

# Time-Series Modelling

ECON20222 - Lecture 9

Ralf Becker and Martyn Andrews

April 2024

# Aim for today

- Understand the basic features of time-series data
- Understand autocorrelation
- Understand the difference between stationary and non-stationary data and the different consequences of dealing with these
- Understand how to build dynamic models that can be used for forecasting

# Purpose of time-series modelling (TS modelling)

There are typically three things econometricians want to achieve with time-series modelling

- 1 **Establishing causal relationships** between time-series.

This is very difficult and in general causal relationships are more difficult to establish with TS modelling. It is not impossible but getting convincing exogenous variation is difficult.

- 2 Understanding the **dynamics in relationships between variables**.

Questions like, “If the Central Bank changes the base rate, how long will it take for this to carry through to mortgage rates?” This is perfectly possible as long as we don’t make strong causal statements (the CB may change base rates because mortgage rates are very low!!!)

- 3 **Forecasting one or several time series**.

This is possibly the most common purpose of TS modelling. We will focus on this.

# Import some data into R

```
rGDP <- pdfetch_ONS("ABMI", "UKEA")
periodicity(rGDP)    # check data frequency
names(rGDP) <- "real GDP" # give a sensible name

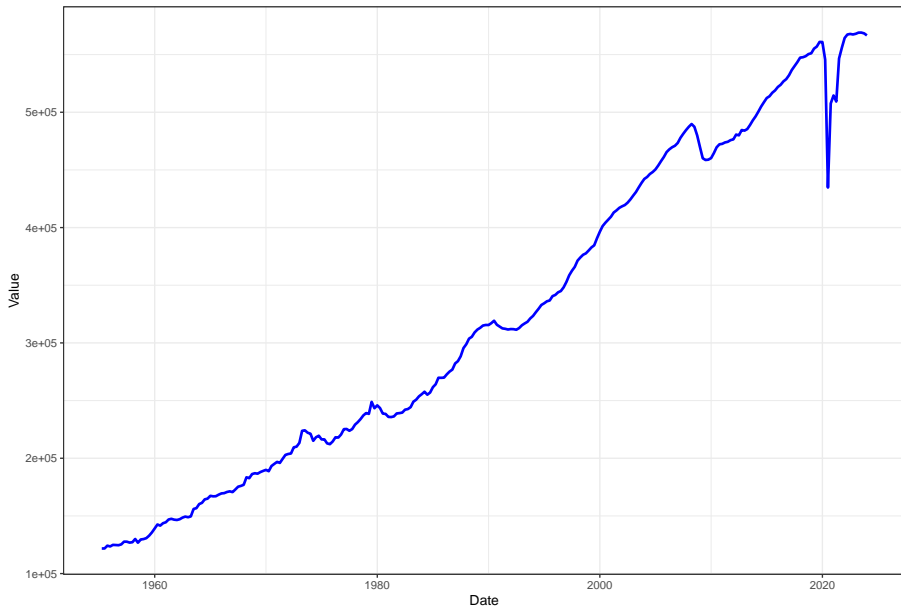
# keep all the data including 2023-Q4
# this was the last observation available at the time this was
# remove this line if you want to use updated data
rGDP <- rGDP["/2023-12"]
```

pdfetch functions allow you to directly tap a number of large data depositories:

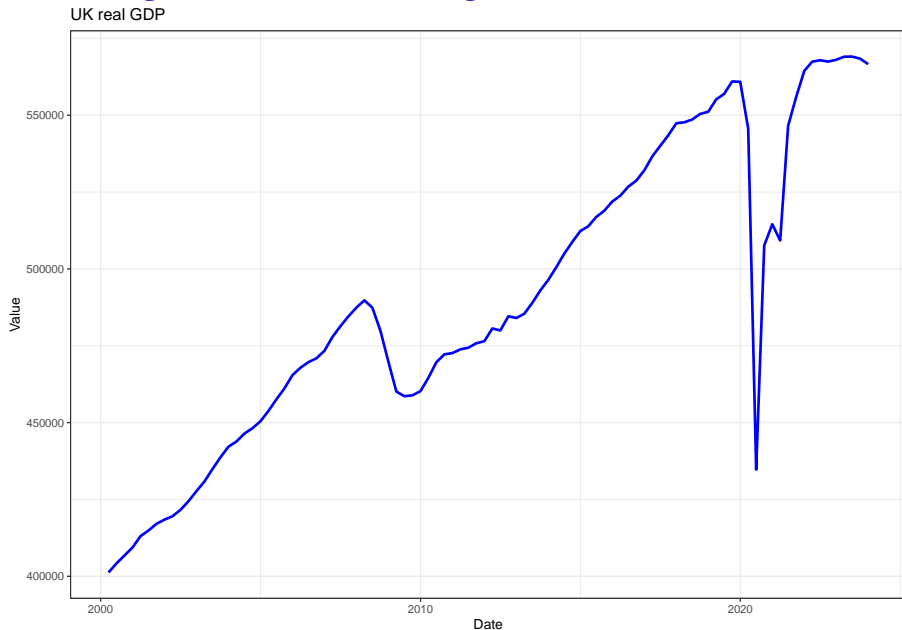
- Bundesbank
- Office for National Statistics (ONS)
- Eurostats
- FRED, etc

