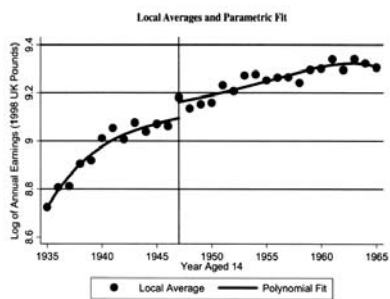
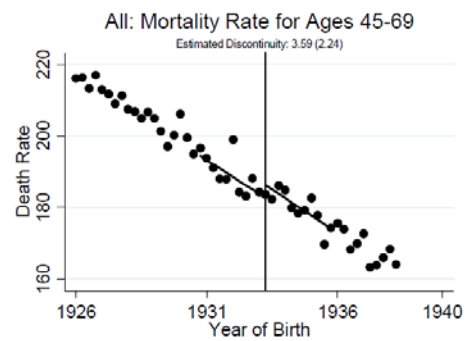


FIGURE 6. AVERAGE ANNUAL LOG EARNINGS BY YEAR AGED 14 (Great Britain)

19



20

Figure 1: % Newborn Discharged Early, Vaginal Deliveries without Complications

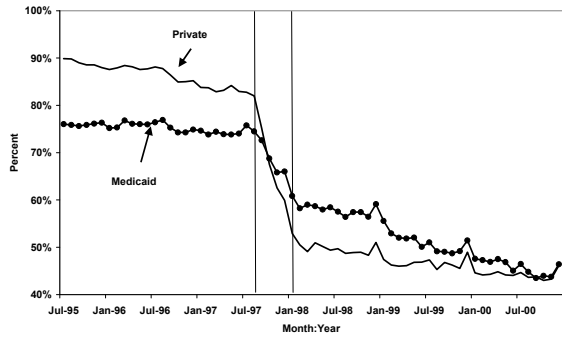
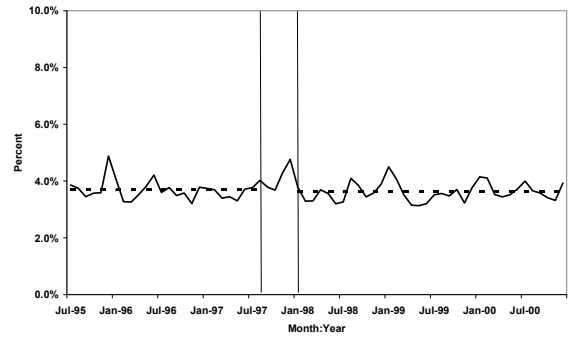
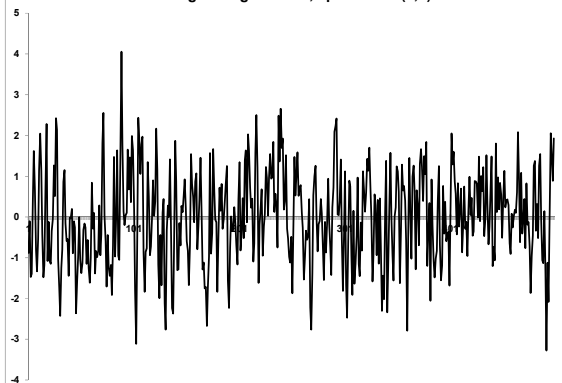


Figure 4: % Newborns Readmitted within 28-Days, Vaginal Deliveries without Complications, Private Insurance

