

# Longitudinal Analysis

*Richard White*

*2018-11-02*



# Contents

<b>1</b>	<b>Syllabus</b>	<b>5</b>
<b>2</b>	<b>Reference</b>	<b>7</b>
2.1	Scope of this course . . . . .	7
2.2	Introduction . . . . .	8
2.3	Method summary . . . . .	9
2.4	Identifying your scenario . . . . .	11
<b>3</b>	<b>Panel data: One area without autocorrelation</b>	<b>21</b>
3.1	Aim . . . . .	21
3.2	Creating the data . . . . .	22
3.3	True data . . . . .	23
3.4	Investigation . . . . .	24
3.5	Seasonality . . . . .	25
<b>4</b>	<b>Panel data: One area with autocorrelation</b>	<b>33</b>
4.1	Aim . . . . .	33
4.2	Creating the data . . . . .	34
4.3	Investigation . . . . .	35
4.4	Regressions . . . . .	37
4.5	Residual analysis . . . . .	39
4.6	(R ONLY) Regression with AR(1) correlation in residuals . . . . .	41
4.7	(STATA ONLY) Regression with robust standard errors . . . . .	45
<b>5</b>	<b>Not panel data: Multiple areas</b>	<b>47</b>
5.1	Aim . . . . .	47
5.2	Creating the data . . . . .	48
5.3	Investigating the data . . . . .	49
5.4	Regression . . . . .	50
<b>6</b>	<b>Panel data: multiple areas without autocorrelation</b>	<b>51</b>
6.1	Aim . . . . .	51
6.2	Creating the data . . . . .	52
6.3	Investigation . . . . .	53
6.4	Regression . . . . .	55
6.5	Residual analysis . . . . .	57
<b>7</b>	<b>Panel data: multiple areas with autocorrelation</b>	<b>61</b>
7.1	Aim . . . . .	61
7.2	Creating the data . . . . .	62
7.3	Investigation . . . . .	63
7.4	Regressions . . . . .	65
7.5	Residual analysis . . . . .	67

7.6	(R ONLY) Regression with AR(1) correlation in residuals . . . . .	69
7.7	Residual analysis . . . . .	70
7.8	(STATA ONLY) Regression with robust standard errors . . . . .	73
<b>8</b>	<b>Exercises</b>	<b>75</b>
8.1	Exercise 1 . . . . .	75
8.2	Exercise 2 . . . . .	77
8.3	Exercise 3 . . . . .	79
<b>9</b>	<b>Solutions</b>	<b>81</b>
9.1	Exercise 1 . . . . .	81
9.2	Exercise 2 . . . . .	84
9.3	Exercise 3 . . . . .	90

# Chapter 1

## Syllabus

**Instructor:** Richard White [richard.white@fhi.no]

**Time:** 09:30 - 15:00, 18th September 2017

**Location:** Main auditorium, L8, Lindern Campus, Folkehelseinstituttet, Oslo

**Language:** English

### Format and Procedures

09:00 - 10:00: Lecture 1

10:00 - 10:10: Break

10:10 - 11:10: Lecture 2

10:10 - 10:15: Break

11:15 - 11:45: Examples from FHI

### Description

This course will provide a basic overview of general statistical methodology that can be useful in the areas of infectious diseases, environmental medicine, and labwork. By the end of this course, students will be able to identify appropriate statistical methods for a variety of circumstances.

This course will **not** teach students how to implement these statistical methods, as there is not sufficient time. The aim of this course is to enable the student to identify which methods are required for their study, allowing the student to identify their needs for subsequent methods courses, self-learning, or external help.

You should register for this course if you are one of the following:

- Have experience with applying statistical methods, but are sometimes confused or uncertain as to whether or not you have selected the correct method.
- Do not have experience with applying statistical methods, and would like to get an overview over which methods are applicable for your projects so that you can then undertake further studies in these areas.

### Lecture 1

1. Identifying continuous, categorical, count, and censored variables
2. Identifying exposure and outcome variables
3. Identifying when t-tests (paired and unpaired) should be used
4. Identifying when non-parametric t-test equivalents should be used
5. Identifying when ANOVA should be used
6. Identifying when linear regression should be used

7. Identifying the similarities between t-tests, ANOVA, and regression
8. Identifying when logistic regression models should be used
9. Identifying when Poisson/negative binomial and cox regression models should be used
10. Identifying when chi-squared/fisher's exact test should be used

**Lecture 2**

1. Identifying when data does not have any dependencies (i.e. all observations are independent of each other) versus when data has complicated dependencies (i.e. longitudinal data, matched data, multiple cohorts)
2. Identifying when mixed effects regression models should be used
3. Identifying when conditional logistic regression models should be used
4. (TBD) Understanding the different imputation methods used when lab data is below the limit of detection (LOD)
5. (TBD) Understanding the best practices for data files and project folders

**Prerequisites**

To participate in this course it is recommended that you have some experience with either research or data.

**Additional information**

For the last 30 minutes of the course we will be going through examples of analyses performed at FHI and identifying which statistical methods are appropriate. If you would like your analysis to be featured/included in this section, please send an email to [richard.white@fhi.no](mailto:richard.white@fhi.no) briefly describing your problem.

# Chapter 2

## Reference

### 2.1 Scope of this course

When dealing with longitudinal data, there are two kinds of analyses that can be performed.

“Time series” analyses generally deal with one variable. The aim is to then predict the future only using the previous observations. A common example would be to predict tomorrow’s temperature, using today’s and yesterday’s temperature as exposures. We will not be focusing on these kinds of analyses in this course.

“Regression analyses” are very similar to ordinary regressions that you have been working with for many years. The only difference is that they have more advanced data structures that your current methods cannot handle. For example, if you want to see how the number of tuberculosis patients (outcome) is affected by the number of immigrants to Norway (exposure) over a 20 year period, then the number of patients in each year might be associated with each other, which might break assumptions of the regression models that you normally use (independent residuals). To account for the advanced structure of the data (correlation between different years) we will use more advanced regression techniques. This is what we will be focusing on in this course.

To recap: this course will let you run “normal regressions” in situations where the data structure would ordinarily prohibit you from running regression models. These situations mostly pertain to clusters of correlated data.

## 2.2 Introduction

There are two important definitions in this course:

- Panel data
- Autocorrelation

Panel data is a set of data with measurements repeated at equally spaced points. For example, weight data recorded every day, or every week, or every year would be considered panel data. A person who records three weight measurements randomly in 2018 would not be considered panel data.

When you have panel data, autocorrelation is the correlation between subsequent observations. For example, if you have daily observations, then the 1 day autocorrelation is the correlation between observations 1 day apart, and likewise the 2 day autocorrelation is the correlation between observations 2 days apart.

In this course we will consider 5 scenarios where we have multiple observations for each geographical area:

- Panel data: One geographical area, no autocorrelation
- Panel data: One geographical area, with autocorrelation
- Not panel data: Multiple geographical areas
- Panel data: Multiple geographical areas, no autocorrelation
- Panel data: Multiple geographical areas, with autocorrelation

Note, the following scenario can be covered by standard regression models:

- Multiple geographical areas, one time point/observation per geographical area



## 2.3 Method summary

### 2.3.1 Panel data: One geographical area, no autocorrelation

```
// STATA CODE
glm y yearminus2000 dailyrainfall cos365 sin365, family(poisson)

# R CODE
fit1 <- glm(y~yearMinus2000 + dailyrainfall + sin365 + cos365, data=d, family=poisson())
residuals(fit1, type = "response")
```

### 2.3.2 Panel data: One geographical area, with autocorrelation

```
// STATA CODE
glm y yearminus2000 cos365 sin365, family(poisson) vce(robust)

# R CODE
fit <- MASS::glmmPQL(y~yearMinus2000+sin365 + cos365, random = ~ 1 | ID,
                    family = poisson, data = d,
                    correlation=nlme::corAR1(form=~dayOfSeries|ID))
r <- residuals(fit1, type = "normalized")
pacf(r)
```

### 2.3.3 Not panel data: Multiple geographical areas

```
// STATA CODE
meglm y x yearMinus2000 || fylke:, family(poisson)

# R CODE
fit <- lme4::glmer(y~x + yearMinus2000 + (1|fylke),data=d,family=poisson())
```

### 2.3.4 Panel data: Multiple geographical areas, no autocorrelation

```
// STATA CODE
meglm y yearminus2000 cos365 sin365 || fylke:, family(poisson)

# R CODE
fit <- MASS::glmmPQL(y~yearMinus2000+sin365 + cos365, random = ~ 1 | fylke,
                    family = poisson, data = d)
r <- residuals(fit1, type = "normalized")
pacf(r)
```

### 2.3.5 Panel data: Multiple geographical areas, with autocorrelation

```
// STATA CODE
meglm y yearminus2000 cos365 sin365 || fylke:, family(poisson) vce(robust)

# R CODE
fit <- MASS::glmmPQL(y~yearMinus2000+sin365 + cos365, random = ~ 1 | fylke,
                    family = poisson, data = d,
                    correlation=nlme::corAR1(form=~dayOfSeries|fylke))
```

```
r <- residuals(fit1, type = "normalized")  
pacf(r)
```

## 2.4 Identifying your scenario

### 2.4.1 Step 1: Do you have panel data?

This step should be fairly simple. If your data has equally spaced time intervals between them, you have panel data.

### 2.4.2 Step 2: Do you have multiple geographical areas?

Again, fairly simple, just look at your data.

### 2.4.3 Step 3: Do you have autocorrelation?

Firstly, you must run a model pretending that you do not have autocorrelation.

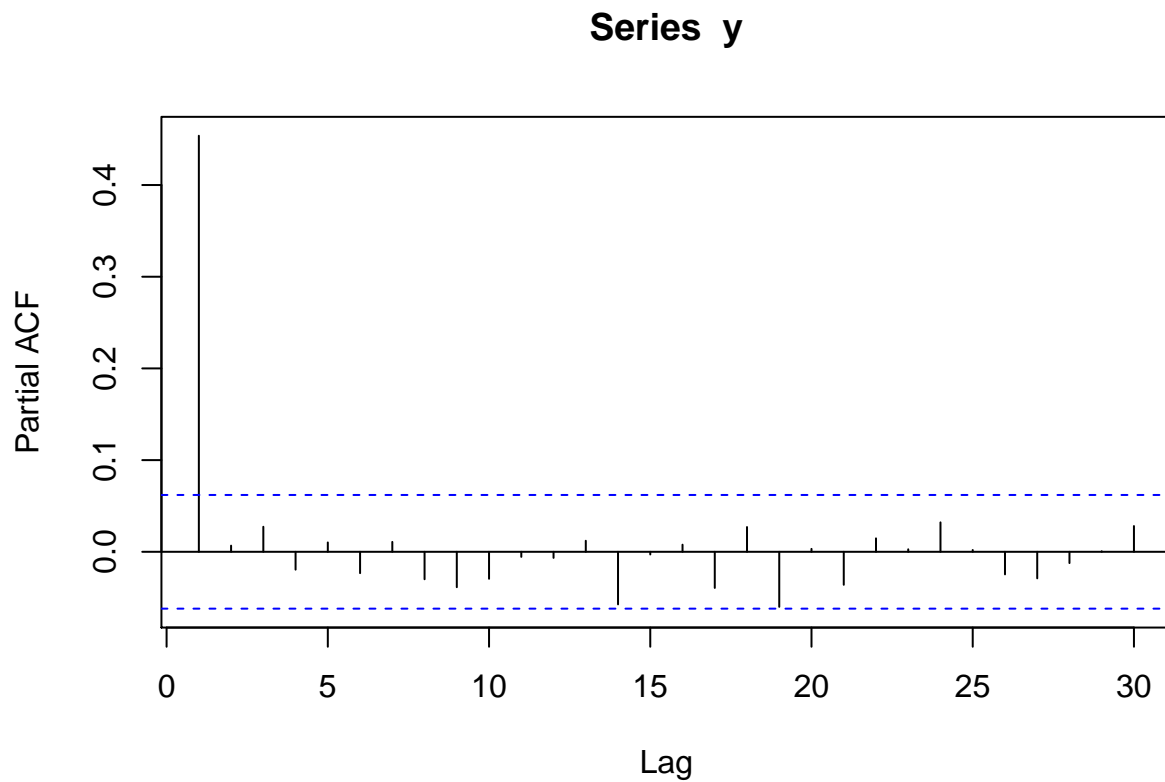
You then inspect the residuals from the model and see if autocorrelation exists. This is done with two statistical procedures: `pacf` (for **autoregressive models**, the most common type of autocorrelation), and `acf` (for **moving average models**, a less common type of autocorrelation).

#### 2.4.4 AR(1) data

```
y <- round(as.numeric(arima.sim(model=list("ar"=c(0.5)), rand.gen = rnorm, n=1000)))
```

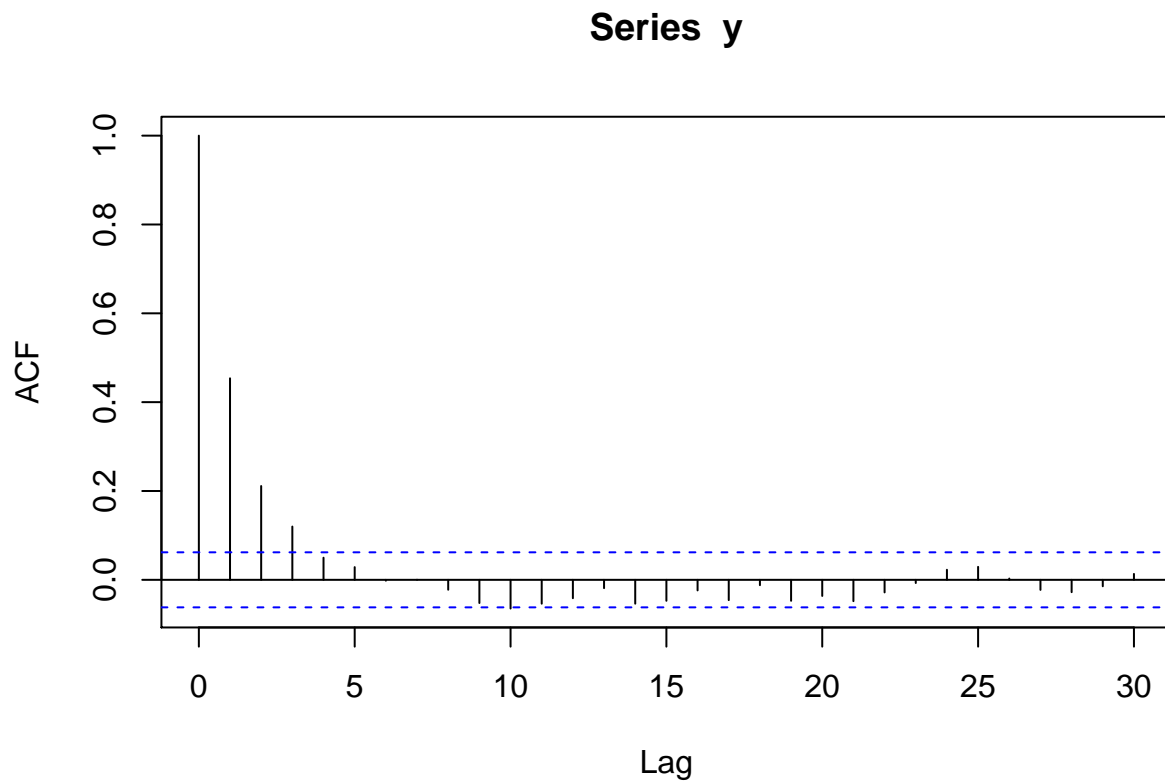
With autoregressive data, a `pacf` plot contains a number of sharp significant lines, indicating how many subsequent observations have autocorrelation. i.e. if one line is significant, it means that each observation is only correlated with its preceding observation (AR(1)). If two lines are significant, it means that each observation is correlated with its two preceding observations (AR(2)). The following plot represents AR(1) data.

```
pacf(y)
```



With autoregressive data, an `acf` plot contains a number of decreasing lines. The following `acf` plot represents some sort of `AR` data. Note that the `acf` plot displays `lag 0` (which is pointless and can be ignored), while the `pacf` plot does not.

```
acf(y)
```