# An example

UK real GDP



# An example - focus on later periods



# The autocorrelation function (ACF)

Recall the correlation coefficient between two random variables  $z_i$  and  $p_i$  for two cross-sectional variables (index  $i$ )

$$\text{Corr}(z_i, p_i) = \frac{\text{Cov}(z_i, p_i)}{\text{sd}(z_i)\text{sd}(p_i)}$$

- a measure that expresses the strength of relationship between  $z_i$  and  $p_i$
- it takes values in the interval  $[-1, 1]$

# The autocorrelation function (ACF)

Consider the  $y_t = rGDP_t$  time series

- we now use the subscript  $t$
- the subscript goes from  $t = 1, \dots, T$  where  $T$  indicates how many observations we have
- here observations are quarterly (but other series can have other frequencies: e.g. annual, monthly, weekly, daily, hourly, etc. )
- here observations are from Q1 1955 to Q4 2023 (276 observations)

The ACF expresses how observations are correlated to observations 1, 2, 3 or  $k$  observations prior.

How can you calculate a correlation of a series with itself?



## The autocorrelation function (ACF)

Let's consider the time series  $y_t$  and the series one period prior,  $y_{t-1}$ . We also call  $y_{t-1}$  a one period lag of  $y_t$ .

Observation	$y_t$	$y_{t-1}$
1	$y_{1955Q1}$	NA
2	$y_{1955Q2}$	$y_{1955Q1}$
3	$y_{1955Q3}$	$y_{1955Q2}$
4	$y_{1955Q4}$	$y_{1955Q3}$
$\vdots$	$\vdots$	$\vdots$
275	$y_{2023Q3}$	$y_{2023Q2}$
276	$y_{2023Q4}$	$y_{2023Q3}$

Now we have “two” series for which we can calculate a correlation coefficient. We call this the first order autocorrelation coefficient  $\rho_1$ .

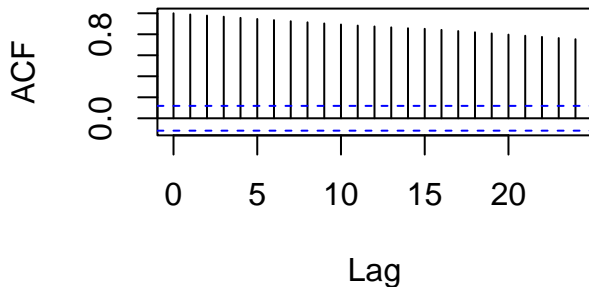
An ACF is a collection of autocorrelation coefficients calculated for longer lags  $k$ ,  $\rho_k$ .

## The autocorrelation function (ACF)

In R this ACF is easily calculated using the `acf` function.

```
temp_acf <- acf(rGDP)
```

### Series rGDP

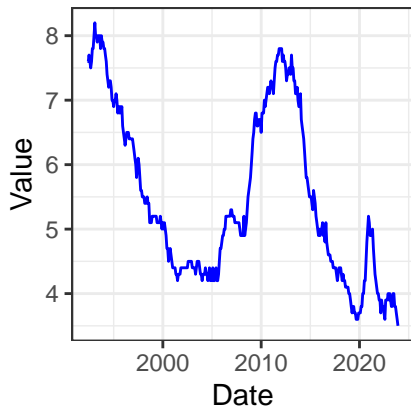


# The autocorrelation function (ACF)

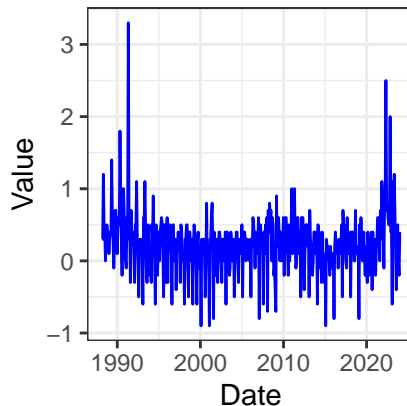
- The GDP series is strongly upward trending for most of the time
- this is common for many macroeconomic series
- this results in an ACF which has large  $\rho_k$  for fairly large values of  $k$  ( $\rho_8 = 0.914$ )
- we call this a persistent series