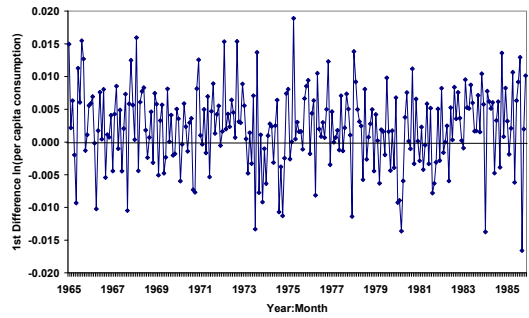


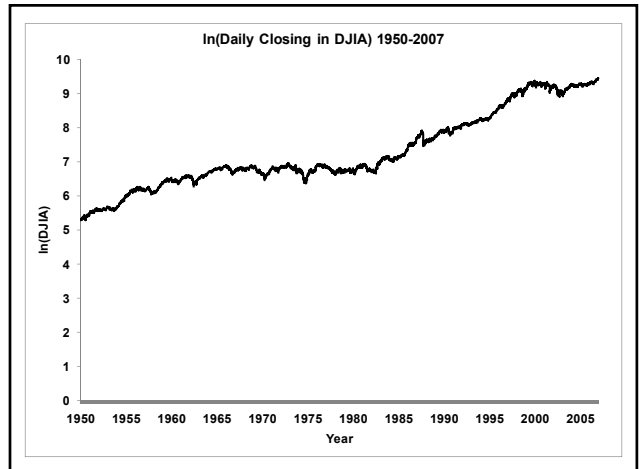
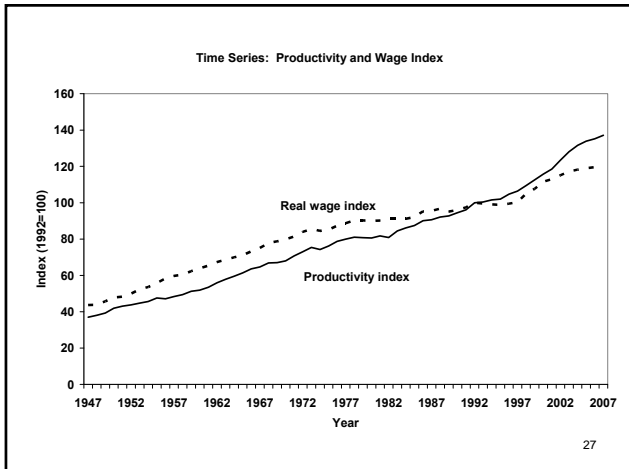
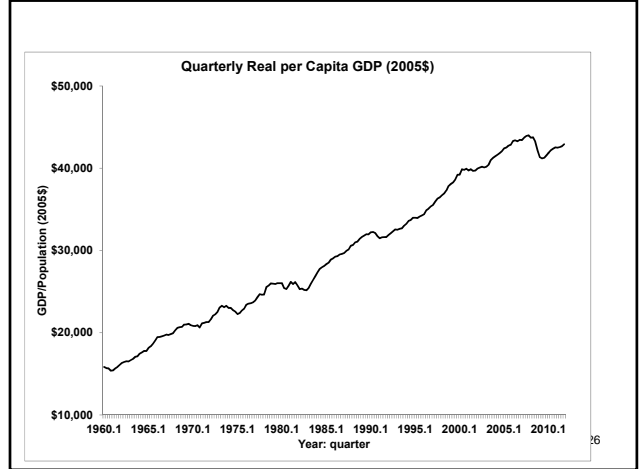
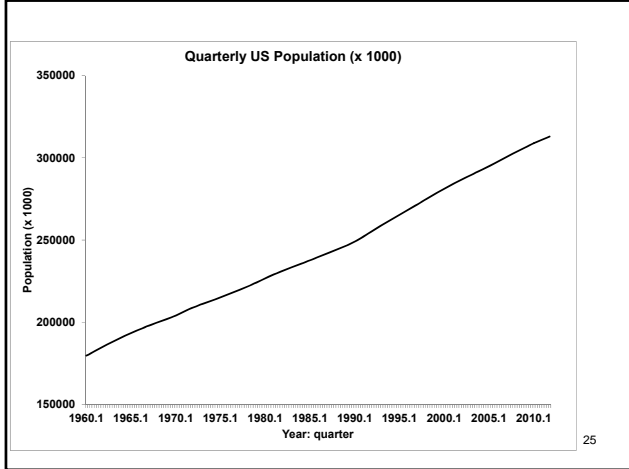
$$y_t = \varepsilon_t + \alpha \varepsilon_{t-1}$$

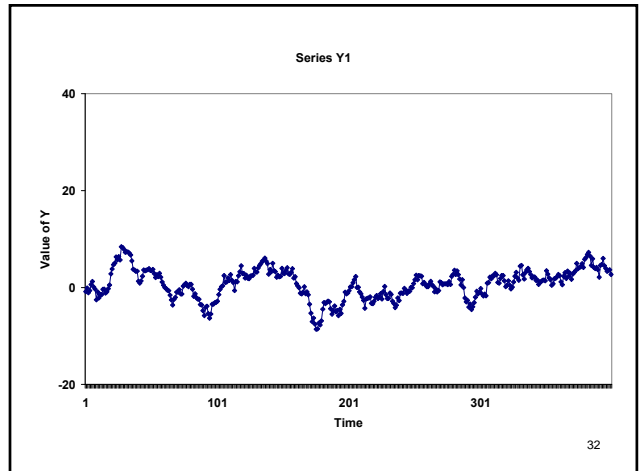
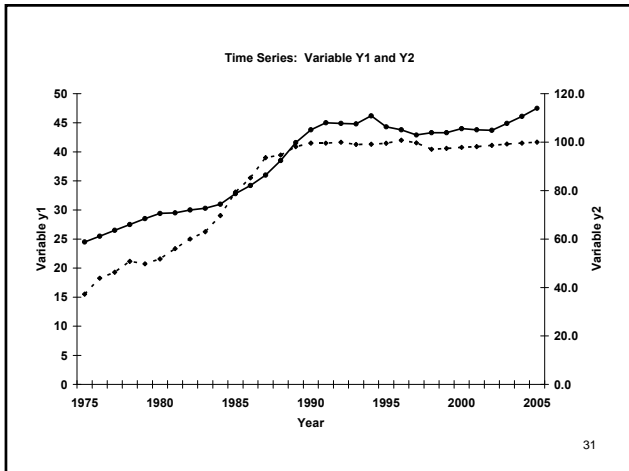
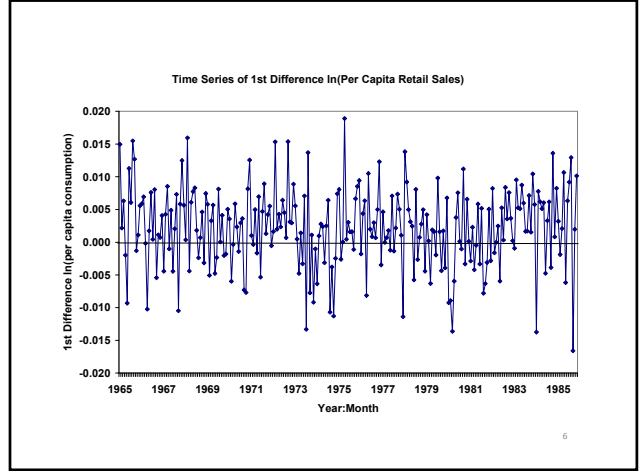
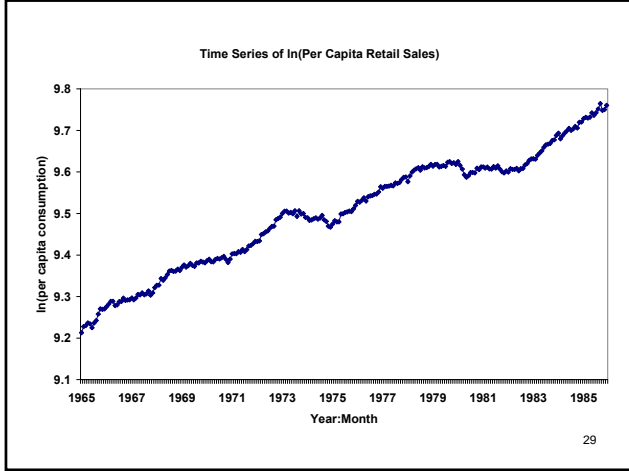
Moving average: $\alpha=0.5$, epsilon $\varepsilon \sim N(0,1)$