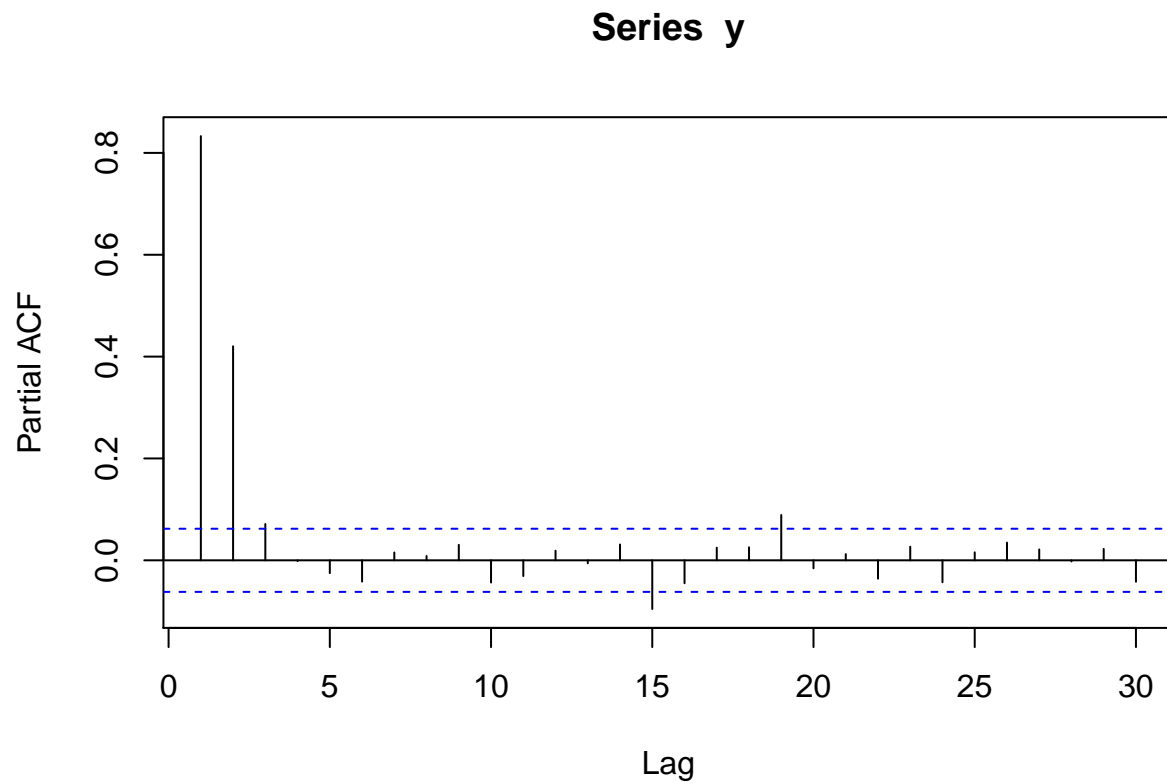


### 2.4.5 AR(2) data

```
y <- round(as.numeric(arima.sim(model=list("ar"=c(0.5,0.4)), rand.gen = rnorm, n=1000)))
```

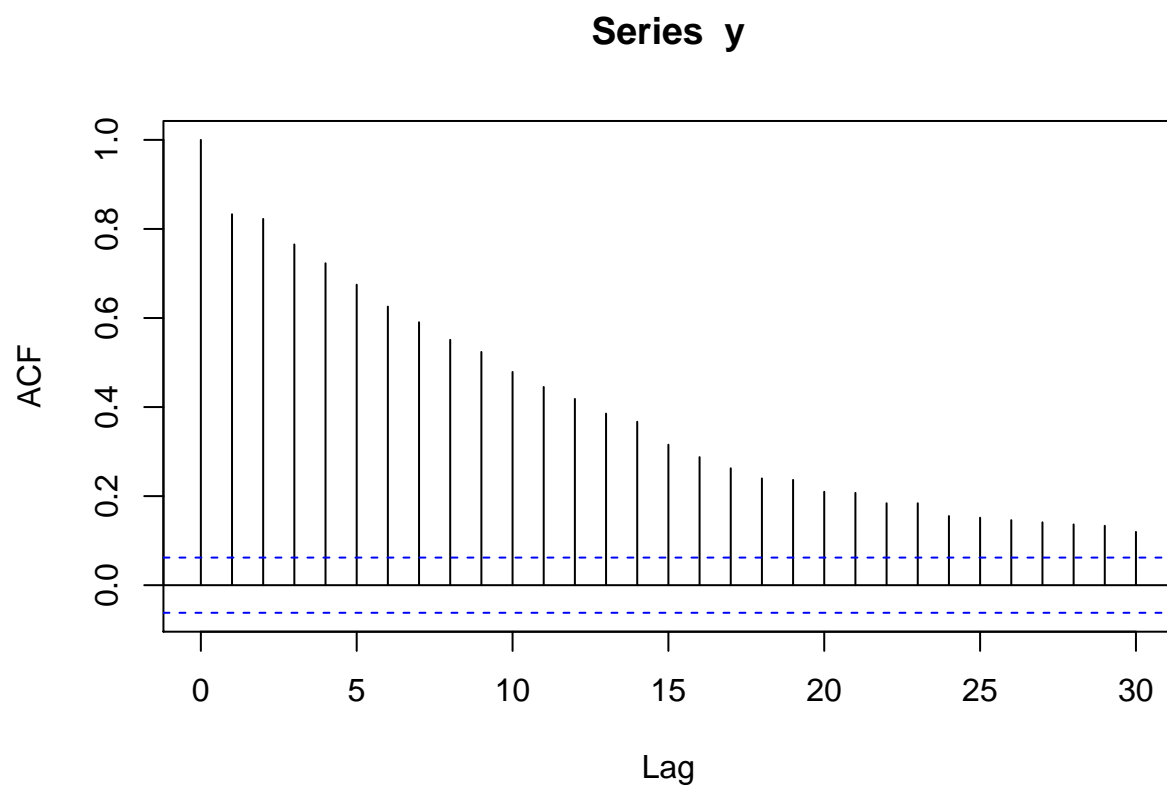
The following `pacf` plot represents AR(2) data. This means that each observation is correlated with its two preceding observations (AR(2)).

```
pacf(y)
```



The following `acf` plot represents some sort of AR data:

```
acf(y)
```

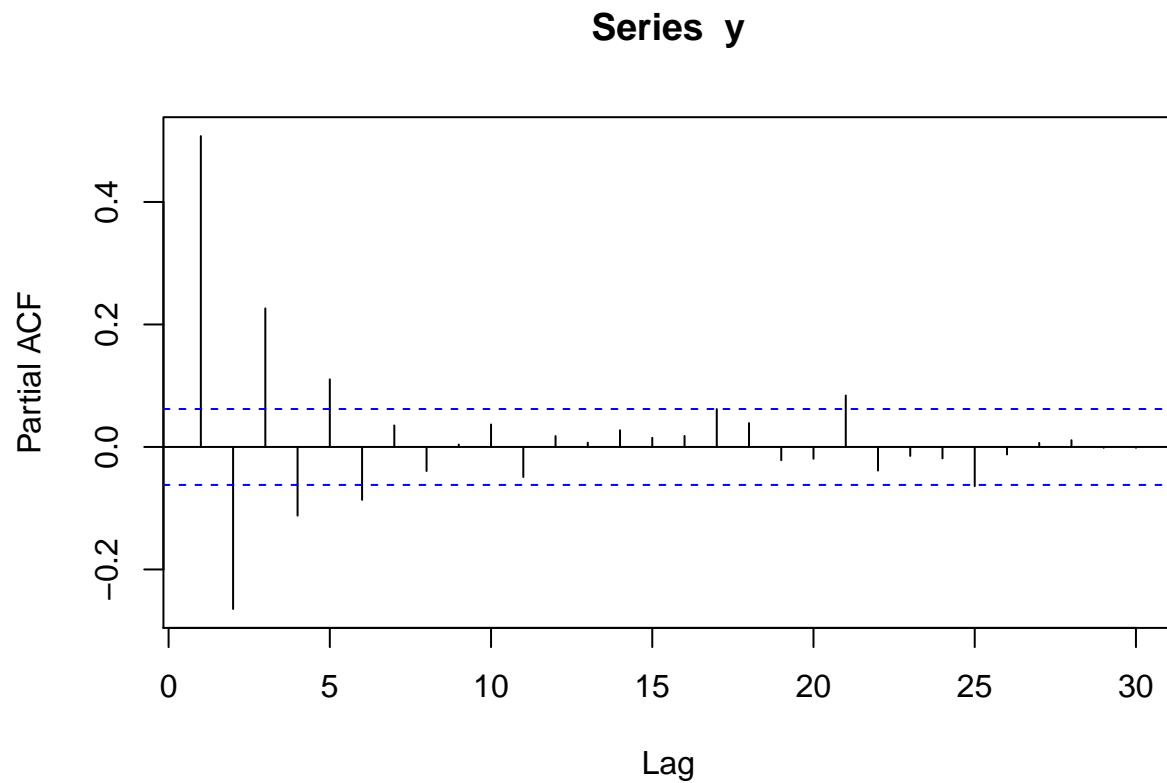


### 2.4.6 MA(1) data

```
y <- round(as.numeric(arima.sim(model=list("ma"=c(0.9)), rand.gen = rnorm, n=1000)))
```

With moving average data, a `pacf` plot contains a number of decreasing lines. The following `pacf` plot represents some sort of MA data.:

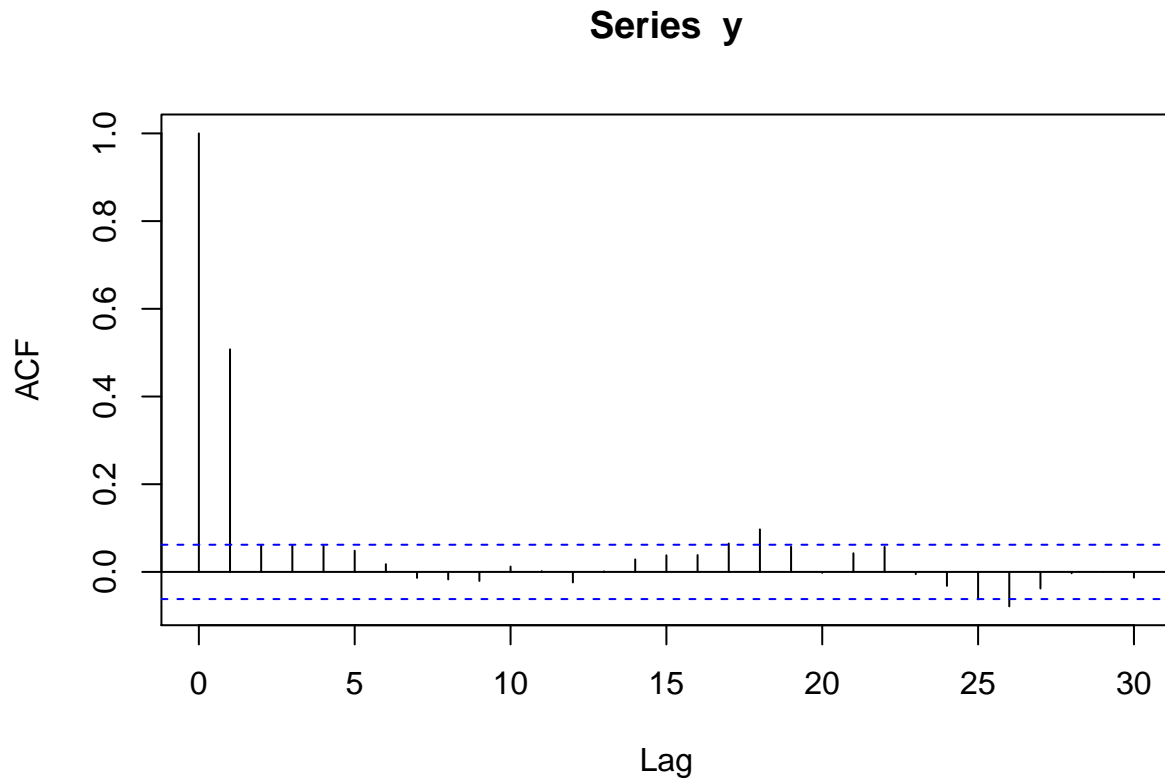
```
pacf(y)
```





With moving average data, an `acf` plot contains a number of sharp significant lines, demarking how many subsequent observations have autocorrelation. i.e. if one line is significant, it means that each observation is only correlated with its preceeding observation. If two lines are significant, it means that each observation is correlated with its two preceeding observations. The following plot represents `MA(1)` data. Note that the `acf` plot displays lag 0 (which is pointless and can be ignored), while the `pacf` plot does not.

```
acf(y)
```

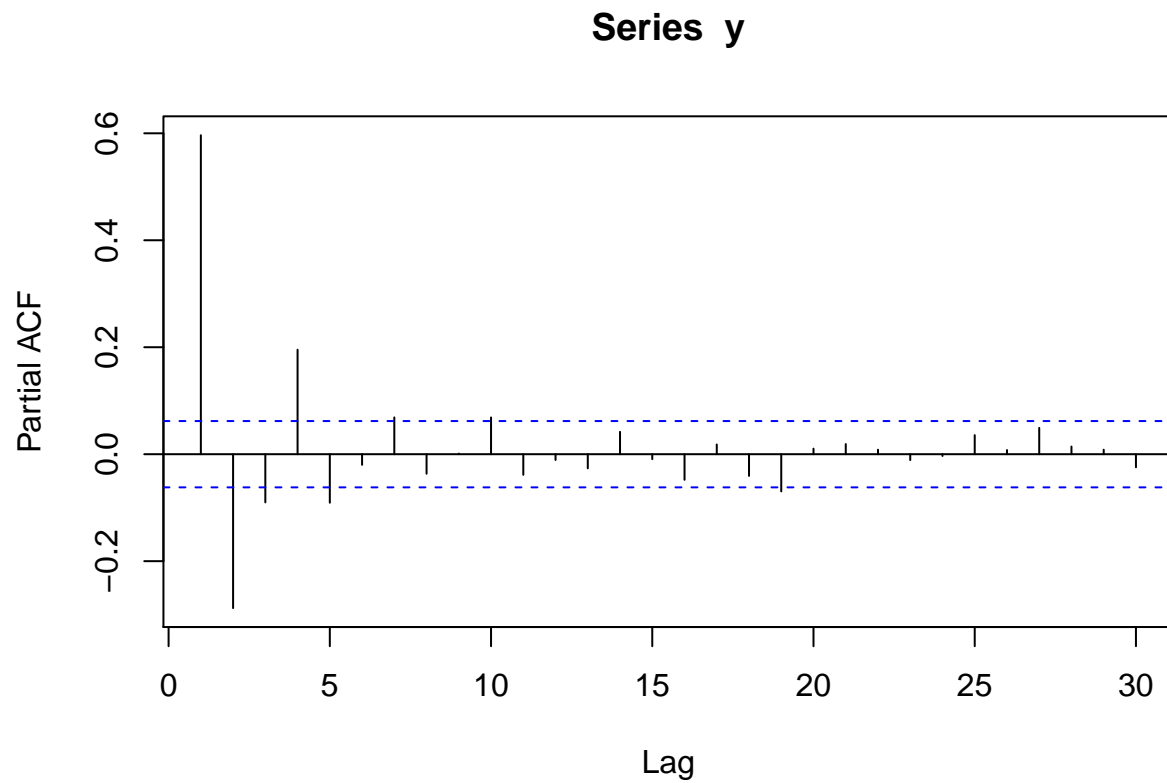


### 2.4.7 MA(2) data

```
y <- round(as.numeric(arima.sim(model=list("ma"=c(0.9,0.6)), rand.gen = rnorm, n=1000)))
```

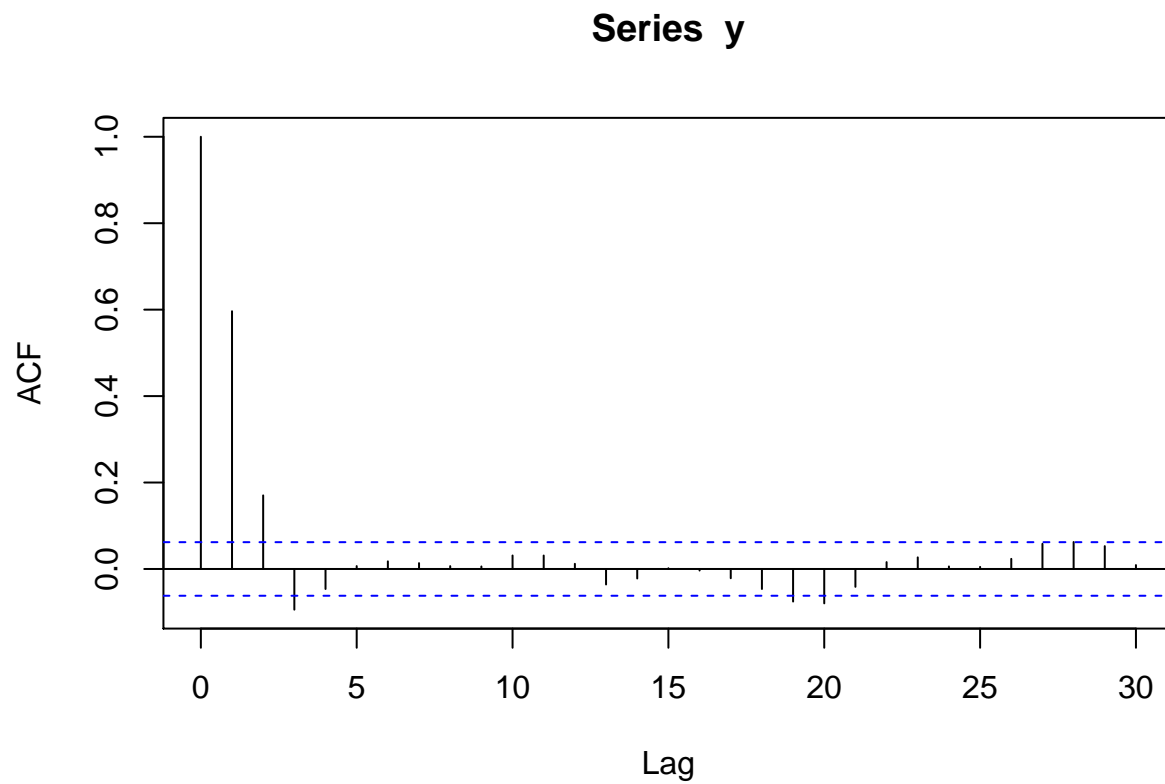
The following `pacf` plot represents some sort of MA data.

```
pacf(y)
```



The following `acf` plot represents `MA(2)` data. This means that each observation is correlated with its two preceding observations.

```
acf(y)
```





## Chapter 3

# Panel data: One area without autocorrelation

### 3.1 Aim

We are given a dataset containing daily counts of diseases from one geographical area. We want to identify:

- Does seasonality exist?
- If seasonality exists, when are the high/low seasons?
- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)
- Is daily rainfall associated with the number of cases?