## Two further examples

### Female Unemployment



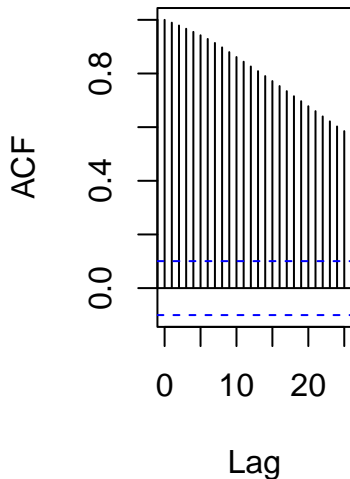
### Monthly Inflation Rate



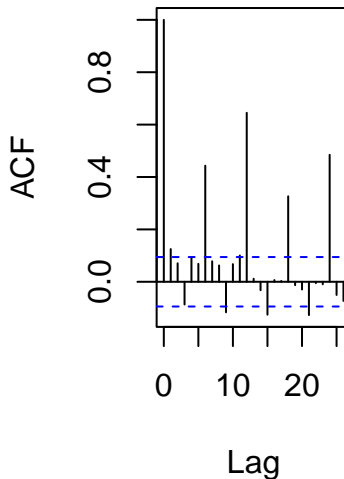
These are data at a monthly frequency

# The ACF

## Unemployment Rate



## Inflation



# The ACF

- the ACF shows that the unemployment rate is also very persistent and the ACF only slowly converges to 0.  $\Rightarrow$  a series can be persistent without time trend
- The inflation rate is not persistent and the autocorrelation quickly drops towards 0.
- But there are peaks of autocorrelation at frequencies 6 and 12 indicating some seasonal variation.

The ACF tells us something about how informative today's observation is for that in 1, 2, 3 or  $k$  periods ahead.

- If the ACF decays quickly to 0 then today's info is not very valuable for forecasting long into the future
- If the ACF decays slowly then today's info is valuable for future observations

# Stationary and Nonstationary Series

The ACF expresses how persistent a series is.

- A series that is extremely persistent is called a **nonstationary** series.
- A series that is not very persistent is called a **stationary** series.

Here: rGDP and unemployment rate are nonstationary. Inflation rate is stationary.

- In general series with a time-trend are nonstationary
- Some series without time trend are also nonstationary (e.g. female unemployment rate)
- **BUT** there is a huge grey area inbetween.

Formal statistical tests exist (e.g. Augmented Dickey-Fuller test) to decide (but they can be contradictory) and are not dealt with here. Here we eye-ball the series and look at how slowly the ACF converges to 0.

# Transformations

An important time-series transformation we consider is that of differencing a series.

Observation	$y_t$	$y_{t-1}$	$\Delta y_t$
2	$y_{1955Q2}$	$y_{1955Q1}$	$y_{1955Q2} - y_{1955Q1} = \Delta y_{1955Q2}$
3	$y_{1955Q3}$	$y_{1955Q2}$	$y_{1955Q3} - y_{1955Q2} = \Delta y_{1955Q3}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$

Often we are actually much more interested in the difference of a series rather than the level. GDP is a case in point, the growth rate is what we are really interested in!