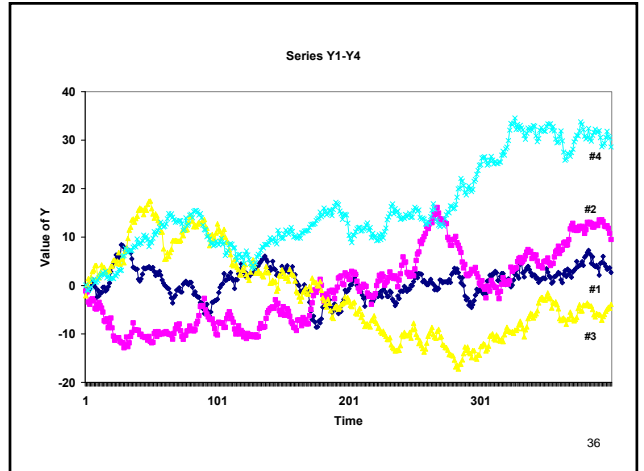
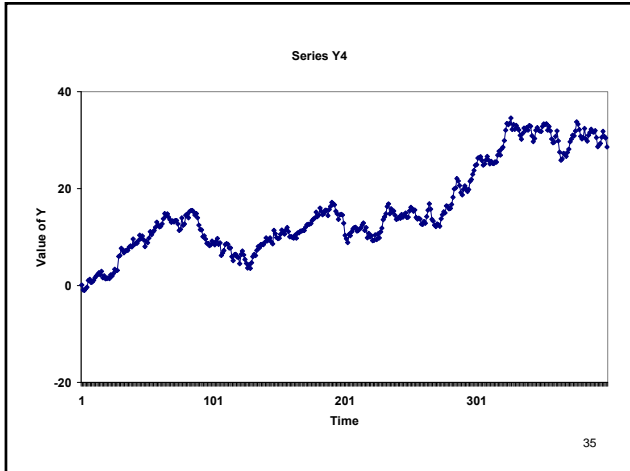
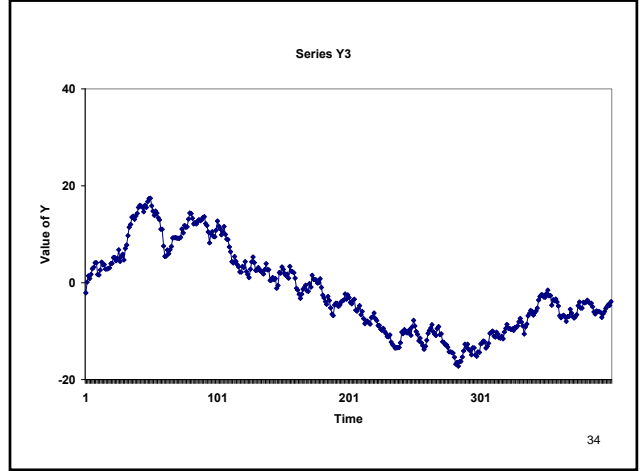
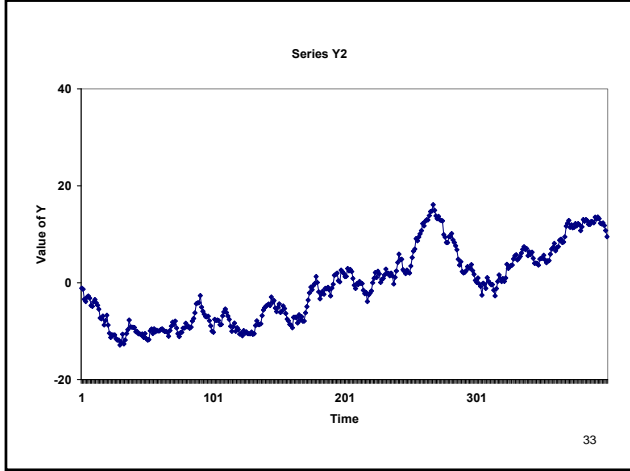


Time Series of 1st Difference ln(Per Capita Retail Sales)





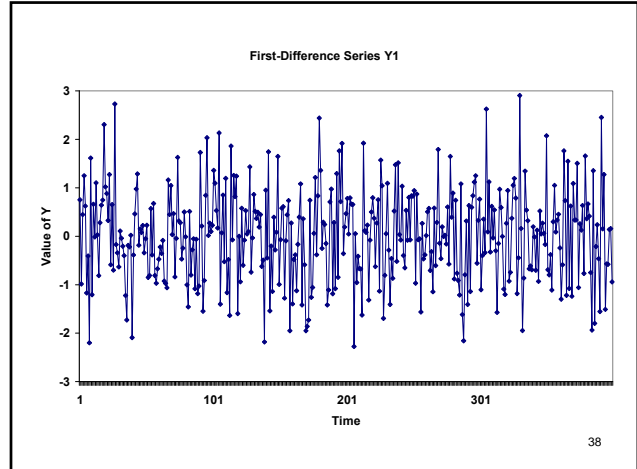




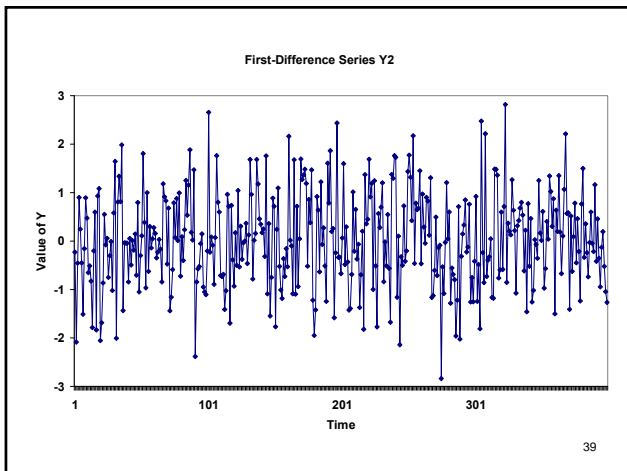
$$y_t = a + \rho y_{t-1} + v_t$$

Series	Coefficient (std error) on ρ	P-value, $H_0: \rho=1$
Y1		
Y2		
Y3		
Y4		

37



38



39

OLS – Levels on Levels

coefficients (standard errors)