## 3.2 Creating the data

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/chapter\\_3.csv](http://rwhite.no/longitudinal_analysis/data/chapter_3.csv)

*# R CODE*

```
dir.create("data")

library(data.table)
library(ggplot2)
set.seed(4)

AMPLITUDE <- 1.5
SEASONAL_HORIZONTAL_SHIFT <- 20

d <- data.table(date=seq.Date(
  from=as.Date("2000-01-01"),
  to=as.Date("2018-12-31"),
  by=1))
d[,year:=as.numeric(format.Date(date,"%G"))]
d[,week:=as.numeric(format.Date(date,"%V"))]
d[,month:=as.numeric(format.Date(date,"%m"))]
d[,yearMinus2000:=year-2000]
d[,dailyrainfall:=runif(.N, min=0, max=10)]

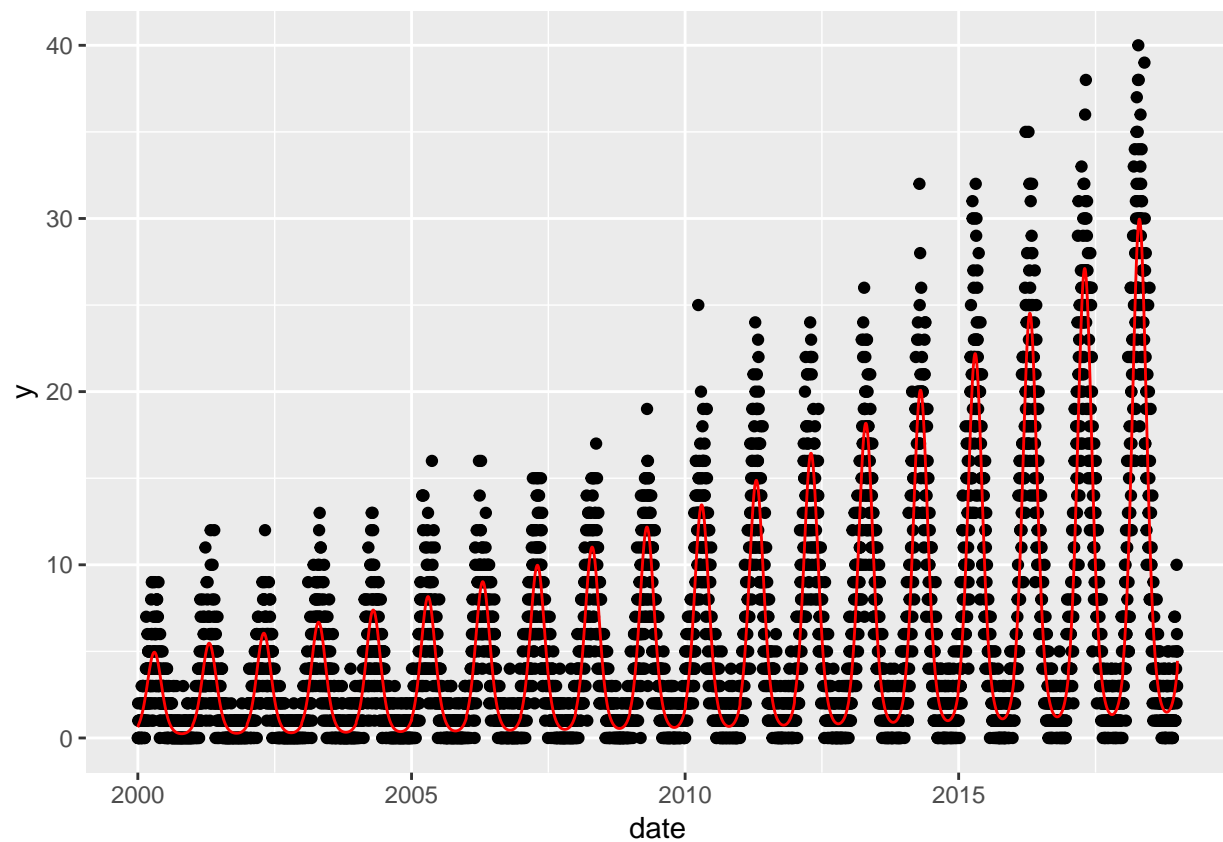
d[,dayOfYear:=as.numeric(format.Date(date,"%j"))]
d[,seasonalEffect:=sin(2*pi*(dayOfYear-SEASONAL_HORIZONTAL_SHIFT)/365)]
d[,mu := exp(0.1 + yearMinus2000*0.1 + seasonalEffect*AMPLITUDE)]
d[,y:=rpois(.N,mu)]

fwrite(d,"data/chapter_3.csv")
```

### 3.3 True data

Here we show the true data, and note that there is an increasing annual trend (the data gets higher as time goes on) and there is a seasonal pattern (one peak/trough per year)

```
q <- ggplot(d,aes(x=date))  
q <- q + geom_point(mapping=aes(y=y))  
q <- q + geom_line(mapping=aes(y=mu),colour="red")  
q
```

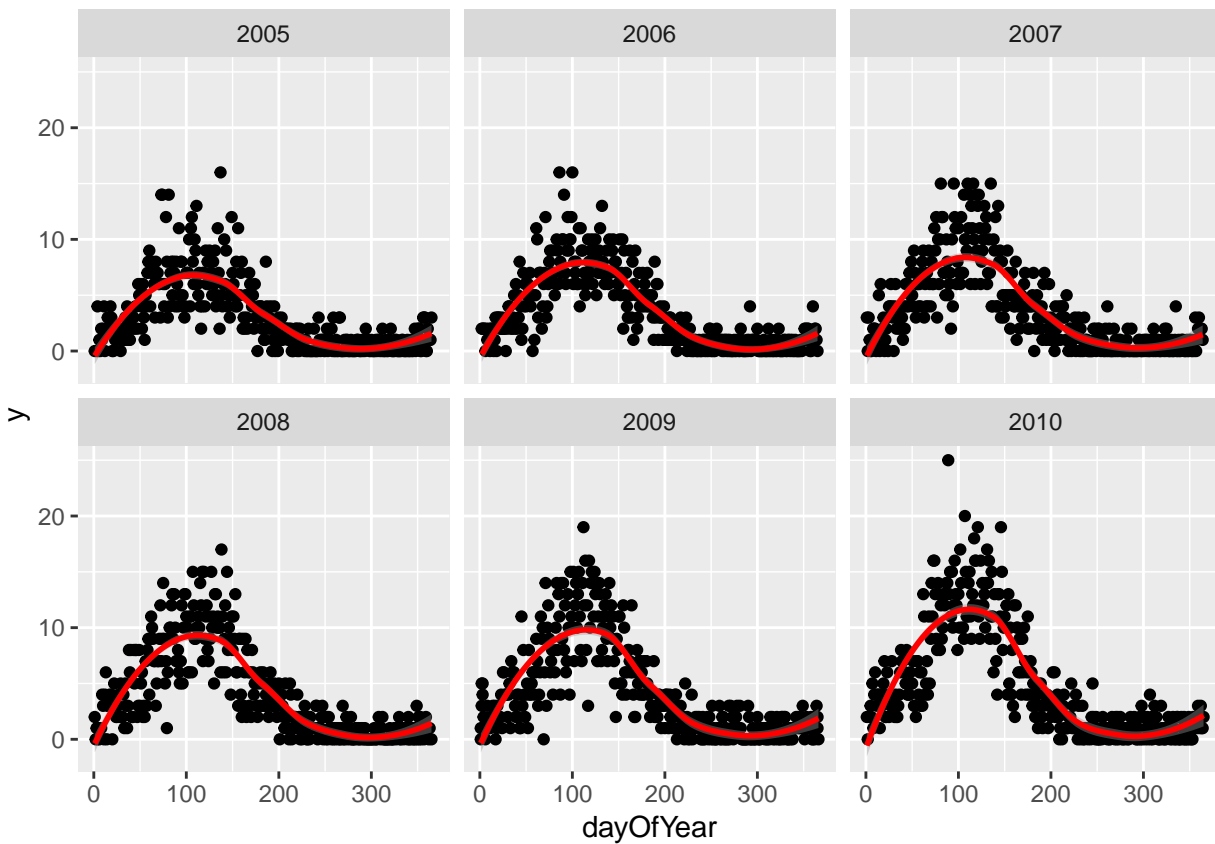


### 3.4 Investigation

Pretending we have no prior knowledge of our dataset, we display the data for few years and see a clear seasonal trend

```
q <- ggplot(d[year %in% c(2005:2010)], aes(x=dayOfYear, y=y))
q <- q + facet_wrap(~year)
q <- q + geom_point()
q <- q + stat_smooth(colour="red")
q
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```





## 3.5 Seasonality

If we want to investigate the seasonality of our data, and identify when are the peaks and troughs, we have a few ways to approach this.

Non-parametric approaches are flexible and easy to implement, but they can lack power and be hard to interpret:

- Create a categorical variable for the seasons (e.g. **spring**, **summer**, **autumn**, **winter**) and include this in the regression model
- Create a categorical variable for the months (e.g. **Jan**, **Feb**, ..., **Dec**) and include this in the regression model

Parametric approaches are more powerful but require more effort:

- Identify the periodicity of the seasonality (how many days between peaks?)
- Using trigonometry, transform **day of year** into variables that appropriately model the observed periodicity
- Obtain coefficient estimates
- Back-transform these estimates into human-understandable values (day of peak, day of trough)

The non-parametric approaches are simple and we will therefore not cover them in this course. We will briefly examine the parametric approach.

*NOTE:* You don't always have to investigate seasonality! It depends entirely on what the purpose of your analysis is!

The Lomb-Scargle Periodogram shows a clear seasonality with a period of 365 days.

```
// STATA CODE STARTS
insheet using "chapter_3.csv", clear

sort date
gen time=_n
tsset time, daily

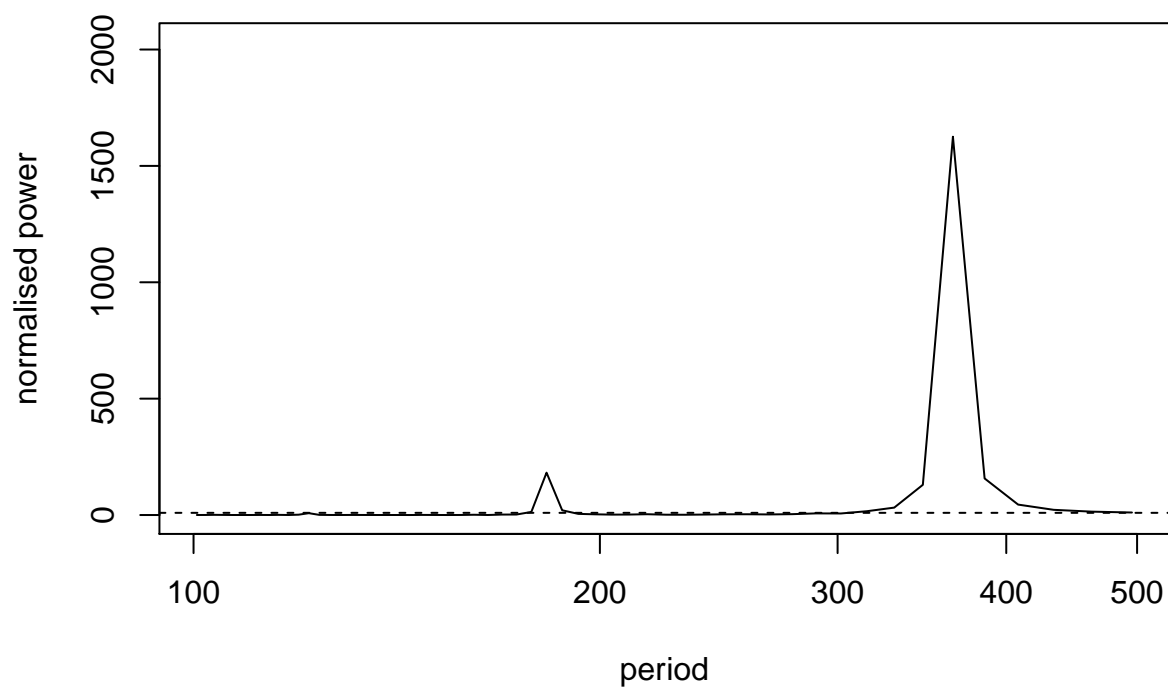
wntestb y

cumsp y, gen(cumulative_spec_dist)
gen period=_N/_n

browse cumulative_spec_dist period
// STATA CODE ENDS

# R CODE
lomb::lsp(d$y,from=100,to=500,ofac=1,type="period")
```

### Lomb–Scargle Periodogram



We then generate two new variables `cos365` and `sin365` and perform a likelihood ratio test to see if they are significant or not. This is done with two simple poisson regressions.

When we do not have autocorrelation, we can use the `glm` function in R and in STATA. Note that it is very important to specify the `family` (as this is how we differentiate between linear/logistic/poisson regressions).

```
// STATA CODE STARTS
gen cos365=cos(dayofyear*2*_pi/365)
gen sin365=sin(dayofyear*2*_pi/365)

glm y yearminus2000 dailyrainfall, family(poisson)
estimates store m1
glm y yearminus2000 dailyrainfall cos365 sin365, family(poisson)
estimates store m2

predict resid, anscombe

lrtest m1 m2
// STATA CODE ENDS

# R CODE
d[,cos365:=cos(dayOfYear*2*pi/365)]
d[,sin365:=sin(dayOfYear*2*pi/365)]

fit0 <- glm(y~yearMinus2000 + dailyrainfall, data=d, family=poisson())
fit1 <- glm(y~yearMinus2000 + dailyrainfall + sin365 + cos365, data=d, family=poisson())

print(lmtest::lrtest(fit0, fit1))

## Likelihood ratio test
##
## Model 1: y ~ yearMinus2000 + dailyrainfall
## Model 2: y ~ yearMinus2000 + dailyrainfall + sin365 + cos365
##   #Df LogLik Df Chisq Pr(>Chisq)
## 1    3 -26904
## 2    5 -12892  2 28024 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We see that the likelihood ratio test for `sin365` and `cos365` was significant, meaning that there is significant seasonality with a 365 day periodicity in our data (which we already strongly suspected due to the periodogram).

We can now run/look at the results of our main regression.

```
print(summary(fit1))

##
## Call:
## glm(formula = y ~ yearMinus2000 + dailyrainfall + sin365 + cos365,
##      family = poisson(), data = d)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -4.0676  -0.9229  -0.1170   0.5861   3.4103
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.0887436  0.0176742   5.021 5.14e-07 ***
## yearMinus2000 0.1016117  0.0010525  96.539 < 2e-16 ***
## dailyrainfall 0.0002287  0.0018476   0.124  0.901
## sin365       1.3972586  0.0103200 135.393 < 2e-16 ***
## cos365      -0.5035265  0.0086308 -58.341 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 45536.8  on 6939  degrees of freedom
## Residual deviance:  7328.5  on 6935  degrees of freedom
## AIC: 25794
##
## Number of Fisher Scoring iterations: 5
```

We also see that the (significant!) coefficient for `year` is 0.1 which means that for each additional year, the outcome increases by  $\exp(0.1)=1.11$ . We also see that the coefficient for `dailyrainfall` was not significant, which means that we did not find a significant association between the outcome and `dailyrainfall`.

*NOTE:* See that this is basically the same as a normal regression.

Through the likelihood ratio test we saw a clear significant seasonal effect. We can now use trigonometry to back-calculate the amplitude and location of peak/troughs from the `cos365` and `sin365` estimates:

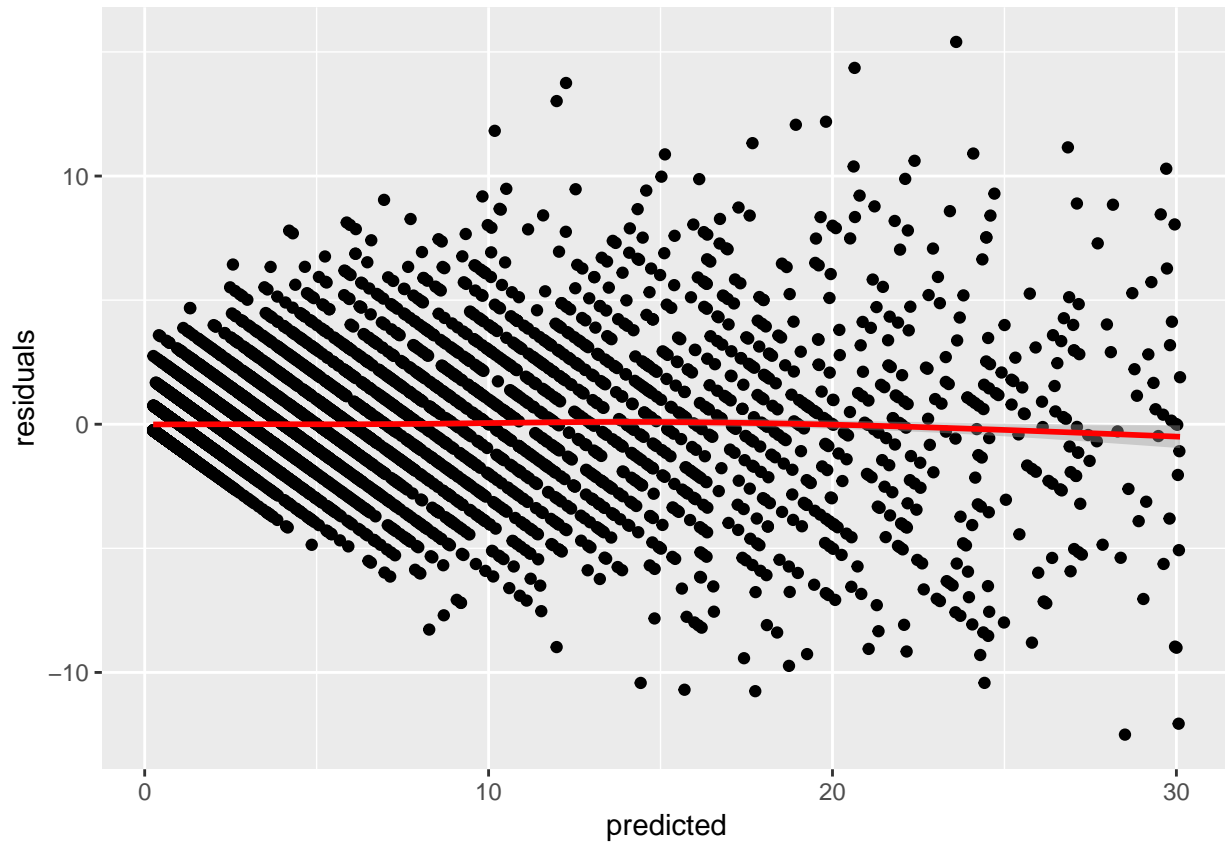
```
b1 <- 1.428417 # sin coefficient
b2 <- -0.512912 # cos coefficient
amplitude <- sqrt(b1^2 + b2^2)
p <- atan(b1/b2) * 365/2/pi
if (p > 0) {
  peak <- p
  trough <- p + 365/2
} else {
  peak <- p + 365/2
  trough <- p + 365
}
if (b1 < 0) {
  g <- peak
  peak <- trough
  trough <- g
}
print(sprintf("amplitude is estimated as %s, peak is estimated as %s, trough is estimated as %s",round(
## [1] "amplitude is estimated as 1.52, peak is estimated as 111, trough is estimated as 294"
print(sprintf("true values are: amplitude: %s, peak: %s, trough: %s",round(AMPLITUDE,2),round(365/4+SEAS
## [1] "true values are: amplitude: 1.5, peak: 111, trough: 294"
```

*NOTE:* An amplitude of 1.5 means that when comparing the average time of year to the peak, the peak is expected to be  $\exp(1.5)=4.5$  times higher than average. We take the exponential because we have run a poisson regression (so think incident rate ratio).