The GDP growth rate can be approximated for small growth rates by (assuming that  $y_t$  is the GDP series)

$$gGDP_t = \frac{y_t - y_{t-1}}{y_{t-1}} \text{ or} \quad (1)$$

$$gGDP_t = \ln(y_t) - \ln(y_{t-1}) \quad (2)$$



# ACF of differenced series

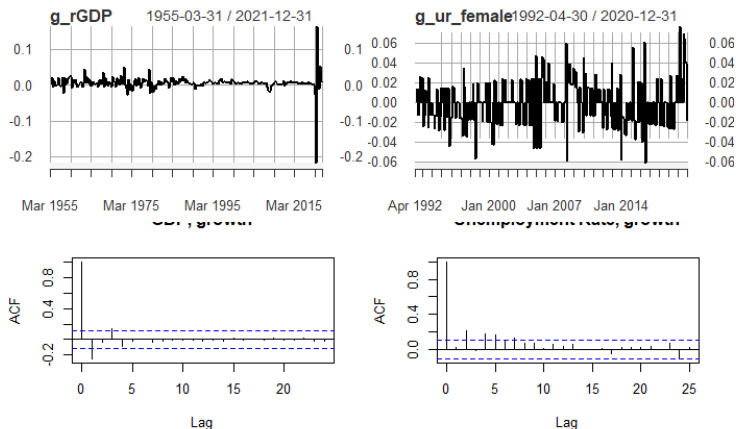


Figure 1: TS and ACF plots for growth in GDP and unemployment rate, ONS

Differencing can turn nonstationary series into a stationary series.

# A simple regression

## Cross-Section Data

$$y_t = \alpha + \beta x_t + u_t \quad (3)$$

$$E(u) = 0 \quad (4)$$

$$E(u|x) = 0 \quad (5)$$

this implied that  $x_t$  was exogenous.

## Time-Series Data

$$ur_t = \alpha + \beta rGDP_t + u_t \quad (6)$$

# A simple regression

let's look at the end of the data table

```
# we multiply by 100 to express in percentage points, i.e. 0.5 is 0.5%
reg_data$d_lgdp <- 100*diff(log(reg_data$real.GDP))
reg_data$d_lur <- 100*diff(log(reg_data$ur_female.Close))
tail(reg_data,10)
```

	real.GDP	ur_female.Close	d_lgdp	d_lur
2021-06-30	546579	4.5	7.07181982	-8.515781
2021-09-30	555956	4.1	1.70103010	-9.309042
2021-12-31	564407	3.9	1.50864680	-5.001042
2022-03-31	567396	3.9	0.52818502	0.000000
2022-06-30	567889	3.7	0.08685044	-5.264373
2022-09-30	567445	3.9	-0.07821487	5.264373
2022-12-31	568034	3.9	0.10374477	0.000000
2023-03-31	569027	3.8	0.17466086	-2.597549
2023-06-30	569076	4.0	0.00861082	5.129329
2023-09-30	568397	3.7	-0.11938746	-7.796154

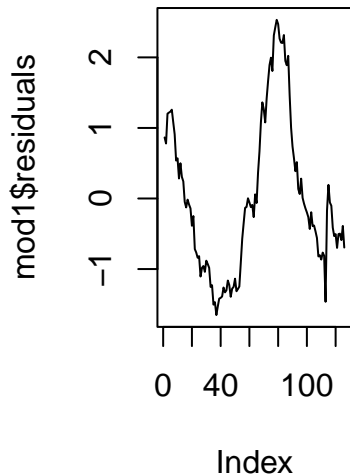
# A simple regression

```
mod1 <- lm(ur_female.Close~real.GDP,data = reg_data)
stargazer_HC(mod1)
```

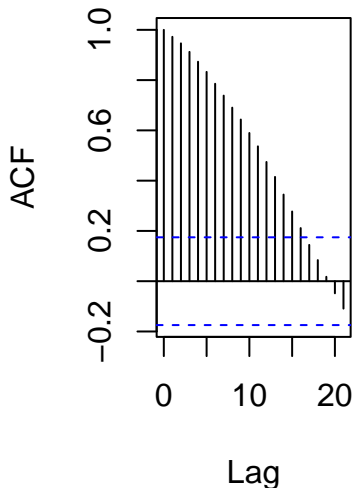
```
=====
                        Dependent variable:
-----
                        ur_female.Close
-----
real.GDP                -0.00001***
                        (0.00000)
Constant                9.791***
                        (0.438)
-----
Observations            126
R2                      0.297
Adjusted R2             0.291
Residual Std. Error    1.102 (df = 124)
F Statistic             52.302*** (df = 1; 124)
=====
Note:                   *p<0.1; **p<0.05; ***p<0.01
                        Robust standard errors in parenthesis
```

# A simple regression - residual autocorrelation

**mod1 – Residuals**



**Series mod1\$residua**



# A simple regression - testing for residual autocorrelation

We can again apply a hypothesis test the Breusch-Godfrey test (`bgtest`). The null hypothesis is that there is no autocorrelation.

```
bgtest(mod1,order=4)
```

Breusch-Godfrey test for serial correlation of order up to

```
data:  mod1
```

```
LM test = 119.76, df = 4, p-value < 2.2e-16
```

## A simple regression - HAC standard errors

When we estimate a regression which has autocorrelated error terms we need to apply a different formula to calculate standard errors for coefficients in a regression model (autoregressive heteroscedasticity consistent - HAC).