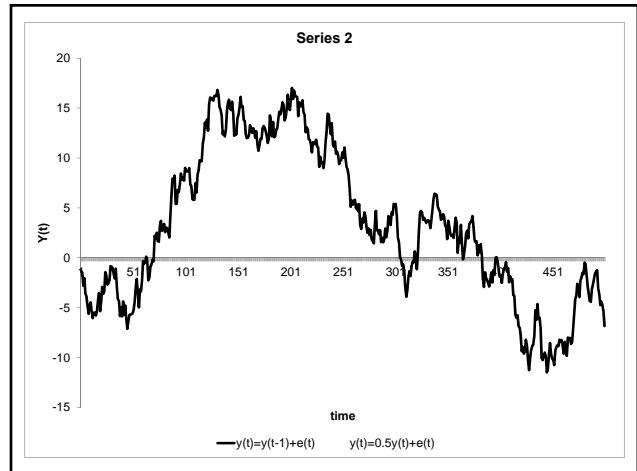
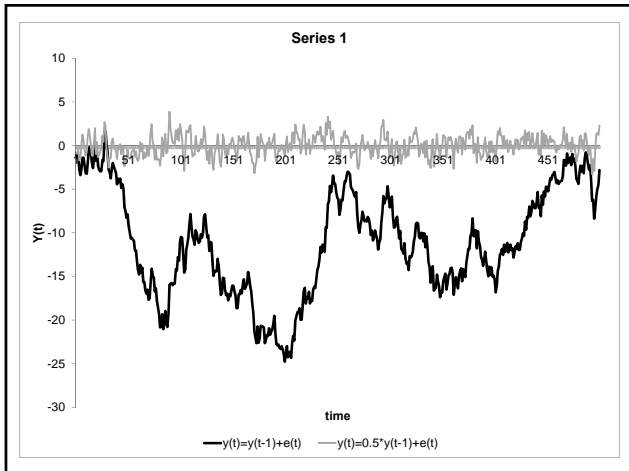
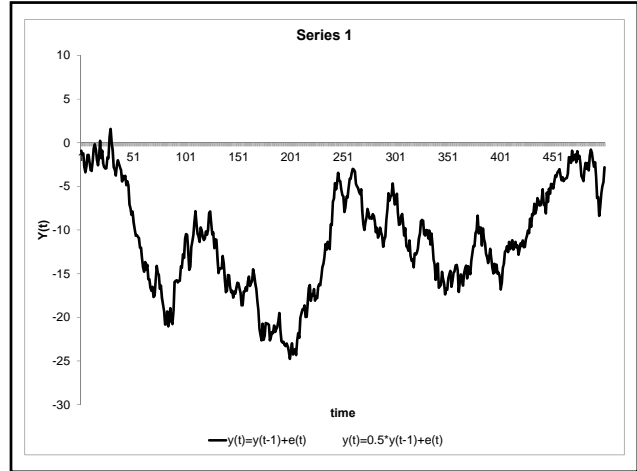
Ind. Var.	Dependent Variable			
	y1	y2	y3	y4
y1		0.134 (0.124)	0.274 (0.140)	0.461 (0.149)
y2	0.022 (0.020)		-0.854 (0.037)	0.823 (0.045)
y3	0.034 (0.018)	-0.668 (0.029)		-0.532 (0.047)
y4	0.051 (0.016)	0.559 (0.030)	-0.463 (0.041)	

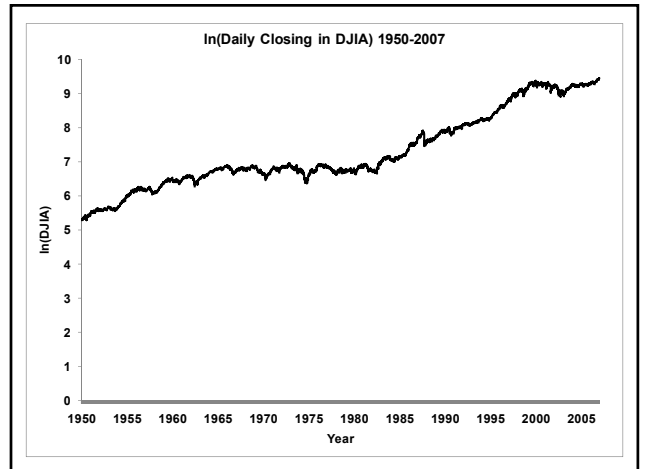
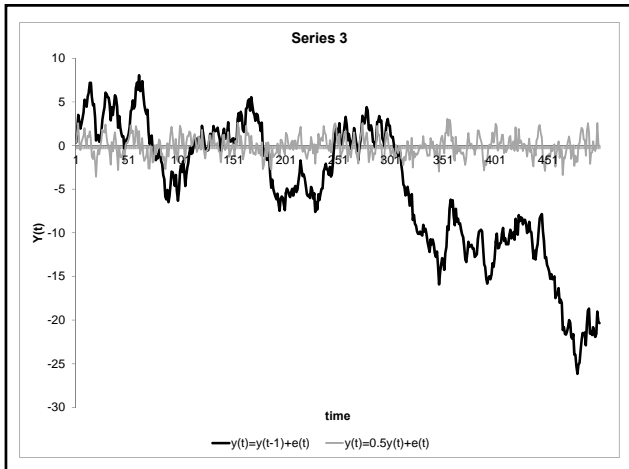
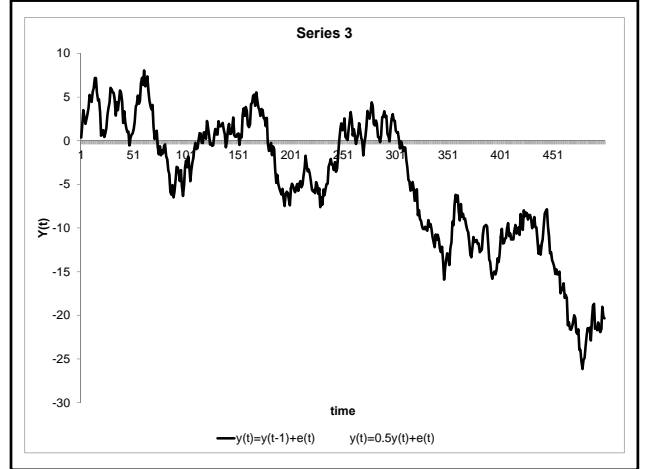
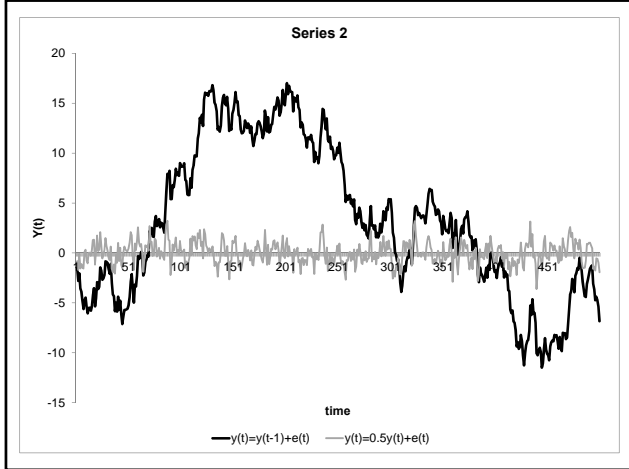
40