We now investigate our residuals to determine if we have a good fit:

```
d[,residuals:=residuals(fit1, type = "response")]
d[,predicted:=predict(fit1, type = "response")]
q <- ggplot(d,aes(x=predicted,y=residuals))
q <- q + geom_point()
q <- q + stat_smooth(colour="red")
q
```

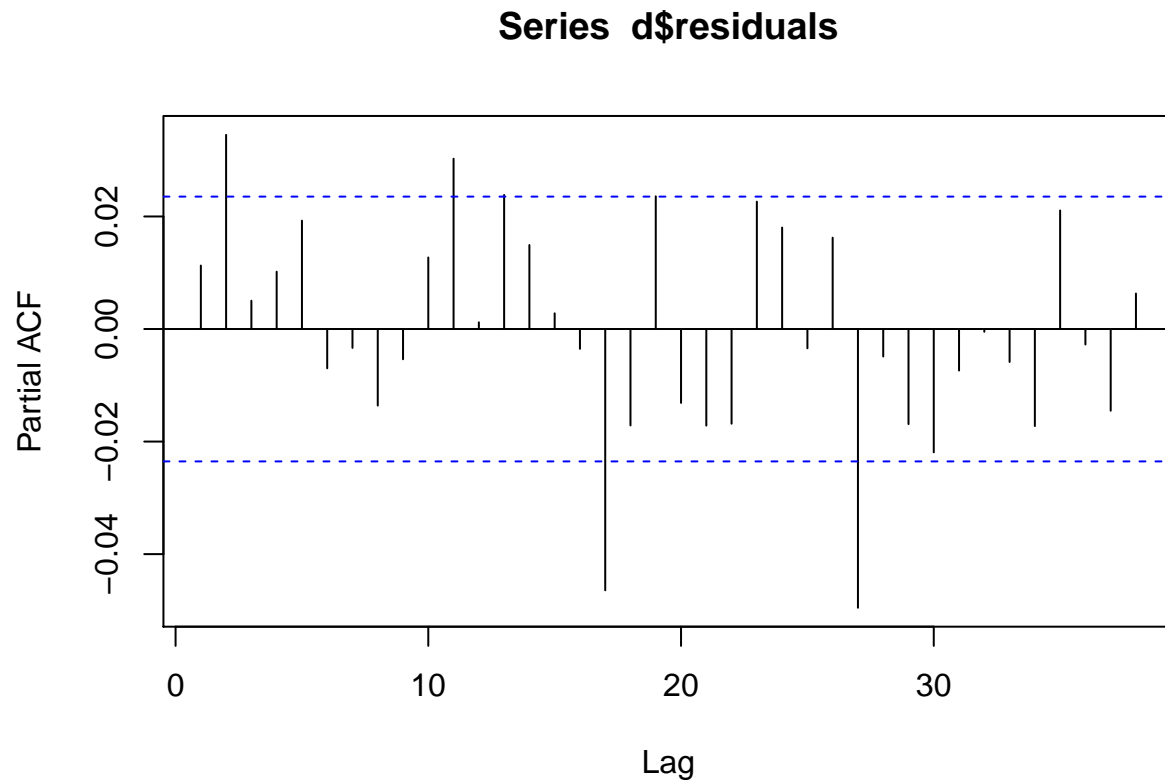
```
## `geom_smooth()` using method = 'gam' and formula 'y ~ s(x, bs = "cs")'
```



We check the `pacf` of the residuals to ensure that it is not **AR**. If we observe **AR** in our residuals, then this model was not appropriate and we need to use a different model.

```
// STATA CODE STARTS  
pac resid  
// STATA CODE ENDS
```

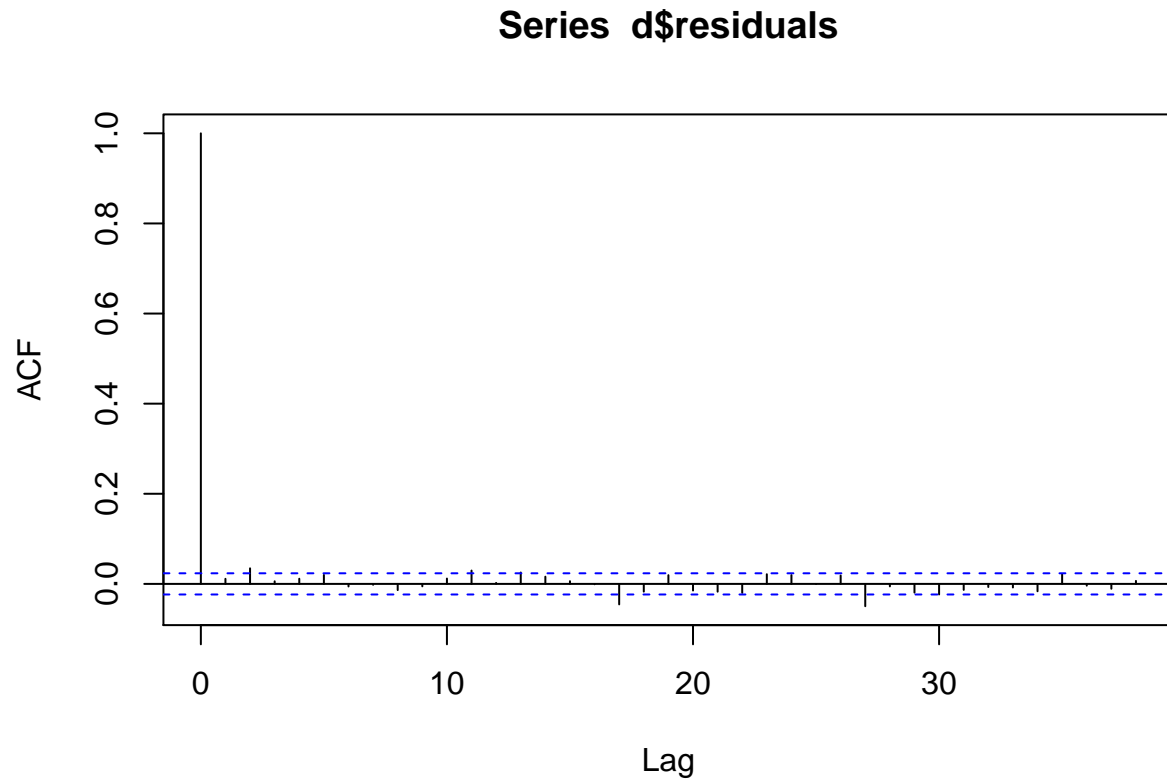
```
# R CODE  
# this is for AR  
pacf(d$residuals)
```



We check the `acf` of the residuals to ensure that it is not MA. If we observe MA in our residuals, then this model was not appropriate and we need to use a different model.

```
// STATA CODE STARTS  
ac resid  
// STATA CODE ENDS
```

```
# R CODE  
# this is for MA  
acf(d$residuals)
```





## Chapter 4

# Panel data: One area with autocorrelation

### 4.1 Aim

We are given a dataset containing daily counts of diseases from one geographical area. We want to identify:

- Does seasonality exist?
- If seasonality exists, when are the high/low seasons?
- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)

(We remove the question about rainfall in order to simplify and streamline the exercise)

## 4.2 Creating the data

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/chapter\\_4.csv](http://rwhite.no/longitudinal_analysis/data/chapter_4.csv)

```
library(data.table)
library(ggplot2)
set.seed(4)

AMPLITUDE <- 1.5
SEASONAL_HORIZONTAL_SHIFT <- 20

d <- data.table(date=seq.Date(
  from=as.Date("2000-01-01"),
  to=as.Date("2018-12-31"),
  by=1))
d[,year:=as.numeric(format.Date(date,"%G"))]
d[,week:=as.numeric(format.Date(date,"%V"))]
d[,month:=as.numeric(format.Date(date,"%m"))]
d[,yearMinus2000:=year-2000]
d[,dayOfSeries:=1:.N]

d[,dayOfYear:=as.numeric(format.Date(date,"%j"))]
d[,seasonalEffect:=sin(2*pi*(dayOfYear-SEASONAL_HORIZONTAL_SHIFT)/365)]
d[,mu := exp(0.1 + yearMinus2000*0.1 + seasonalEffect*AMPLITUDE)]
d[,y:=rpois(.N,mu)]
d[,y:=round(as.numeric(arima.sim(model=list("ar"=c(0.5)), rand.gen = rpois, n=nrow(d), lambda=mu)))]

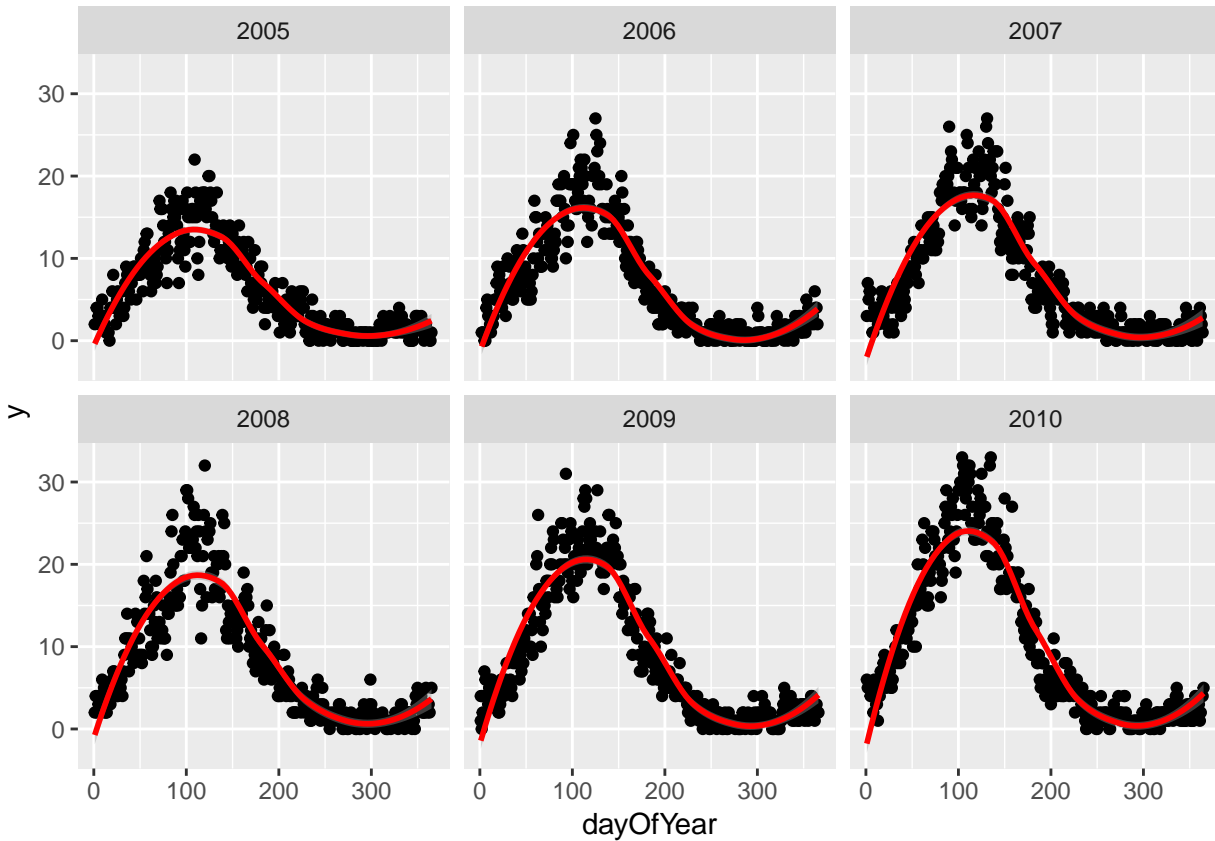
fwrite(d,"data/chapter_4.csv")
```

## 4.3 Investigation

We display the data for few years and see a clear seasonal trend

```
q <- ggplot(d[year %in% c(2005:2010)], aes(x=dayOfYear, y=y))  
q <- q + facet_wrap(~year)  
q <- q + geom_point()  
q <- q + stat_smooth(colour="red")  
q
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```



The Lomb-Scargle Periodogram shows a clear seasonality with a period of 365 days

```
// STATA CODE STARTS
insheet using "chapter_4.csv", clear

sort date
gen time=_n
tsset time, daily

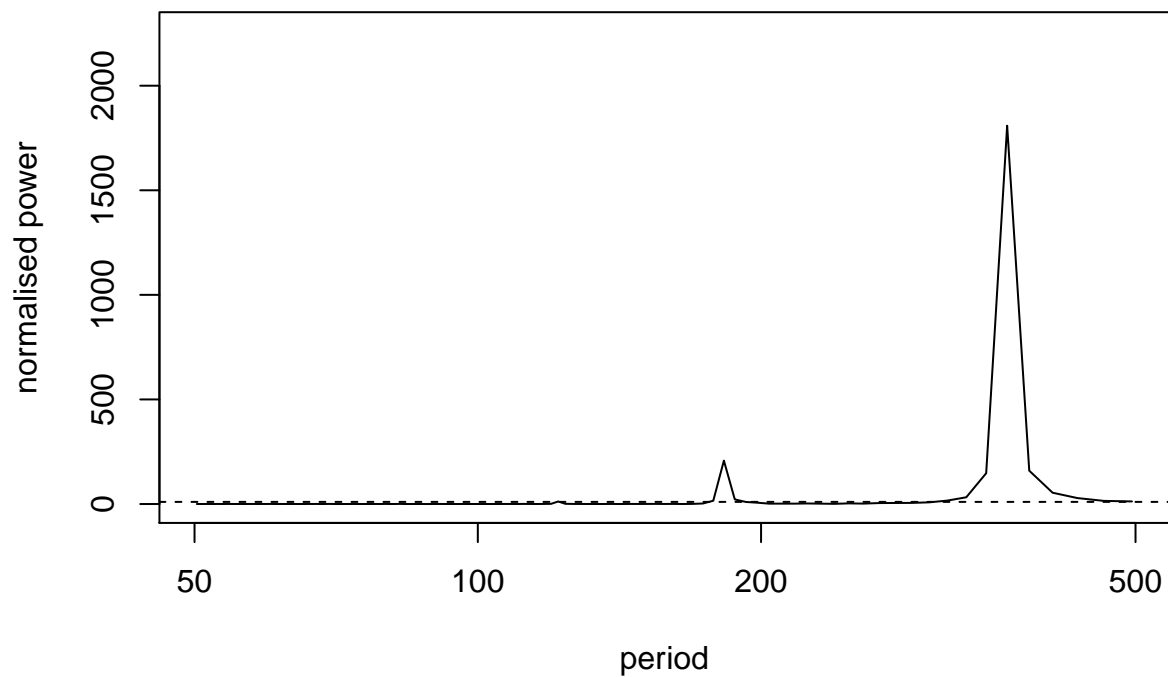
wntestb y

cumsp y, gen(cumulative_spec_dist)
gen period=_N/_n

browse cumulative_spec_dist period
// STATA CODE ENDS

# R CODE
lomb::lsp(d$y,from=50,to=500,ofac=1,type="period")
```

### Lomb–Scargle Periodogram



## 4.4 Regressions

We then generate two new variables `cos365` and `sin365` and perform a likelihood ratio test to see if they are significant or not. This is done with two simple poisson regressions.

```
// STATA CODE STARTS
gen cos365=cos(dayofyear*2*_pi/365)
gen sin365=sin(dayofyear*2*_pi/365)

glm y yearminus2000, family(poisson)
estimates store m1
glm y yearminus2000 cos365 sin365, family(poisson)
estimates store m2

predict resid, anscombe

lrtest m1 m2
// STATA CODE ENDS

# R CODE
d[,cos365:=cos(dayOfYear*2*pi/365)]
d[,sin365:=sin(dayOfYear*2*pi/365)]

fit0 <- glm(y~yearMinus2000, data=d, family=poisson())
fit1 <- glm(y~yearMinus2000+sin365 + cos365, data=d, family=poisson())

print(lmtest::lrtest(fit0, fit1))

## Likelihood ratio test
##
## Model 1: y ~ yearMinus2000
## Model 2: y ~ yearMinus2000 + sin365 + cos365
##   #Df LogLik Df Chisq Pr(>Chisq)
## 1    2 -43124
## 2    4 -14542  2 57163 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We see that the likelihood ratio test for `sin365` and `cos365` was significant, meaning that there is significant seasonality with a 365 day periodicity in our data (which we already strongly suspected due to the periodogram).