- They are called Newey-West standard errors.
- They are implemented in `stargazer_HAC.r`.
- This will not change the coefficient estimates.
- Will change the standard errors to the coefficients and hence inference (which will be incorrect if you don't use them).
- If you have time-series data, in doubt, use Newey-West standard errors
- But crucial problems remain (see next slides)

In this example the standard errors only change marginally and hence are not shown here.

# Spurious Regression

We will explore what can happen if we run a regression involving nonstationary variables.

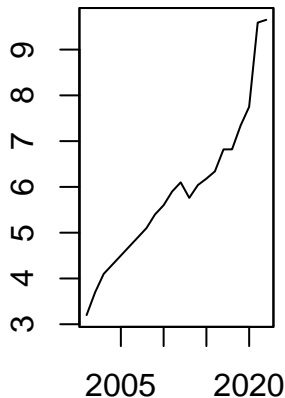
Let's get some datasets from EUROSTAT.

- % of agricultural area Total fully converted and under conversion to organic farming in Germany [**Org**]
- Thousands of passengers travelling to and from Norway by boat [**Pass**]
- % of population with tertiary education in Italy [**Tert**]
- Primary Energy Consumption, Million tons of oil equivalent, in Poland [**Ene**]

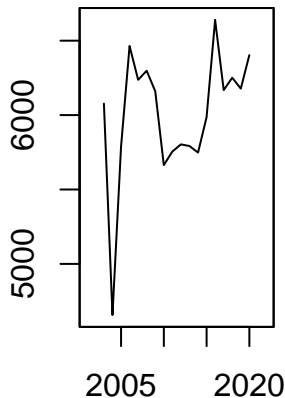


## A set of time-series

**Organic Farming  
GER**

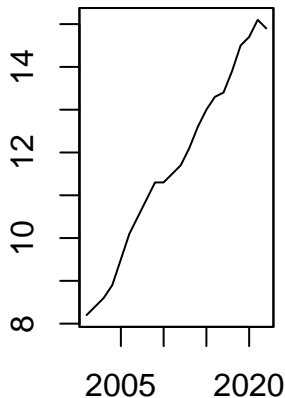


**Sea Passengers  
NOR**

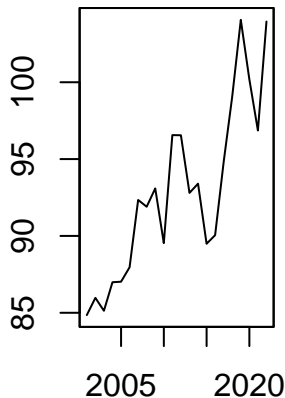


## A set of time-series

**Tertiary Education  
ITA**



**Energy Consumption  
POL**



# Spurious Regression - An example

$$Alc_t = \alpha + \gamma Org_t + u_t \quad (7)$$

```
mod_sr <- lm(Ene.Cons.PO~Organic.Farming.GE, data = data_sr)
stargazer_HAC(mod_sr)
```

```
=====
                        Dependent variable:
                        -----
                        Ene.Cons.PO
                        -----
Organic.Farming.GE      2.879***
                        (0.412)
Constant                75.858***
                        (2.521)
=====
Observations            22
R2                      0.710
Adjusted R2             0.695
Residual Std. Error     3.151 (df = 20)
F Statistic             48.874*** (df = 1; 20)
=====
Note:                   *p<0.1; **p<0.05; ***p<0.01
                        Robust standard errors in parenthesis
```

# Spurious Regression

All possible combinations of simple regressions between the four variables.