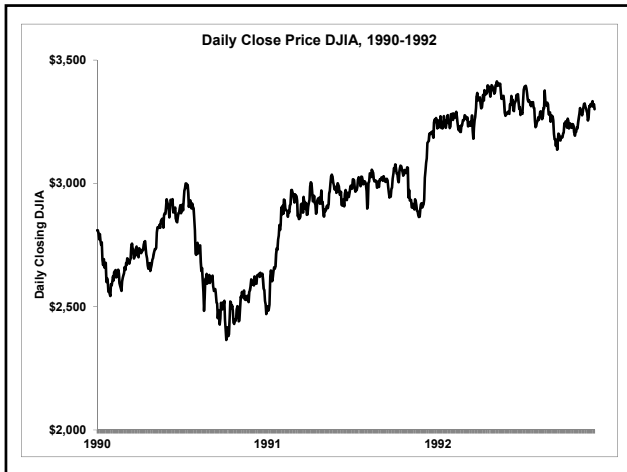
OLS – Differences on Differences coefficients (standard errors)

Ind. Var.	Dependent Variable			
	Δy_1	Δy_2	Δy_3	Δy_4
Δy_1		-0.046 (0.051)	0.045 (0.051)	0.006 (0.047)
Δy_2	-0.044 (0.049)		0.011 (0.051)	0.005 (0.047)
Δy_3	0.043 (0.049)	0.011 (0.049)		-0.003 (0.047)
Δy_4	0.067 (0.052)	0.006 (0.053)	-0.004 (0.054)	

41







```

. * notice the variable time. the first day is
. * 1, the second day 2, etc.
.
. * define the data as time series
. * the variable time is the index
. tsset time
      time variable: time, 1 to 14362
      delta: 1 unit

```

Generating lags in STATA

- Data set needs to be identified as a time-series (tsset)
- Given variable z
- One period lag
gen zlag1=z[_n-1]
- Five period lag
gen zlag5=z[_n-5]

Model (1)

$$y_t = \alpha + \rho y_{t-1} + \varepsilon_t$$

$$H_o : \rho = 1$$

$$H_a : \rho < 1$$

Model (2)

$$y_t - y_{t-1} = \alpha + \rho y_{t-1} - y_{t-1} + \varepsilon_t$$

$$\Delta y_t = \alpha + (\rho - 1)y_{t-1} + \varepsilon_t$$

$$\Delta y_t = \alpha + \theta y_{t-1} + \varepsilon_t$$

$$H_o : \theta = 0$$

$$H_a : \theta < 0$$

```

* take the ln of the daily closing price
gen ln_close=ln(close)

* get the 1st difference
gen ln_close_1=ln_close[_n-1]

* get the 1st difference
gen d_ln_close=ln_close-ln_close_1

* get the lag of the 1st difference
gen d_ln_close_1=d_ln_close[_n-1]

```

```

. * test for random walk
. * run a regression of change ln(closing price)
. * on one period lag
. reg ln_close ln_close1

```

Source	SS	df	MS	Number of obs = 14361
Model	18001.7298	1	18001.7298	F(1, 14359) = 0.0000
Residual	1.18771395	14359	.000082716	Prob > F = 0.9999
Total	18002.9175	14360	1.25368507	R-squared = 0.9999
				Adj R-squared = 0.9999
				Root MSE = .00909

```

. ln_close | Coef. Std. Err. t P>|t| [95% Conf. Interval]
-----+-----
ln_close1 | .9999874 .0000678 . 0.000 .9998545 1.00012
 _cons | .003808 .0004995 0.76 0.446 -0.005984 .00136

```

```

. test ln_close1=1
(1) ln_close1 = 1
F( 1, 14359) = 0.03
Prob > F = 0.8520

```

```

. * now run model where null is transformed into 0
. reg d_ln_close ln_close_1

```

Source	SS	df	MS	Number of obs = 14361
Model	2.8777e-06	1	2.8777e-06	F(1, 14359) = 0.03
Residual	1.18771395	14359	.000082716	Prob > F = 0.8520
Total	1.18771682	14360	.00008271	R-squared = -0.0001
				Adj R-squared = -0.0001
				Root MSE = .00909

```

. d_ln_close | Coef. Std. Err. t P>|t| [95% Conf. Interval]
-----+-----
ln_close_1 | -.0000126 .0000678 -0.19 0.852 -.0001455 .0001202
 _cons | .0003808 .0004995 0.76 0.446 -.0005984 .00136