We can now run/look at the results of our main regression.

```
print(summary(fit1))

##
## Call:
## glm(formula = y ~ yearMinus2000 + sin365 + cos365, family = poisson(),
##      data = d)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.6774  -0.6738  -0.0503   0.4920   3.5820
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.7981246  0.0105300   75.80  <2e-16 ***
## yearMinus2000 0.0991480  0.0007416  133.70  <2e-16 ***
## sin365        1.4074818  0.0073418  191.71  <2e-16 ***
## cos365       -0.5390314  0.0061513  -87.63  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 81832.6  on 6939  degrees of freedom
## Residual deviance:  5217.8  on 6936  degrees of freedom
## AIC: 29093
##
## Number of Fisher Scoring iterations: 4
```

We also see that the coefficient for year is 0.1 which means that for each additional year, the outcome increases by  $\exp(0.1)=1.11$ .

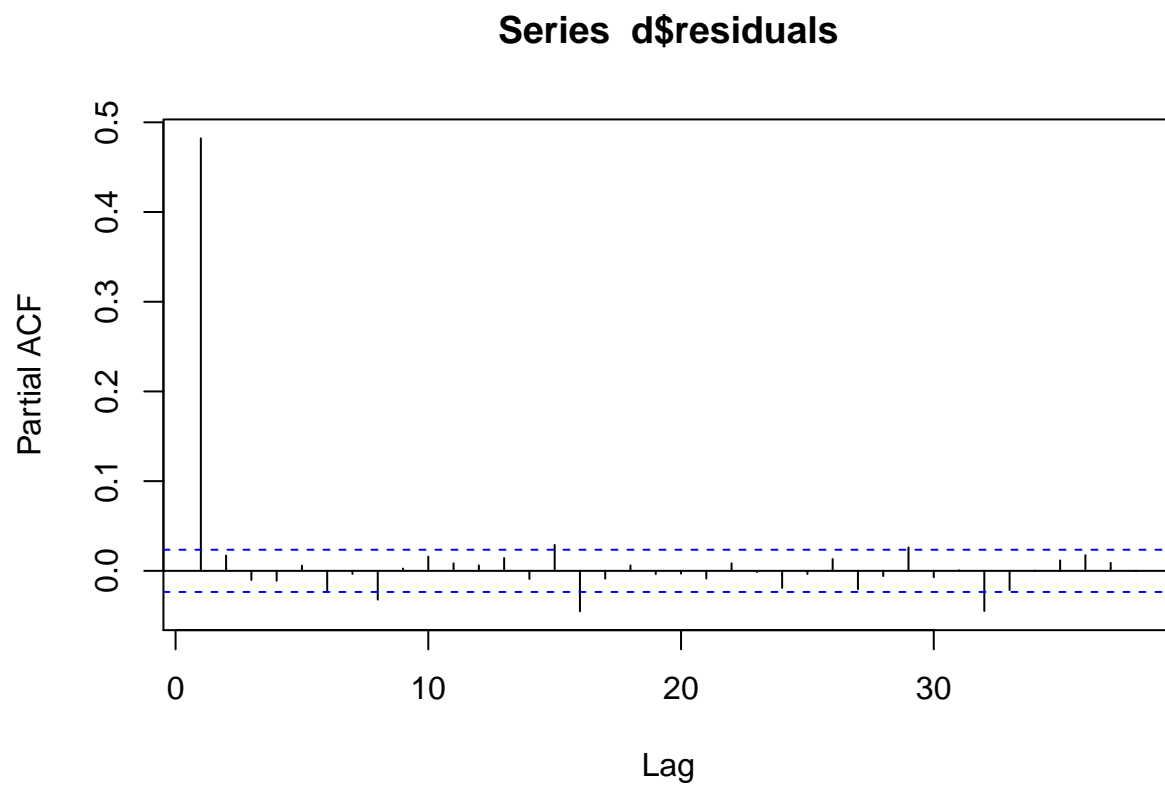
## 4.5 Residual analysis

```
d[,residuals:=residuals(fit1, type = "response")]  
d[,predicted:=predict(fit1, type = "response")]
```

We can see a clear AR(1) pattern in our residuals.

```
// STATA CODE STARTS  
pac resid  
// STATA CODE ENDS
```

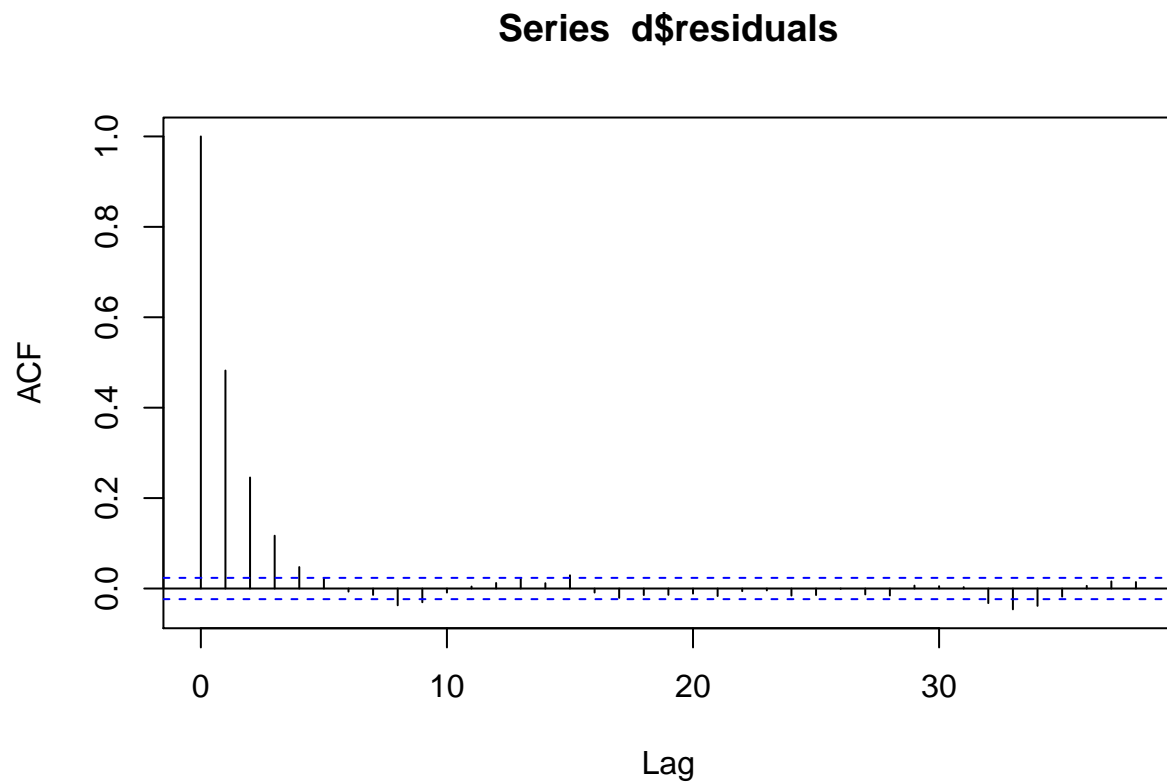
```
# R CODE  
# this is for AR  
pacf(d$residuals)
```



And again we see some sort of AR pattern in our residuals.

```
// STATA CODE STARTS
ac resid
// STATA CODE ENDS

# R CODE
# this is for MA
acf(d$residuals)
```



This means our model is bad, we have autocorrelation. We now need to change our model to account for this AR(1) autocorrelation!



## 4.6 (R ONLY) Regression with AR(1) correlation in residuals

First we create an `id` variable. This generally corresponds to geographical locations, or people. In this case, we only have one geographical location, so our `id` for all observations is 1. This lets the computer know that all data belongs to the same group.

When we have autocorrelation in the residuals, we can use the `MASS::glmpQL` function in R.

```
d[,ID:=1]
# this is for MA
fit <- MASS::glmpQL(y~yearMinus2000+sin365 + cos365, random = ~ 1 | ID,
                    family = poisson, data = d,
                    correlation=nlme::corAR1(form=~dayOfSeries|ID))
```

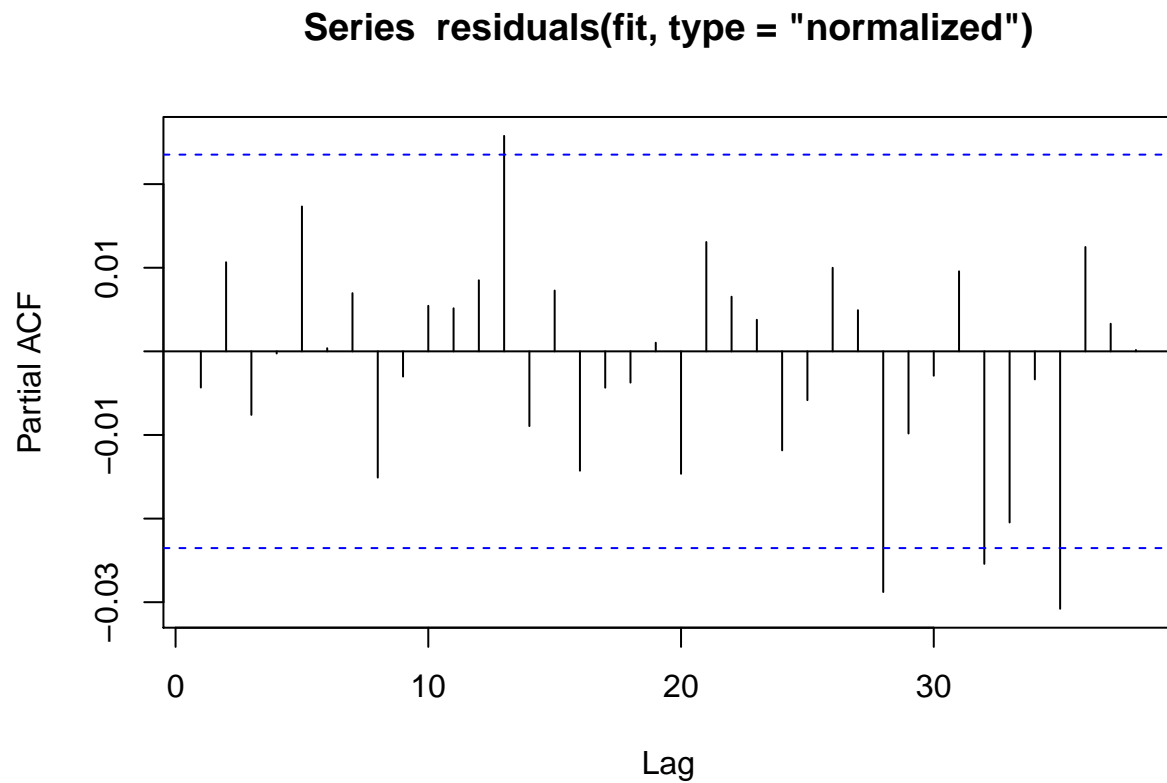
```
## iteration 1
```

```
summary(fit)
```

```
## Linear mixed-effects model fit by maximum likelihood
## Data: d
## AIC BIC logLik
## NA NA NA
##
## Random effects:
## Formula: ~1 | ID
## (Intercept) Residual
## StdDev: 1.149087e-05 0.841689
##
## Correlation Structure: AR(1)
## Formula: ~dayOfSeries | ID
## Parameter estimate(s):
## Phi
## 0.4926123
## Variance function:
## Structure: fixed weights
## Formula: ~invwt
## Fixed effects: y ~ yearMinus2000 + sin365 + cos365
## Value Std.Error DF t-value p-value
## (Intercept) 0.7980540 0.015203158 6936 52.49265 0
## yearMinus2000 0.0991582 0.001070583 6936 92.62077 0
## sin365 1.4074339 0.010596650 6936 132.81876 0
## cos365 -0.5389807 0.008876448 6936 -60.72031 0
## Correlation:
## (Intr) yM2000 sin365
## yearMinus2000 -0.832
## sin365 -0.409 0.000
## cos365 0.186 0.000 -0.158
##
## Standardized Within-Group Residuals:
## Min Q1 Med Q3 Max
## -2.89886750 -0.75775061 -0.05982255 0.60730689 6.49964489
##
## Number of Observations: 6940
## Number of Groups: 1
```

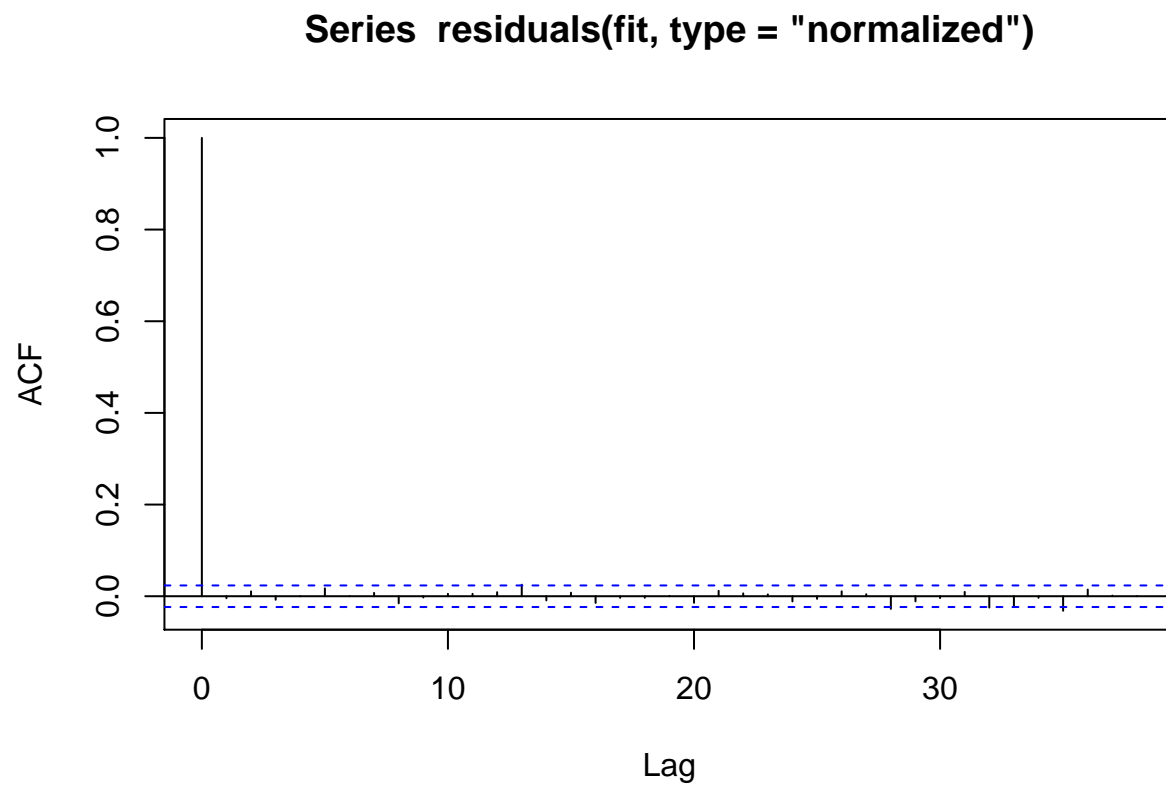
We can see that the residuals no longer display any signs of autocorrelation.

```
pacf(residuals(fit, type = "normalized")) # this is for AR
```



We can see that the residuals no longer display any signs of autocorrelation.

```
acf(residuals(fit, type = "normalized")) # this is for MA
```



We also obtain the same estimates that we did in the last chapter.

```

b1 <- 1.3936185 # sin coefficient
b2 <- -0.5233866 # cos coefficient
amplitude <- sqrt(b1^2 + b2^2)
p <- atan(b1/b2) * 365/2/pi
if (p > 0) {
  peak <- p
  trough <- p + 365/2
} else {
  peak <- p + 365/2
  trough <- p + 365
}
if (b1 < 0) {
  g <- peak
  peak <- trough
  trough <- g
}
print(sprintf("amplitude is estimated as %s, peak is estimated as %s, trough is estimated as %s",round(
## [1] "amplitude is estimated as 1.49, peak is estimated as 112, trough is estimated as 295"
print(sprintf("true values are: amplitude: %s, peak: %s, trough: %s",round(AMPLITUDE,2),round(365/4+SEA
## [1] "true values are: amplitude: 1.5, peak: 111, trough: 294"

```

## 4.7 (STATA ONLY) Regression with robust standard errors

In STATA it is not possible to explicitly model autocorrelation in the residuals (with the exception of linear regression). Since most of our work deals with logistic and poisson regressions, we will be focusing on modelling strategies that work with all kinds of regressions.

The STATA approach to autocorrelation is to estimate more **robust** standard errors. That is, STATA makes the standard errors larger to account for the model misspecification. This is done through the `vce(robust)` option.

```
// STATA CODE STARTS
glm y yearminus2000 cos365 sin365, family(poisson) vce(robust)
// STATA CODE ENDS
```



## Chapter 5

# Not panel data: Multiple areas

### 5.1 Aim

We are given a dataset containing counts of diseases from multiple geographical areas. We want to identify:

- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)
- Is variable  $x$  associated with the outcome?

## 5.2 Creating the data

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/chapter\\_5.csv](http://rwhite.no/longitudinal_analysis/data/chapter_5.csv)

```
library(data.table)
library(lme4)

## Loading required package: Matrix
set.seed(4)

fylkeIntercepts <- data.table(fylke=1:20,fylkeIntercepts=rnorm(20))

d <- data.table(fylke=rep(1:20,each=100))
d <- merge(d,fylkeIntercepts,by="fylke")
d[,mainIntercept:=3]
d[,x:=runif(.N)]
d[,year:=sample(c(1950:2018),.N,replace=T)]
d[,mu := exp(mainIntercept + fylkeIntercepts + 3*x)]
d[,y:=rpois(.N,mu)]

fwrite(d,"data/chapter_5.csv")
```



## 5.3 Investigating the data

We can see from the data that we have 20 geographical areas (`fylke`) with 100 observations for each fylke, but the sampling did not happen consistently (some years have multiple measurements, other years have no measurements).