Table 1: Regression statistics

Dep.Var	Exp. Var			
	Org	Pass	Tert	Ene
Org	$\hat{\gamma}_1$	0.001	0.714***	0.246***
	$se_{\hat{\gamma}_1}$	(0.001)	(0.057)	(0.035)
	$R^2$	0.143	0.888	0.710
Pass	$\hat{\gamma}_1$	159.463	106.734*	23.876
	$se_{\hat{\gamma}_1}$	(97.459)	(53.471)	(20.838)
	$R^2$	0.143	0.199	0.076
Tert	$\hat{\gamma}_1$	1.244***	0.002*	0.329***
	$se_{\hat{\gamma}_1}$	(0.099)	(0.001)	(0.045)
	$R^2$	0.888	0.199	0.725
Ene	$\hat{\gamma}_1$	2.879***	0.003	2.205***
	$se_{\hat{\gamma}_1}$	(0.412)	(0.003)	(0.304)
	$R^2$	0.710	0.076	0.725

Note: \* $p < 0.1$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$ ,  
Newey-West standard errors in parenthesis

# (Unmasking of a) Spurious Regression

One way how you can unmask the spuriousness, if both series are trending is to include a time trend

```
mod_sr2 <- lm(Ene.Cons.P0~Organic.Farming.GE+index(data_sr),  
              data = data_sr)
```

or estimate a model in the differences of variables rather than the levels

```
mod_sr3 <- lm(diff(Ene.Cons.P0)~diff(Organic.Farming.GE),  
              data = data_sr)
```

# (Unmasking of a) Spurious Regression

```
stargazer_HAC(mod_sr,mod_sr2,mod_sr3, type_out = "text", omit.stat = "f")
```

Dependent variable:			
	Ene.Cons.PO (1)	(2)	diff(Ene.Cons.PO) (3)
Organic.Farming.GE	2.879*** (0.482)	1.145 (1.872)	
index(data_sr)		0.001 (0.001)	
diff(Organic.Farming.GE)			-1.242 (1.813)
Constant	75.858*** (2.518)	66.628*** (10.384)	1.291 (1.064)
Observations	22	22	21
R2	0.710	0.737	0.020
Adjusted R2	0.695	0.709	-0.031
Residual Std. Error	3.151 (df = 20)	3.077 (df = 19)	3.547 (df = 19)
Note:			
*p<0.1; **p<0.05; ***p<0.01			
Newey-West standard errors in parenthesis			

# Spurious Regression - A summary

If you estimate a regression between two nonstationary series:

- Do not mistake very significant coefficients or large  $R^2$  values for evidence of a substantial link between two series
- when regressing nonstationary series (in particular but not only when the series have a time-trend) will very likely deliver a spurious correlation **beyond** the common time trend
- When series are nonstationary but don't have a time-trend then you should consider estimating a regression in differences (recall differencing can turn nonstationary series into stationary ones)

For all these reasons, when considering time-series data, we need to **use stationary data**. **If not estimated coefficients will be, in general, neither unbiased nor consistent** and cannot be interpreted.

## A simple regression - but better

From the discussion on spurious regressions we have learned that estimating a model in differences can protect you against spurious regressions.

$$\Delta ur_t = \alpha + \beta \Delta rGDP_t + u_t$$

where we use growth rates (or log differences)

```
# Data (real.GDP, ur_female.Close) are in reg_data  
# we multiply by 100 to express in percentage points,  
# i.e. 0.5 is 0.5% or 0.005  
reg_data$d_lgdp <- 100*diff(log(reg_data$real.GDP))  
reg_data$d_lur <- 100*diff(log(reg_data$ur_female.Close))  
mod4 <- lm(d_lur~d_lgdp,data = reg_data)
```



# A simple regression - but better

```
=====
                        Dependent variable:
                        -----
                                d_lur
                        -----
d_lgdp                                -0.214
                                      (0.142)

Constant                             -0.483
                                      (0.378)

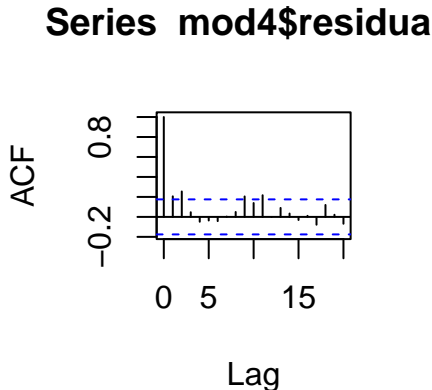
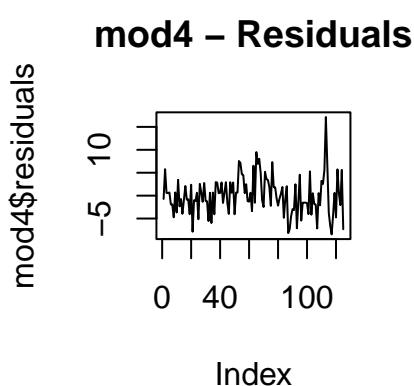
-----
Observations                          125
R2                                    0.018
Adjusted R2                           0.010
Residual Std. Error                   4.151 (df = 123)
F Statistic                           2.281 (df = 1; 123)

=====
Note:                                *p<0.1; **p<0.05; ***p<0.01
                                Robust standard errors in parenthesis
```