```

```

. * now get dickey fuller test
. dfuller ln_close

```

Dickey-Fuller test for unit root

Test Statistic	Interpolated Dickey-Fuller		
	1% Critical Value	5% Critical Value	10% Critical Value
Z(t)	-0.187	-3.430	-2.860

MacKinnon approximate p-value for Z(t) = 0.9401

```

gen d_ln_close_1=d_ln_close[_n-1]
(* missing value generated)

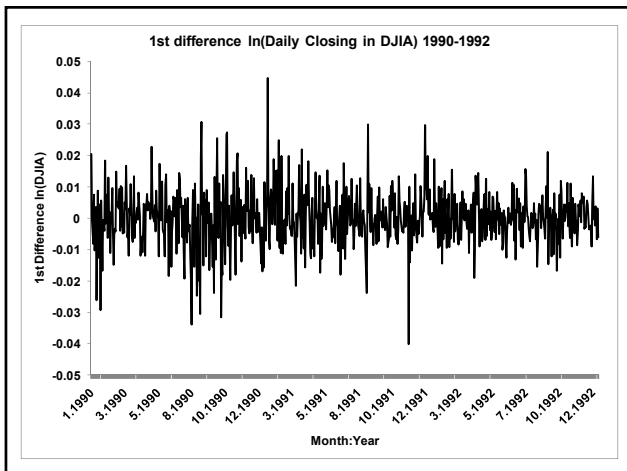
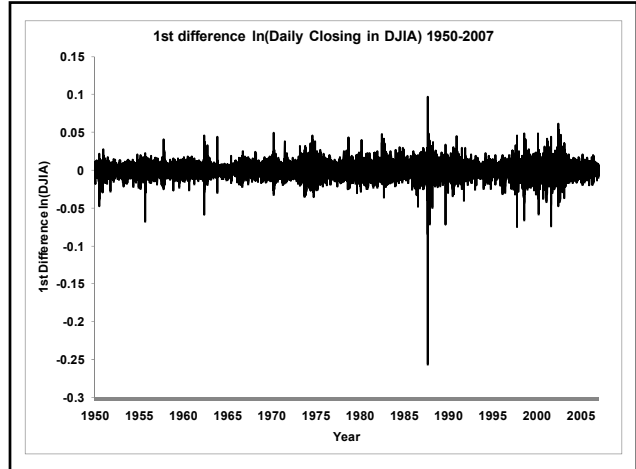
. reg d_ln_close d_ln_close_1

-----+-----
Source |      SS          df       MS          Number of obs = 14361
-----+-----
Model |   .005717215      1   .005717215          F( 1, 14359) = 69.45
Residual |  1.18199961 14359   .00082318          Prob > F      = 0.0000
Total |  1.18771682 14360   .00082271          R-squared     = 0.0048
                                          Adj R-squared = 0.0047
                                          Root MSE    = .00907

-----+-----
d_ln_close |      Coef.   Std. Err.      P>|t|   [95% Conf. Interval]
-----+-----
d_ln_close_1 |  .0693792   .008325   8.33   0.000   .0530611   .0856972
Constant |  -.0026877   .0007167  3.75   0.000   -.0041202  -.0004171

```

57



productivity.dta

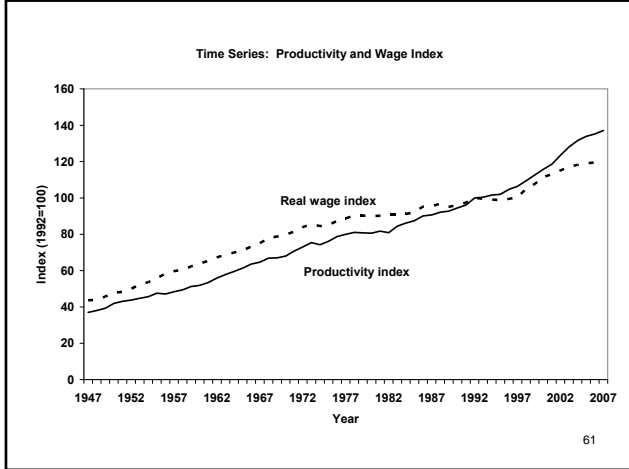
- Annual data on productivity index and mean real wages, 1947-2007
- Index variable time =1 in 1947, 2=1948....
- Need to “tsset” time series data and identify the index that orders time

```

. tsset time
.         time variable:  time, 1 to 61
.         delta: 1 unit

```

60



```

* generate ln of outcomes
gen rwage1=ln(rwage)
gen prod1=ln(productivity)

reg rwage1 time prod1

```

Notice the excellent fit

Source	SS	df	MS	Number of obs = 61		
Model	4.65023272	2	2.32511636	F(2, 58)	= 2083.25	
Residual	.06473372	58	.001116099	Prob > F	= 0.0000	
Total	4.71496644	60	.078582774	R-squared	= 0.9863	
				Adj R-squared	= 0.9858	
				Root MSE	= .03341	

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
time	-.0147661	.0020467	-7.21	0.000	-.018863	-.0106692
prod1	1.479204	.1000978	14.78	0.000	1.278837	1.679572
_cons	-1.531055	.3689574	-4.15	0.000	-2.269604	-.7925068

Interpret the coef. on prod1

62

```

* regress levels on 1st differences
reg rwage1 time rwage1l

```

Source	SS	df	MS	Number of obs = 60		
Model	4.32027231	2	2.16013616	F(2, 57)	= 12947.36	
Residual	.009509872	57	.00016684	Prob > F	= 0.0000	
Total	4.32978218	59	.073386139	R-squared	= 0.9978	
				Adj R-squared	= 0.9977	
				Root MSE	= .01292	