This means we have:

- multiple geographical areas
- multiple observations in each geographical area
- not panel data

```
print(d)
```

```
##      fylke fylkeIntercepts mainIntercept      x year      mu      y
##    1:      1      0.2167549            3 0.93831909 1966 416.42739 392
##    2:      1      0.2167549            3 0.24217109 1981  51.58692  51
##    3:      1      0.2167549            3 0.56559453 1972 136.12022 135
##    4:      1      0.2167549            3 0.18089910 1950  42.92490  39
##    5:      1      0.2167549            3 0.90449929 1951 376.24959 367
##    ---
## 1996:     20     -0.2834446            3 0.89237059 1995 220.00872 209
## 1997:     20     -0.2834446            3 0.80522348 2006 169.39375 157
## 1998:     20     -0.2834446            3 0.59989167 1955  91.49007  96
## 1999:     20     -0.2834446            3 0.04148228 1996  17.13293  18
## 2000:     20     -0.2834446            3 0.77673920 2002 155.51980 152
```

## 5.4 Regression

For this scenario, we use the `lme4::glmer` function in R. We need to introduce a `(1|fylke)` term to identify the geographical areas (i.e. clusters). In STATA we use the `meglm` function and introduce a `|| fylke:` term to identify the geographical areas (i.e. clusters).

```
// STATA CODE STARTS
insheet using "chapter_5.csv", clear

gen yearMinus2000 = year-2000
meglm y x yearMinus2000 || fylke:, family(poisson)
// STATA CODE ENDS

# R CODE
d[,yearMinus2000:=year-2000]
summary(fit <- lme4::glmer(y~x + yearMinus2000 + (1|fylke),data=d,family=poisson()))

## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: poisson ( log )
## Formula: y ~ x + yearMinus2000 + (1 | fylke)
## Data: d
##
##      AIC      BIC   logLik deviance df.resid
## 15415.5 15437.9 -7703.8 15407.5     1996
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -3.0448 -0.6432 -0.0067  0.6452  4.2338
##
## Random effects:
## Groups Name          Variance Std.Dev.
## fylke (Intercept) 0.6114  0.7819
## Number of obs: 2000, groups: fylke, 20
##
## Fixed effects:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  3.375e+00  1.749e-01  19.295  <2e-16 ***
## x            3.002e+00  5.994e-03 500.874  <2e-16 ***
## yearMinus2000 -9.943e-07  7.192e-05 -0.014    0.989
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##              (Intr) x
## x            -0.025
## yearMns2000  0.007 -0.030
## convergence code: 0
## Model is nearly unidentifiable: very large eigenvalue
## - Rescale variables?
## Model is nearly unidentifiable: large eigenvalue ratio
## - Rescale variables?
```

You can see that the format of the results is the same as an ordinary regression.

## Chapter 6

# Panel data: multiple areas without autocorrelation

### 6.1 Aim

We are given a dataset containing daily counts of diseases from multiple geographical areas. We want to identify:

- Does seasonality exist?
- If seasonality exists, when are the high/low seasons?
- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)

## 6.2 Creating the data

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/chapter\\_6.csv](http://rwhite.no/longitudinal_analysis/data/chapter_6.csv)

```
library(data.table)
library(ggplot2)
set.seed(4)

AMPLITUDE <- 1.5
SEASONAL_HORIZONTAL_SHIFT <- 20

fylkeIntercepts <- data.table(fylke=1:20,fylkeIntercepts=rnorm(20))

d <- data.table(date=seq.Date(
  from=as.Date("2010-01-01"),
  to=as.Date("2015-12-31"),
  by=1))
d[,year:=as.numeric(format.Date(date,"%G"))]
d[,week:=as.numeric(format.Date(date,"%V"))]
d[,month:=as.numeric(format.Date(date,"%m"))]

temp <- vector("list",length=20)
for(i in 1:20){
  temp[[i]] <- copy(d)
  temp[[i]][,fylke:=i]
}
d <- rbindlist(temp)

d[,yearMinus2000:=year-2000]
d[,dayOfSeries:=1:.N]

d[,dayOfYear:=as.numeric(format.Date(date,"%j"))]
d[,seasonalEffect:=sin(2*pi*(dayOfYear-SEASONAL_HORIZONTAL_SHIFT)/365)]
d[,mu := exp(0.1 + yearMinus2000*0.1 + seasonalEffect*AMPLITUDE)]
d[,y:=rpois(.N,mu)]

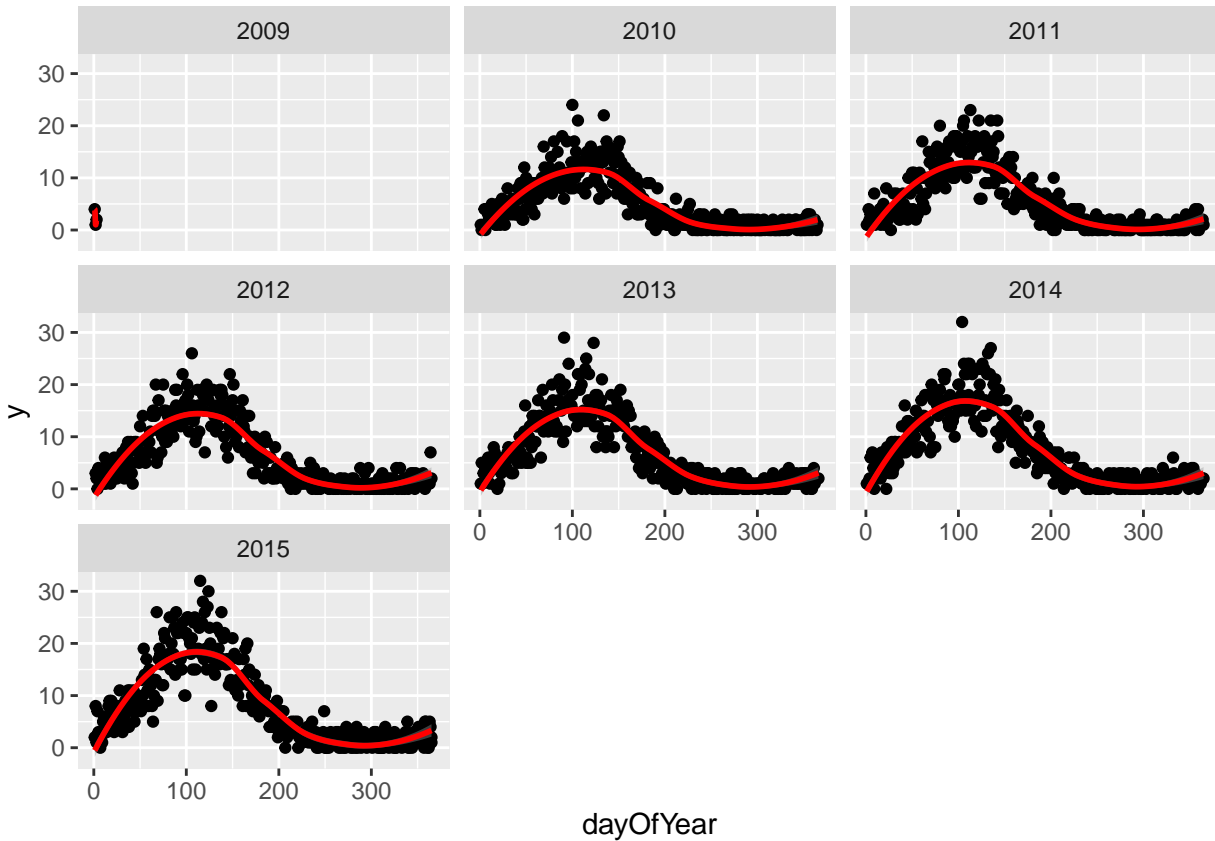
fwrite(d,"data/chapter_6.csv")
```

## 6.3 Investigation

We then drill down into a few years for fylke 1, and see a clear seasonal trend

```
q <- ggplot(d[fylke==1], aes(x=dayOfYear, y=y))
q <- q + facet_wrap(~year)
q <- q + geom_point()
q <- q + stat_smooth(colour="red")
q
```

## `geom\_smooth()` using method = 'loess' and formula 'y ~ x'



The Lomb-Scargle Periodogram shows a clear seasonality with a period of 365 days

```
// STATA CODE STARTS
insheet using "chapter_6.csv", clear

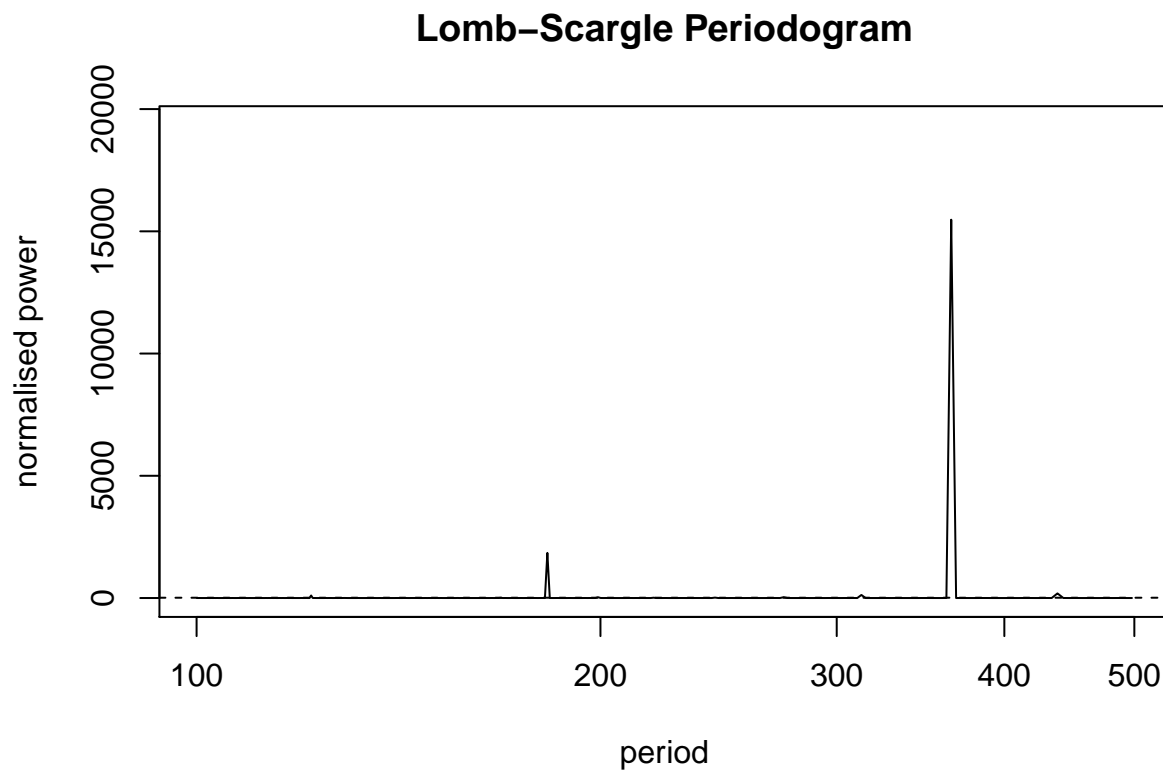
sort fylke date
by fylke: gen time=_n
tsset fylke time, daily

wntestb y if fylke==1

cumsp y if fylke==1, gen(cumulative_spec_dist)
by fylke: gen period=_N/_n

browse cumulative_spec_dist period
// STATA CODE ENDS

# RCODE
lomb::lsp(d$y,from=100,to=500,ofac=1,type="period")
```



## 6.4 Regression

First we create an `id` variable. This generally corresponds to geographical locations, or people. In this case, we only have one geographical location, so our `id` for all observations is 1. This lets the computer know that all data belongs to the same group.

When we have panel data with multiple areas, we use the `MASS::glmPQL` function in R and the `meglm` function in STATA. In R we identify the geographical areas with `random = ~ 1 | fylke` and in STATA with `|| fylke:`.

```
// STATA CODE STARTS
gen cos365=cos(dayofyear*2*_pi/365)
gen sin365=sin(dayofyear*2*_pi/365)

meglm y yearminus2000 || fylke:, family(poisson) iter(10)
estimates store m1
meglm y yearminus2000 cos365 sin365 || fylke:, family(poisson) iter(10)
estimates store m2

predict resid, anscombe

lrtest m1 m2
// STATA CODE ENDS
```

```
# R CODE
d[,cos365:=cos(dayOfYear*2*pi/365)]
d[,sin365:=sin(dayOfYear*2*pi/365)]
fit0 <- MASS::glmPQL(y~yearMinus2000, random = ~ 1 | fylke,
                     family = poisson, data = d)

## iteration 1
fit1 <- MASS::glmPQL(y~yearMinus2000 + sin365 + cos365, random = ~ 1 | fylke,
                     family = poisson, data = d)

## iteration 1
print(lmtest::lrtest(fit0, fit1))

## Likelihood ratio test
##
## Model 1: y ~ yearMinus2000
## Model 2: y ~ yearMinus2000 + sin365 + cos365
##   #Df LogLik Df Chisq Pr(>Chisq)
## 1    4
## 2    6      2
```

We see that the likelihood ratio test for `sin365` and `cos365` was significant, meaning that there is significant seasonality with a 365 day periodicity in our data (which we already strongly suspected due to the periodogram).

We can now run/look at the results of our main regression.

```
print(summary(fit1))

## Linear mixed-effects model fit by maximum likelihood
## Data: d
##   AIC BIC logLik
##   NA  NA   NA
##
## Random effects:
## Formula: ~1 | fylke
##          (Intercept) Residual
## StdDev: 1.583894e-05 0.9976713
##
## Variance function:
## Structure: fixed weights
## Formula: ~invwt
## Fixed effects: y ~ yearMinus2000 + sin365 + cos365
##              Value Std.Error DF t-value p-value
## (Intercept)  0.1122536 0.014488403 43797 7.7478 0
## yearMinus2000 0.0989047 0.001109477 43797 89.1453 0
## sin365        1.4095095 0.003695341 43797 381.4288 0
## cos365        -0.5109375 0.003083683 43797 -165.6907 0
## Correlation:
##              (Intr) yM2000 sin365
## yearMinus2000 -0.979
## sin365        -0.150 0.000
## cos365        0.065 -0.001 -0.151
##
## Standardized Within-Group Residuals:
##              Min          Q1          Med          Q3          Max
## -3.19682240 -0.82387498 -0.07501834 0.63400484 5.82452468
##
## Number of Observations: 43820
## Number of Groups: 20
```

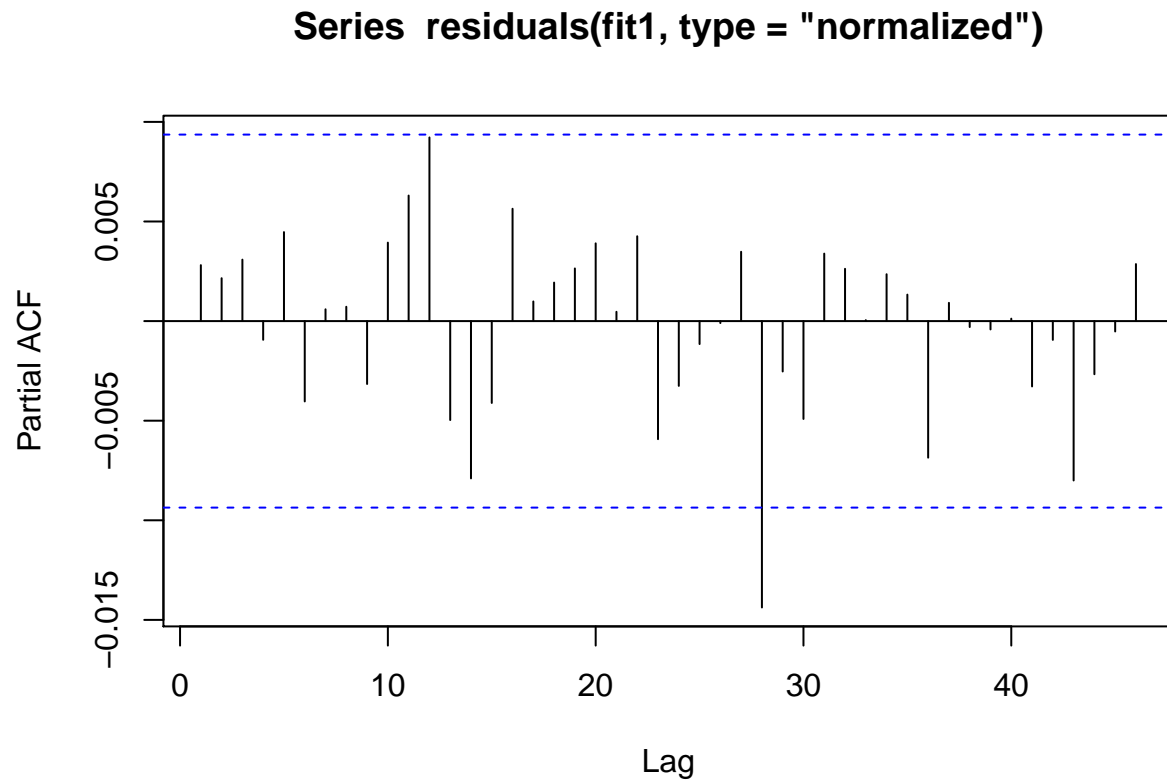


## 6.5 Residual analysis

We see that there is no evidence of autoregression in the residuals

```
// STATA CODE STARTS
pac resid if fylke==1
// STATA CODE ENDS

# R CODE
pacf(residuals(fit1, type = "normalized")) # this is for AR
```