## A simple regression - but better

Let's have a look at the residuals.

```
par(mfrow=c(1,2))  
plot(mod4$residuals, type = "l", main = "mod4 - Residuals")  
acf(mod4$residuals)
```



## A simple regression - but better

- We can see that, at lag 2, there is a small amount of autocorrelation in the residuals.
- Breusch-Godfrey test: p-value of 0.0097
- What is the consequence? As residuals are stationary, HAC standard errors will deal with inference issues

# Adding dynamic effects

In the model we estimated:

$$\Delta ur_t = \alpha + \beta \Delta rGDP_t + u_t$$

all the action happened in one time period ( $t$ ).

We need to consider that effects in the economy may take some time. This fact will motivate two major generalisations.

- Include lags of the explanatory variable
- Include lags of the dependent variable

$$\begin{aligned} \Delta ur_t = & \alpha_0 + \alpha_1 \Delta ur_{t-1} + \alpha_2 \Delta ur_{t-2} + \dots + \alpha_p \Delta ur_{t-p} + \\ & \beta_0 \Delta rGDP_t + \beta_1 \Delta rGDP_{t-1} + \dots + \beta_k \Delta rGDP_{t-k} + u_t \end{aligned}$$

# Autoregressive Distributed Lag (ADL) models

```
mod5 <- lm(d_lur~lag(d_lur,1)+lag(d_lur,2)+  
           d_lgdp+lag(d_lgdp,1)+lag(d_lgdp,2),  
           data = reg_data)
```

Here we use the `lag(series, k)` function which calculates the `k` period lag of `series`.

# Autoregressive Distributed Lag (ARDL) models

```
stargazer_HAC(mod4,mod5, omit.stat = "f")
```

Dependent variable:		
	d_lur	
	(1)	(2)
lag(d_lur, 1)		0.019 (0.098)
lag(d_lur, 2)		0.218*** (0.068)
d_lgdp	-0.214 (0.239)	-0.433*** (0.127)
lag(d_lgdp, 1)		-0.718*** (0.074)
lag(d_lgdp, 2)		-0.558*** (0.072)
Constant	-0.483 (0.423)	0.338 (0.323)
Observations	125	123
R2	0.018	0.328
Adjusted R2	0.010	0.299
Residual Std. Error	4.151 (df = 123)	3.495 (df = 117)
Note: *p<0.1; **p<0.05; ***p<0.01		

# Forecasting models

Are the above models useful forecasting models?

$$\Delta ur_t = \alpha_0 + \alpha_1 \Delta ur_{t-1} + \alpha_2 \Delta ur_{t-2} + \dots + \alpha_p \Delta ur_{t-p} + \beta_0 \Delta rGDP_t + \beta_1 \Delta rGDP_{t-1} + \dots + \beta_k \Delta rGDP_{t-k} + u_t$$

- Say you estimated the above models using data up to Q3 2018.
- Hence you have estimated coefficients
- Could you use it to forecast the value in Q4 2018?

No, if we want to forecast  $\Delta ur_{2018Q4}$  we would need  $\Delta rGDP_{2018Q4}$ ! But we don't have that. We would first need a forecast for  $\Delta rGDP_{2018Q4}$  to then forecast  $\Delta ur_{2018Q4}$ . We call these conditional forecasts.

# Forecasting models

To build a useful forecasting model we remove all contemporaneous terms from the right hand side, such that we can produce forecasts for period  $t$  only having information at time (say)  $t - 1$ .