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
time	.0004911	.0003717	1.32	0.192	-.0002532	.0012354
rwage1l	.9446265	.0233767	40.41	0.000	.8978154	.9914375
_cons	.2443952	.0912928	2.68	0.010	.0615845	.4272059

```

reg prod1 time prod1l

```

Source	SS	df	MS	Number of obs = 60		
Model	7.39083745	2	3.69541873	F(2, 57)	= 18572.67	
Residual	.011341335	57	.000198971	Prob > F	= 0.0000	
Total	7.40217879	59	.125460657	R-squared	= 0.9985	
				Adj R-squared	= 0.9984	
				Root MSE	= .01411	

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
time	.001841	.0008647	2.13	0.038	.0001095	.0035726
prod1l	.9011966	.0422653	21.32	0.000	.8165617	.9858315
_cons	.389113	.1549295	2.51	0.015	.0788721	.699354

63

```

* regress 1st differences on lags and a trend
reg drwage1 time rwage1l

```

Source	SS	df	MS	Number of obs = 60		
Model	.003258489	2	.001629244	F(2, 57)	= 9.77	
Residual	.009509872	57	.00016684	Prob > F	= 0.0002	
Total	.012768361	59	.000216413	R-squared	= 0.2552	
				Adj R-squared	= 0.2291	
				Root MSE	= .01292	

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
time	.0004911	.0003717	1.32	0.192	-.0002532	.0012354
rwage1l	-.0553735	.0233767	-2.37	0.021	-.1021846	-.0085625
_cons	.2443952	.0912928	2.68	0.010	.0615845	.4272059

```

reg dprod1 time prod1l

```

Source	SS	df	MS	Number of obs = 60		
Model	.001579428	2	.000789714	F(2, 57)	= 3.97	
Residual	.011341335	57	.000198971	Prob > F	= 0.0243	
Total	.012920763	59	.000218996	R-squared	= 0.1222	
				Adj R-squared	= 0.0914	
				Root MSE	= .01411	

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
time	.001841	.0008647	2.13	0.038	.0001095	.0035726
dprod1l	-.0988034	.0422653	-2.34	0.023	-.1834383	-.0141685
_cons	.389113	.1549295	2.51	0.015	.0788721	.699354

64


```

. * now get dickey fuller tests with a trend
. dfuller rwagel, trend

Dickey-Fuller test for unit root          Number of obs   =      60

----- Interpolated Dickey-Fuller -----
      Test          1% Critical      5% Critical      10% Critical
Statistic          Value            Value            Value
-----
Z(t)              -2.369            -4.128            -3.490            -3.174
-----
MacKinnon approximate p-value for Z(t) = 0.3964

. dfuller prodl, trend

Dickey-Fuller test for unit root          Number of obs   =      60

----- Interpolated Dickey-Fuller -----
      Test          1% Critical      5% Critical      10% Critical
Statistic          Value            Value            Value
-----
Z(t)              -2.338            -4.128            -3.490            -3.174
-----
MacKinnon approximate p-value for Z(t) = 0.4132

```

65

```

. * now difference data
.
. gen drwagel=rwagel-rwagel1
(1 missing value generated)
. gen dprodl=prodl-prodl1
(1 missing value generated)

```

66

```

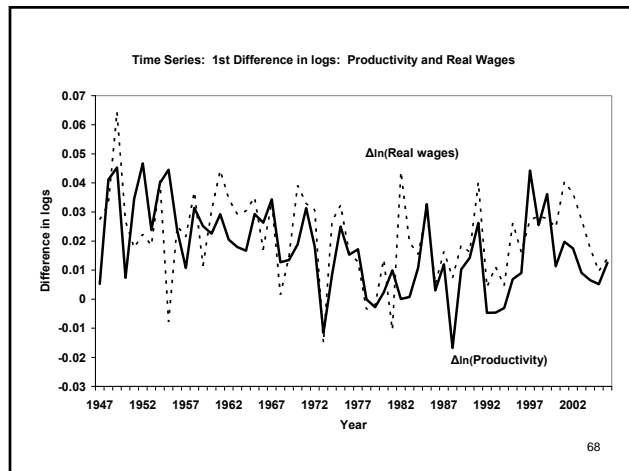
. * run regression of diff on diff
reg drwagel time dprodl

Source |      SS      df      MS              Number of obs =      60
-----+-----+-----+-----+-----
Model | .004894165    2  .002447082          F( 2, 57) = 17.71
Residual | .007874197   57  .000138144          Prob > F      = 0.0000
Total | .012768361   59  .000216413          R-squared     = 0.3833
                                           Adj R-squared = 0.3617
                                           Root MSE    = .01175

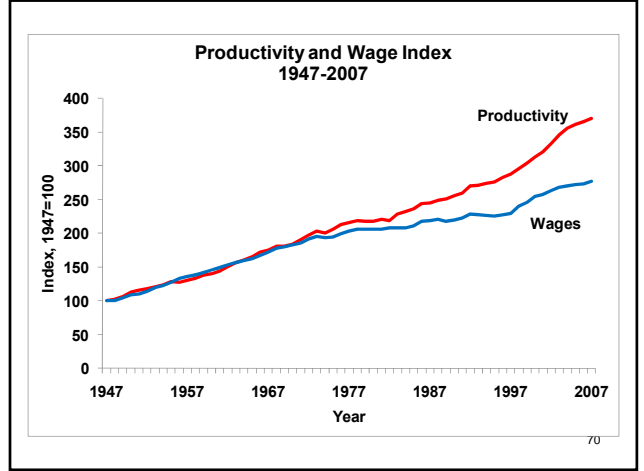
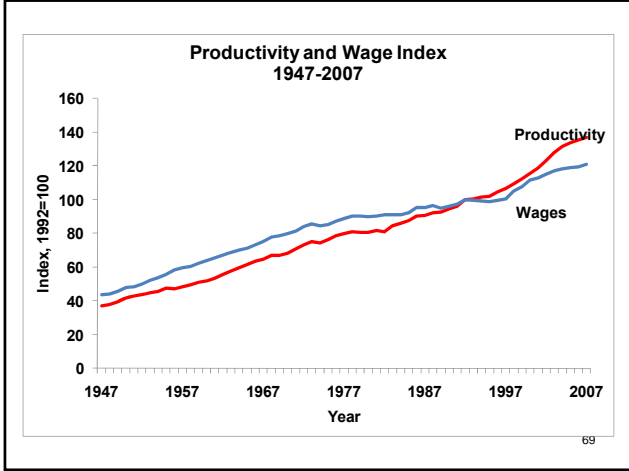
-----+-----+-----+-----+-----
drwagel |      Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
-----+-----+-----+-----+-----
time    |   -.000284   .0000893    -3.18  0.002   -.0004629   -.0001051
dprodl  |   .4548905   .1054274    4.31  0.000   .2437759   .6660052
_cons   |   .0159993   .0042484    3.77  0.000   .0074921   .0245065