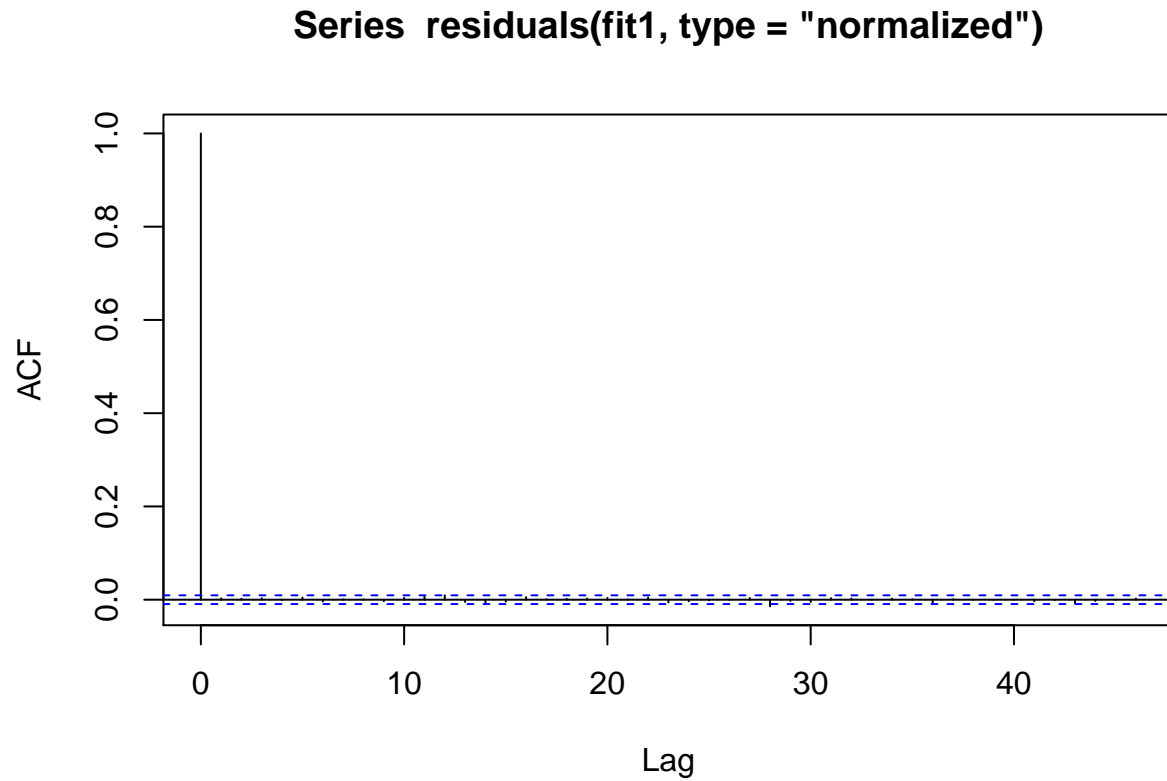


We see that there is no evidence of autoregression in the residuals

```
// STATA CODE STARTS  
ac resid if fylke==1  
// STATA CODE ENDS
```

```
# R CODE
```

```
acf(residuals(fit1, type = "normalized")) # this is for MA
```



We also obtain the same estimates that we did in the last chapter.

```
b1 <- 1.4007640 # sin coefficient
b2 <- -0.5234863 # cos coefficient
amplitude <- sqrt(b1^2 + b2^2)
p <- atan(b1/b2) * 365/2/pi
if (p > 0) {
  peak <- p
  trough <- p + 365/2
} else {
  peak <- p + 365/2
  trough <- p + 365
}
if (b1 < 0) {
  g <- peak
  peak <- trough
  trough <- g
}
print(sprintf("amplitude is estimated as %s, peak is estimated as %s, trough is estimated as %s",round(
## [1] "amplitude is estimated as 1.5, peak is estimated as 112, trough is estimated as 295"
print(sprintf("true values are: amplitude: %s, peak: %s, trough: %s",round(AMPLITUDE,2),round(365/4+SEA
## [1] "true values are: amplitude: 1.5, peak: 111, trough: 294"
```



## Chapter 7

# Panel data: multiple areas with autocorrelation

### 7.1 Aim

We are given a dataset containing daily counts of diseases from multiple geographical areas. We want to identify:

- Does seasonality exist?
- If seasonality exists, when are the high/low seasons?
- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)

## 7.2 Creating the data

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/chapter\\_7.csv](http://rwhite.no/longitudinal_analysis/data/chapter_7.csv)

```
library(data.table)
library(ggplot2)
set.seed(4)

AMPLITUDE <- 1.5
SEASONAL_HORIZONTAL_SHIFT <- 20

fylkeIntercepts <- data.table(fylke=1:20,fylkeIntercepts=rnorm(20))

d <- data.table(date=seq.Date(
  from=as.Date("2010-01-01"),
  to=as.Date("2015-12-31"),
  by=1))
d[,year:=as.numeric(format.Date(date,"%G"))]
d[,week:=as.numeric(format.Date(date,"%V"))]
d[,month:=as.numeric(format.Date(date,"%m"))]

temp <- vector("list",length=20)
for(i in 1:20){
  temp[[i]] <- copy(d)
  temp[[i]][,fylke:=i]
}
d <- rbindlist(temp)

d[,yearMinus2000:=year-2000]
d[,dayOfSeries:=1:.N]

d[,dayOfYear:=as.numeric(format.Date(date,"%j"))]
d[,seasonalEffect:=sin(2*pi*(dayOfYear-SEASONAL_HORIZONTAL_SHIFT)/365)]
d[,mu := round(exp(0.1 + yearMinus2000*0.1 + seasonalEffect*AMPLITUDE))]
d[,y:=rpois(.N,mu)]
d[,y:=mu+round(as.numeric(arima.sim(model=list("ar"=c(0.5)), rand.gen = rpois, n=nrow(d), lambda=mu)))]

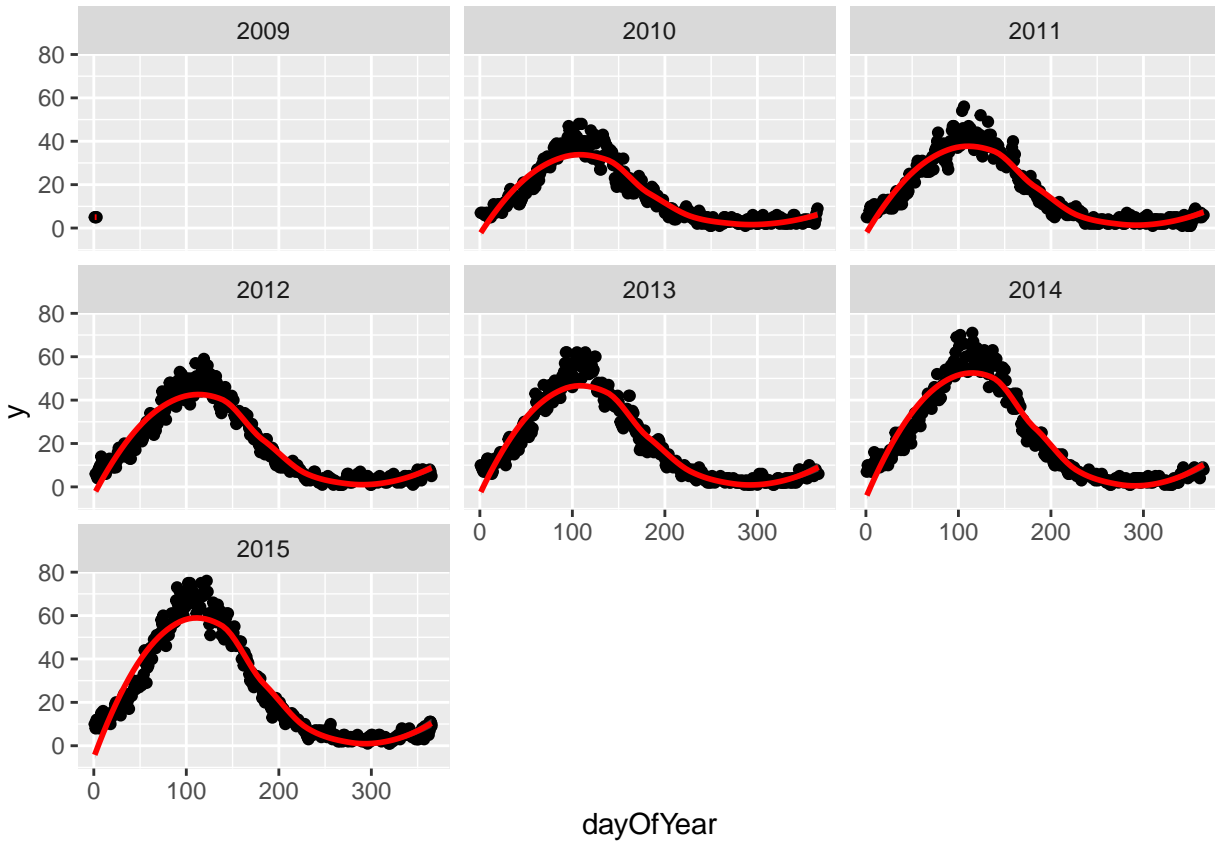
fwrite(d,"data/chapter_7.csv")
```

## 7.3 Investigation

We drill down into a few years in fylke 1, and see a clear seasonal trend

```
q <- ggplot(d[fylke==1], aes(x=dayOfYear, y=y))
q <- q + facet_wrap(~year)
q <- q + geom_point()
q <- q + stat_smooth(colour="red")
q
```

## `geom\_smooth()` using method = 'loess' and formula 'y ~ x'



The Lomb-Scargle Periodogram shows a clear seasonality with a period of 365 days

```
// STATA CODE STARTS
insheet using "chapter_7.csv", clear

sort fylke date
by fylke: gen time=_n
tsset fylke time, daily

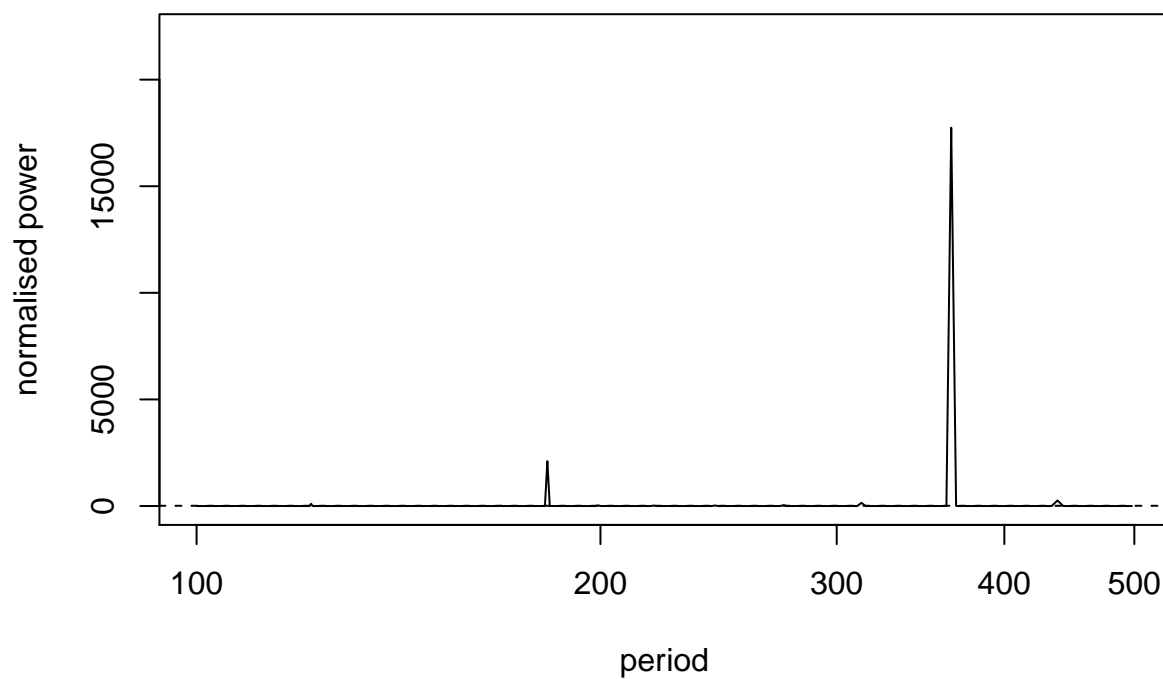
wntestb y if fylke==1

cumsp y if fylke==1, gen(cumulative_spec_dist)
by fylke: gen period=_N/_n

browse cumulative_spec_dist period
// STATA CODE ENDS

# R CODE
lomb::lsp(d$y,from=100,to=500,ofac=1,type="period")
```

### Lomb–Scargle Periodogram





## 7.4 Regressions

First we create an `id` variable. This generally corresponds to geographical locations, or people. In this case, we only have one geographical location, so our `id` for all observations is 1. This lets the computer know that all data belongs to the same group.

When we have panel data with multiple areas, we use the `MASS::glmPQL` function in R and the `meglm` function in STATA. In R we identify the geographical areas with `random = ~ § | fylke` and in STATA with `|| fylke:`.

```
// STATA CODE STARTS
gen cos365=cos(dayofyear*2*_pi/365)
gen sin365=sin(dayofyear*2*_pi/365)

meglm y yearminus2000 || fylke:, family(poisson) iter(10)
estimates store m1
meglm y yearminus2000 cos365 sin365 || fylke:, family(poisson) iter(10)
estimates store m2

predict resid, anscombe

lrtest m1 m2
// STATA CODE ENDS
```

```
# R CODE
d[,cos365:=cos(dayOfYear*2*pi/365)]
d[,sin365:=sin(dayOfYear*2*pi/365)]
fit0 <- MASS::glmPQL(y~yearMinus2000, random = ~ 1 | fylke,
                    family = poisson, data = d)

## iteration 1
fit1 <- MASS::glmPQL(y~yearMinus2000 + sin365 + cos365, random = ~ 1 | fylke,
                    family = poisson, data = d)

## iteration 1
## iteration 2
print(lmtest::lrtest(fit0, fit1))

## Likelihood ratio test
##
## Model 1: y ~ yearMinus2000
## Model 2: y ~ yearMinus2000 + sin365 + cos365
##   #Df LogLik Df Chisq Pr(>Chisq)
## 1    4
## 2    6      2
```

We see that the likelihood ratio test for `sin365` and `cos365` was significant, meaning that there is significant seasonality with a 365 day periodicity in our data (which we already strongly suspected due to the periodogram).

We can now run/look at the results of our main regression.

```
print(summary(fit1))

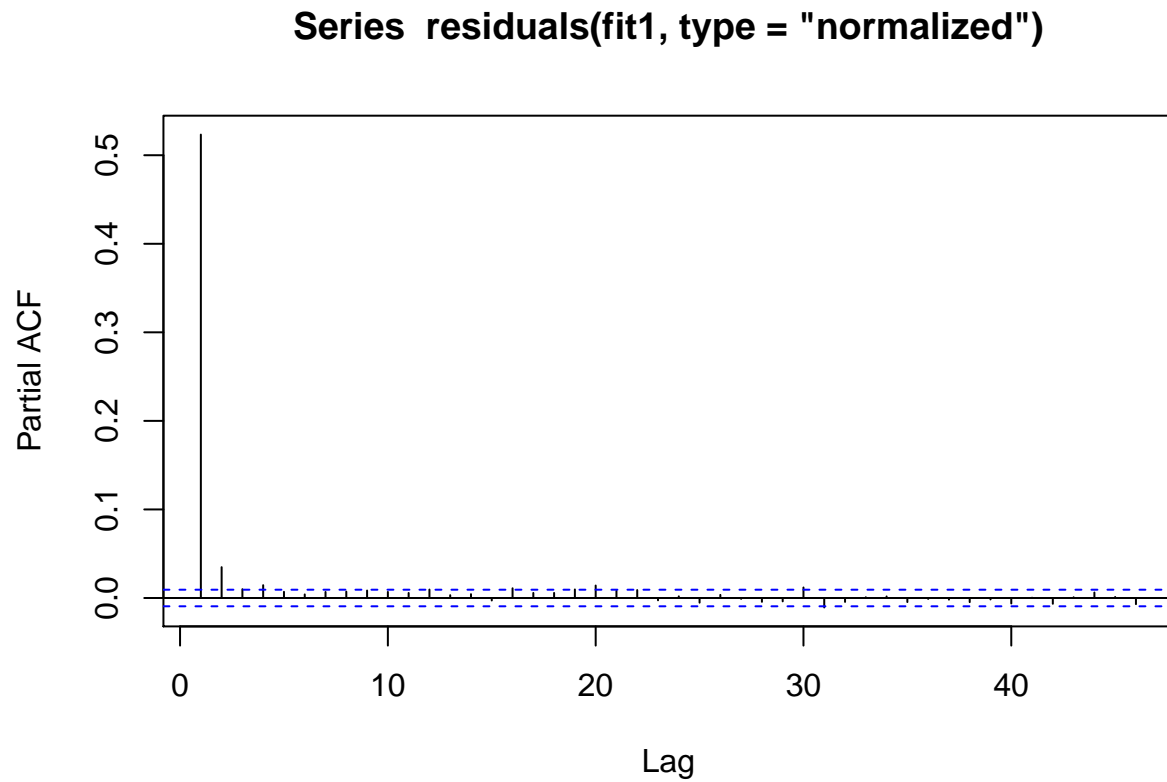
## Linear mixed-effects model fit by maximum likelihood
## Data: d
##   AIC BIC logLik
##   NA  NA   NA
##
## Random effects:
## Formula: ~1 | fylke
##          (Intercept) Residual
## StdDev: 0.004579768 0.7191519
##
## Variance function:
## Structure: fixed weights
## Formula: ~invwt
## Fixed effects: y ~ yearMinus2000 + sin365 + cos365
##              Value Std.Error DF t-value p-value
## (Intercept)  1.2189925 0.006110555 43797 199.4896      0
## yearMinus2000 0.0987374 0.000461394 43797 213.9980      0
## sin365        1.3990267 0.001531179 43797 913.6928      0
## cos365        -0.5171211 0.001282191 43797 -403.3106      0
## Correlation:
##              (Intr) yM2000 sin365
## yearMinus2000 -0.966
## sin365        -0.147 0.000
## cos365         0.065 -0.001 -0.152
##
## Standardized Within-Group Residuals:
##           Min           Q1           Med           Q3           Max
## -3.00864057 -0.70228031 -0.06334676  0.64274011  5.21710224
##
## Number of Observations: 43820
## Number of Groups: 20
```