We remove the  $\beta_0 \Delta rGDP_t$  term from our model:

```
mod6 <- lm(d_lur~lag(d_lur,1)+lag(d_lur,2)+  
           lag(d_lgdp,1)+lag(d_lgdp,2),data = reg_data)
```



# Forecasting models

Dependent variable:			
	(1)	d_lur (2)	(3)
lag(d_lur, 1)		0.019 (0.098)	0.057 (0.097)
lag(d_lur, 2)		0.218*** (0.068)	0.230*** (0.072)
d_lgdp	-0.214 (0.239)	-0.433*** (0.127)	
lag(d_lgdp, 1)		-0.718*** (0.074)	-0.564*** (0.077)
lag(d_lgdp, 2)		-0.558*** (0.072)	-0.468*** (0.056)
Constant	-0.483 (0.423)	0.338 (0.323)	0.039 (0.340)
Observations	125	123	123
R2	0.018	0.328	0.261
Adjusted R2	0.010	0.299	0.236
Residual Std. Error	4.151 (df = 123)	3.495 (df = 117)	3.649 (df = 118)
Note:			
*p<0.1; **p<0.05; ***p<0.01			
Newey-West standard errors in parenthesis			

# Autoregressive models

- We are interested in forecasting the unemployment rate changes,  $\Delta ur_t$ , but we still need observations for  $rGDP_{t-1}$  and further lags.
- Realistically there may be other series we may want to consider: e.g. interest rate, inflation, wages, etc.
- Can we forecast  $\Delta ur_t$  with nothing else but the history of  $\Delta ur_t$ ?

Yes, we call these **autoregressive (AR) models**

$$\Delta ur_t = \alpha_0 + \alpha_1 \Delta ur_{t-1} + \alpha_2 \Delta ur_{t-2} + \dots + \alpha_p \Delta ur_{t-p} + u_t$$

We “merely” need to choose the lag length  $p$ !

# Autoregressive models

For starters we use  $p = 2$  as above.

```
mod7 <- lm(d_lur~lag(d_lur,1)+lag(d_lur,2),data = reg_data)
```

# Autoregressive models

Dependent variable:				
	d_lur			
	(1)	(2)	(3)	(4)
lag(d_lur, 1)		0.019 (0.098)	0.057 (0.097)	0.192* (0.114)
lag(d_lur, 2)		0.218*** (0.068)	0.230*** (0.072)	0.247** (0.100)
d_lgdp	-0.214 (0.239)	-0.433*** (0.127)		
lag(d_lgdp, 1)		-0.718*** (0.074)	-0.564*** (0.077)	
lag(d_lgdp, 2)		-0.558*** (0.072)	-0.468*** (0.056)	
Constant	-0.483 (0.423)	0.338 (0.323)	0.039 (0.340)	-0.385 (0.415)
Observations	125	123	123	123
R2	0.018	0.328	0.261	0.122
Adjusted R2	0.010	0.299	0.236	0.108
Residual Std. Error	4.151 (df = 123)	3.495 (df = 117)	3.649 (df = 118)	3.943 (df = 120)

Note: \*p<0.1; \*\*p<0.05; \*\*\*p<0.01  
Newey-West standard errors in parenthesis

# Lag Length - Information Criterion

The last piece in the puzzle is how we determine the best lag length

- Every additional lag will improve the in-sample fit of your model  
⇒ Should we include as many lags as possible?
- Will “over-fit” the model (to features in the data which are not generic but specific to the in-sample period)  
⇒ Tends to make the forecasts more erratic/volatile.

This trade-off needs to be optimised. ⇒ **Information Criteria**

# Lag Length - Information Criterion

Trade-off between too many variables in model (bad) v in-sample fit (good). Information criteria quantify this trade-off.

One example **Akaike Information Criterion (AIC)**

In R:

Where

- `mod6` and `mod6_4` are the ADL models with 2 and 4 lags respectively
- `mod7` and `mod7_4` are the AR models with 2 and 4 lags respectively

# Lag Length - Information Criterion

Table 2: AIC for ADL and AR with 2 and 4 lags

Model	N. of para	AIC
ADL (2 lags)	6	674.3955
AR (2 lags)	4	691.5328
ADL (4 lags)	10	669.0942
AR (4 lags)	6	682.9785

*Note:* Number of parameters includes  
linear parameters and residual variance

The optimal model (as per the trade off in the AIC criterion) is the model with the **lowest AIC**.

Here ADL(4).

# Summary

We learned that

- the ACF encapsulates how persistent a time-series is
- Time-series which are very persistent are called nonstationary
- Using nonstationary series in simple regressions can lead to very misleading (spurious) results
- Estimating models with either a time-trend or in differences can protect you against misleading results
- To build forecasting models we need to ensure that we use explanatory variables which are available at the time of forecasting
- AR models can be a convenient tool for forecasting
- Information criteria can help us to select the right model for forecasting