```

67

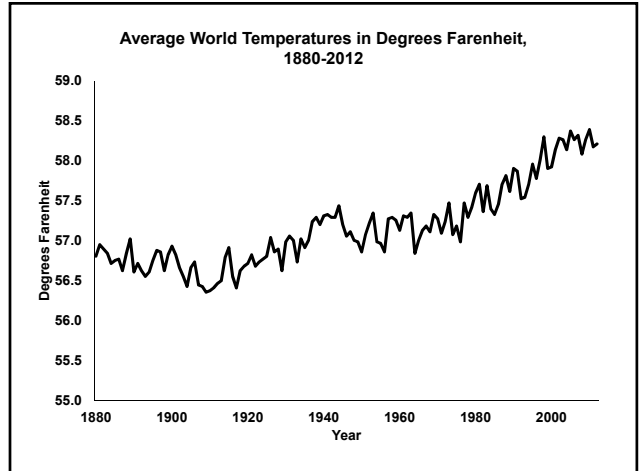


68



Global Warming

71



Two variables:

year
avg_temp_f

```
* set the time series
tsset year

* generate some time trends
gen time=year-1979
gen time2=time*time
gen time3=time2*time/1000
gen time4=time3*time/1000000
gen time5=time4*time/100000000

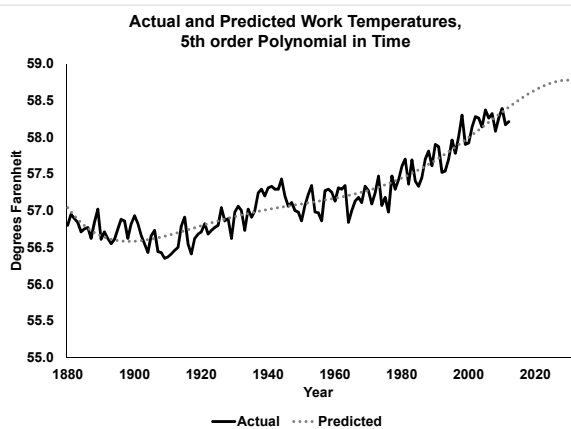
* run high-order polynomial
reg avg_temp_f time*
```

73

Source	SS	df	MS	Number of obs	=	133
Model	31.4589038	5	6.29178077	F(5, 127)	=	197.40
Residual	4.04789829	127	.031873215	Prob > F	=	0.0000
				R-squared	=	0.8860
				Adj R-squared	=	0.8815
Total	35.5068021	132	.268990925	Root MSE	=	.17853

avg_temp_f	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
time	.0195374	.0022199	8.80	0.000	.0151447 .0239301
time2	.0003794	.0000608	6.24	0.000	.000259 .0004998
time3	.0016513	.002339	0.71	0.481	-.0029771 .0062797
time4	-84.83526	52.42066	-1.62	0.108	-188.5663 18.89576
time5	-7.97e+07	3.11e+07	-2.56	0.012	-1.41e+08 -1.81e+07
_cons	57.42266	.0309648	1854.45	0.000	57.36139 57.48394

74



```
dfuller avg_temp_f
Dickey-Fuller test for unit root          Number of obs = 132
----- Interpolated Dickey-Fuller -----
Test Statistic      1% Critical      5% Critical      10% Critical
                   Value              Value              Value
Z(t)                 -1.652             -3.499             -2.888             -2.578
MacKinnon approximate p-value for Z(t) = 0.4558
```

76

```

* get lag
gen avg_temp_f_l=avg_temp_f[_n-1]

* get 1st different
gen d_avg_temp_f=avg_temp_f-avg_temp_f_l

* get dickey_fuller of 1st difference
dfuller d_avg_temp_f

* run regression
reg d_avg_temp_f time

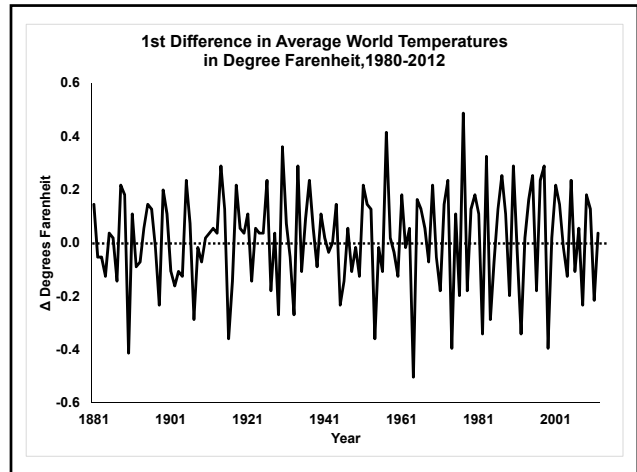
. * get dickey_fuller of 1st difference
. dfuller d_avg_temp_f

Dickey-Fuller test for unit root           Number of obs   =       131

----- Interpolated Dickey-Fuller -----
      Test          1% Critical   5% Critical   10% Critical
Statistic          Value         Value         Value
-----
Z(t)             -15.129         -3.500         -2.888         -2.578

```

77



Question

What does it mean if temperatures
follow a random walk?