## 7.5 Residual analysis

We see that there is an AR(1) autocorrelation in the residuals, meaning that our model is not appropriate.

```
// STATA CODE STARTS  
pac resid if fylke==1  
// STATA CODE ENDS
```

```
# R CODE  
pacf(residuals(fit1, type = "normalized")) # this is for AR
```

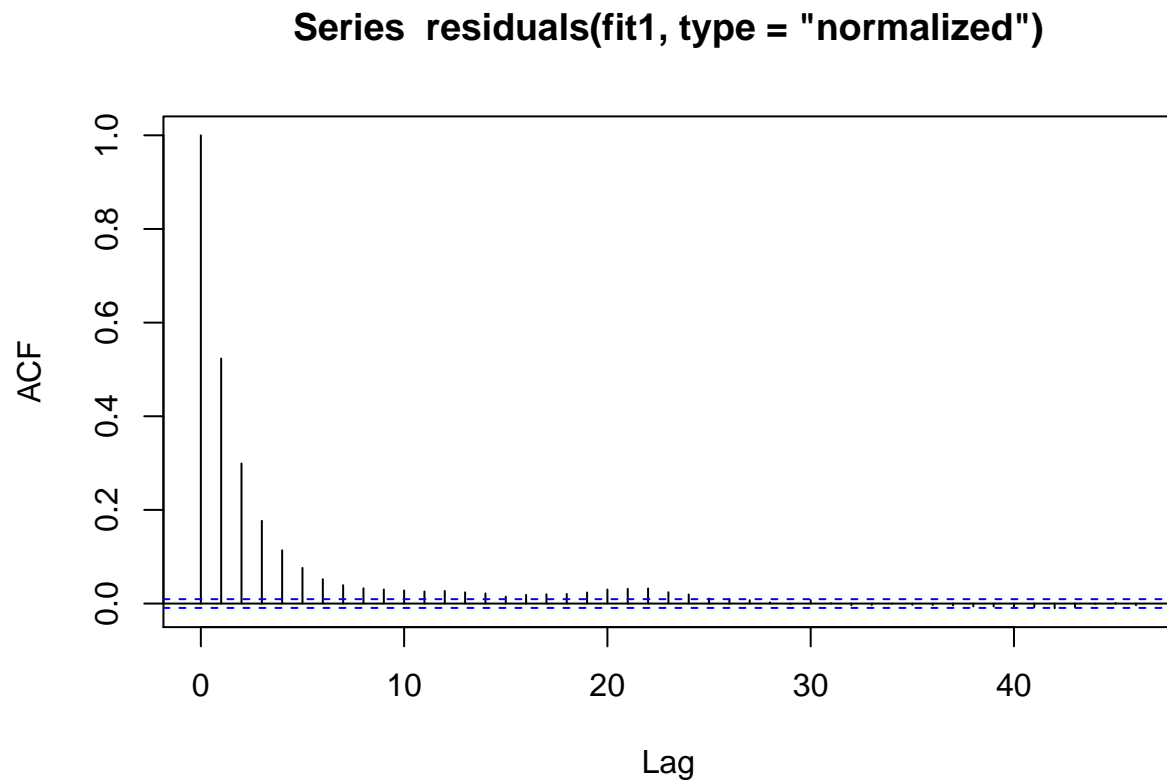


We see that there is some sort of AR autocorrelation in the residuals, meaning that our model is not appropriate.

```
// STATA CODE STARTS  
ac resid if fylke==1  
// STATA CODE ENDS
```

```
# R CODE
```

```
acf(residuals(fit1, type = "normalized")) # this is for MA
```



## 7.6 (R ONLY) Regression with AR(1) correlation in residuals

We include `correlation=nlme::corAR1(form=~dayOfSeries|fylke)` or in other words `correlation=nlme::corAR1(form=~dayOfSeries|fylke)` to let the computer know what is the time variable and what is the group variable.

```
fit1 <- MASS::glmPQL(y~yearMinus2000+sin365 + cos365, random = ~ 1 | fylke,
                    family = poisson, data = d,
                    correlation=nlme::corAR1(form=~dayOfSeries|fylke))
```

```
## iteration 1
```

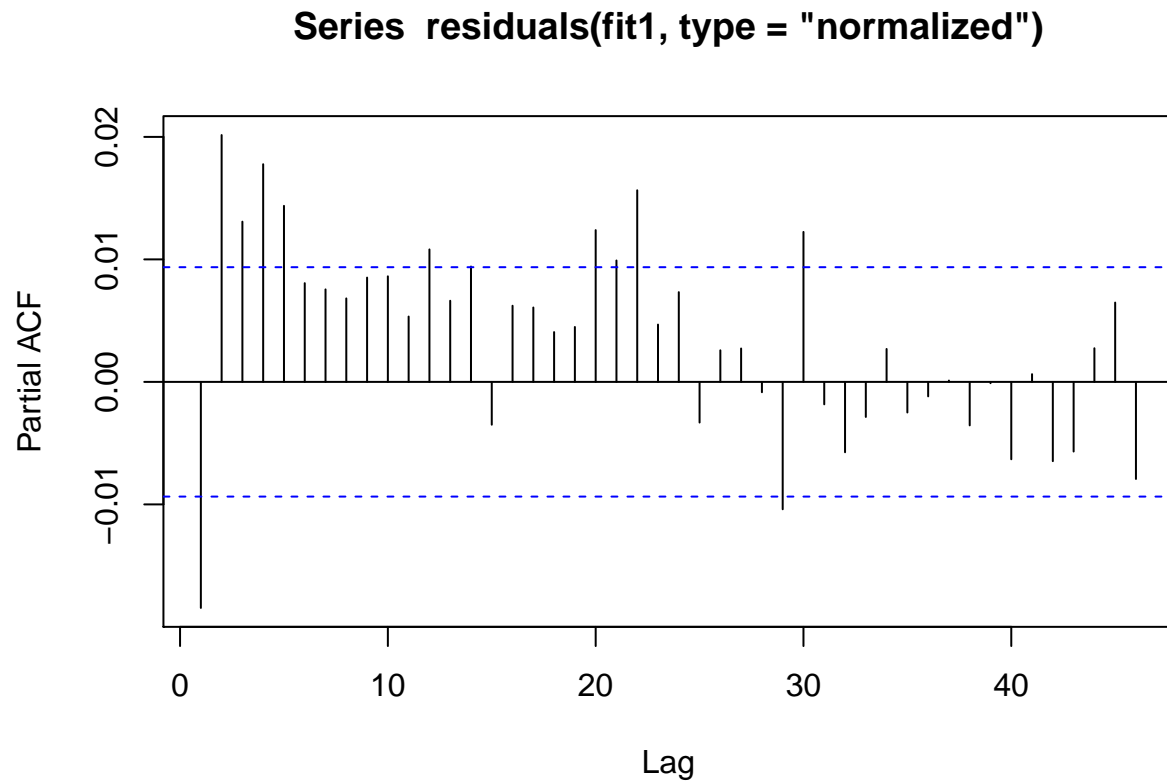
```
summary(fit1)
```

```
## Linear mixed-effects model fit by maximum likelihood
## Data: d
## AIC BIC logLik
## NA NA NA
##
## Random effects:
## Formula: ~1 | fylke
## (Intercept) Residual
## StdDev: 2.40488e-05 0.7195239
##
## Correlation Structure: AR(1)
## Formula: ~dayOfSeries | fylke
## Parameter estimate(s):
## Phi
## 0.5240054
## Variance function:
## Structure: fixed weights
## Formula: ~invwt
## Fixed effects: y ~ yearMinus2000 + sin365 + cos365
## Value Std.Error DF t-value p-value
## (Intercept) 1.2195477 0.010774796 43797 113.1852 0
## yearMinus2000 0.0987065 0.000825226 43797 119.6115 0
## sin365 1.3988945 0.002739109 43797 510.7116 0
## cos365 -0.5169579 0.002292465 43797 -225.5030 0
## Correlation:
## (Intr) yM2000 sin365
## yearMinus2000 -0.979
## sin365 -0.149 0.001
## cos365 0.066 -0.001 -0.151
##
## Standardized Within-Group Residuals:
## Min Q1 Med Q3 Max
## -2.99731654 -0.70249782 -0.06736726 0.64264790 5.20296607
##
## Number of Observations: 43820
## Number of Groups: 20
```

## 7.7 Residual analysis

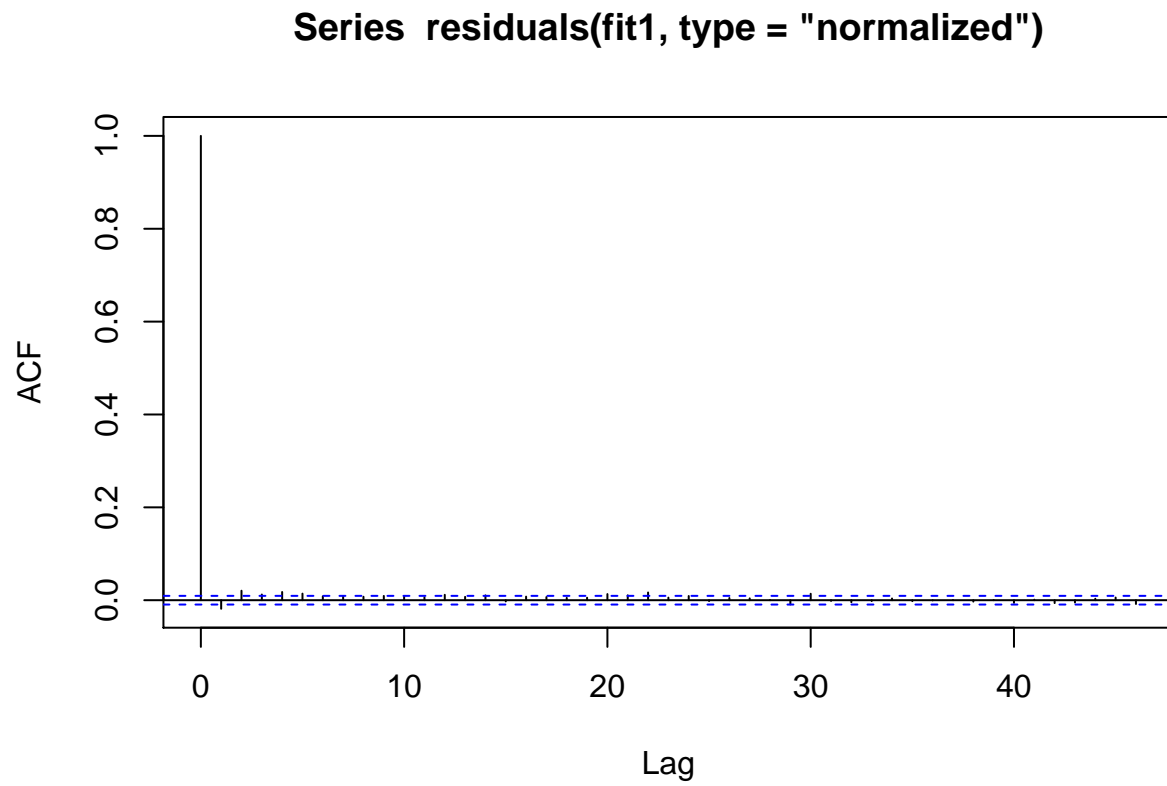
We see that the vast majority of the autoregression in the residuals has been removed.

```
pacf(residuals(fit1, type = "normalized")) # this is for AR
```



We see that the vast majority of the autoregression in the residuals has been removed.

```
acf(residuals(fit1, type = "normalized")) # this is for MA
```



We obtain the same estimates that we did in the last chapter.

```

b1 <- 1.4007640 # sin coefficient
b2 <- -0.5234863 # cos coefficient
amplitude <- sqrt(b1^2 + b2^2)
p <- atan(b1/b2) * 365/2/pi
if (p > 0) {
  peak <- p
  trough <- p + 365/2
} else {
  peak <- p + 365/2
  trough <- p + 365
}
if (b1 < 0) {
  g <- peak
  peak <- trough
  trough <- g
}
print(sprintf("amplitude is estimated as %s, peak is estimated as %s, trough is estimated as %s",round(
## [1] "amplitude is estimated as 1.5, peak is estimated as 112, trough is estimated as 295"
print(sprintf("true values are: amplitude: %s, peak: %s, trough: %s",round(AMPLITUDE,2),round(365/4+SEA
## [1] "true values are: amplitude: 1.5, peak: 111, trough: 294"

```



## 7.8 (STATA ONLY) Regression with robust standard errors

In STATA it is not possible to explicitly model autocorrelation in the residuals (with the exception of linear regression). Since most of our work deals with logistic and poisson regressions, we will be focusing on modelling strategies that work with all kinds of regressions.

The STATA approach to autocorrelation is to estimate more **robust** standard errors. That is, STATA makes the standard errors larger to account for the model misspecification. This is done through the `vce(robust)` option.

```
// STATA CODE STARTS  
meglm y yearminus2000 cos365 sin365 || fylke:, family(poisson) iter(10) vce(robust)  
// STATA CODE ENDS
```



## Chapter 8

# Exercises

### 8.1 Exercise 1

We are given a dataset containing daily counts of diseases  $y$  from one geographical area. We want to identify:

- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)
- Does seasonality exist (use the categorical variable “season”)?
- What season has the most cases? (Spring/summer/autumn/winter?)
- Is `numberOfCows` associated with the outcome  $y$ ?

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/exercise\\_1.csv](http://rwhite.no/longitudinal_analysis/data/exercise_1.csv)

```
library(data.table)
set.seed(4)

d <- data.table(date=seq.Date(
  from=as.Date("2010-01-01"),
  to=as.Date("2015-12-31"),
  by=1))

d[,numberOfCows:=rpois(.N,5)]

d[,year:=as.numeric(format.Date(date,"%G"))]
d[,week:=as.numeric(format.Date(date,"%V"))]
d[,month:=as.numeric(format.Date(date,"%m"))]
d[,season:="Winter"]
d[month %in% c(3:5), season:="Spring"]
d[month %in% c(6:8), season:="Summer"]
d[month %in% c(9:11), season:="Autumn"]

d[,seasonIntercept:=0]
d[season=="Spring",seasonIntercept:=1]
d[season=="Summer",seasonIntercept:=2]

d[,yearMinus2000:=year-2000]
d[,dayOfSeries:=1:.N]

d[,mu := round(exp(0.1 + yearMinus2000*0.2 + seasonIntercept + 0.2*numberOfCows))]
d[,y:=rpois(.N,mu)]
```

```
dir.create("data")  
  
## Warning in dir.create("data"): 'data' already exists  
fwrite(d,"data/exercise_1.csv")
```

## 8.2 Exercise 2

We are given a dataset containing daily counts of diseases  $y$  from three geographical areas ( $fylke$ ). We want to identify:

- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)
- Does seasonality exist (use the categorical variable “season”)?
- What season has the most cases? (Spring/summer/autumn/winter?)
- Is `numberOfCows` associated with the outcome  $y$ ?

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/exercise\\_2.csv](http://rwhite.no/longitudinal_analysis/data/exercise_2.csv)

```
library(data.table)
set.seed(4)

d <- data.table(date=seq.Date(
  from=as.Date("2010-01-01"),
  to=as.Date("2015-12-31"),
  by=1))

temp <- vector("list",length=3)
for(i in 1:3){
  temp[[i]] <- copy(d)
  temp[[i]][,fylke:=i]
}
d <- rbindlist(temp)

d[,numberOfCows:=rpois(.N,5)]

d[,year:=as.numeric(format.Date(date,"%G"))]
d[,week:=as.numeric(format.Date(date,"%V"))]
d[,month:=as.numeric(format.Date(date,"%m"))]
d[,season:="Winter"]
d[month %in% c(3:5), season:="Spring"]
d[month %in% c(6:8), season:="Summer"]
d[month %in% c(9:11), season:="Autumn"]

d[,seasonIntercept:=0]
d[season=="Spring",seasonIntercept:=1]
d[season=="Summer",seasonIntercept:=2]

d[,yearMinus2000:=year-2000]
d[,dayOfSeries:=1:.N,by=fylke]

d[,mu := round(exp(0.1 + yearMinus2000*0.2 + seasonIntercept + 0.0*numberOfCows + 0.1*(fylke-2)))]
d[,y:=rpois(.N,mu)]
for(i in 1:3) d[fylke==i,y:=round(as.numeric(arima.sim(model=list("ar"=c(0.5)), rand.gen = rpois, n=.N,

## Warning in `[.data.table`(d, fylke == i, `:=`(y,
## round(as.numeric(arima.sim(model = list(ar = c(0.5))), : Coerced double RHS
## to integer to match the type of the target column (column 12 named 'y').
## The RHS values contain no fractions so would be more efficiently created
## as integer. Consider using R's 'L' postfix (typeof(0L) vs typeof(0))
## to create constants as integer and avoid this warning. Wrapping the RHS
## with as.integer() will avoid this warning too but it's better if possible
```

```
## to create the RHS as integer in the first place so that the cost of the
## coercion can be avoided.

## Warning in `[.data.table`(d, fylke == i, `:=`(y,
## round(as.numeric(arima.sim(model = list(ar = c(0.5))), : Coerced double RHS
## to integer to match the type of the target column (column 12 named 'y').
## The RHS values contain no fractions so would be more efficiently created
## as integer. Consider using R's 'L' postfix (typeof(0L) vs typeof(0))
## to create constants as integer and avoid this warning. Wrapping the RHS
## with as.integer() will avoid this warning too but it's better if possible
## to create the RHS as integer in the first place so that the cost of the
## coercion can be avoided.

## Warning in `[.data.table`(d, fylke == i, `:=`(y,
## round(as.numeric(arima.sim(model = list(ar = c(0.5))), : Coerced double RHS
## to integer to match the type of the target column (column 12 named 'y').
## The RHS values contain no fractions so would be more efficiently created
## as integer. Consider using R's 'L' postfix (typeof(0L) vs typeof(0))
## to create constants as integer and avoid this warning. Wrapping the RHS
## with as.integer() will avoid this warning too but it's better if possible
## to create the RHS as integer in the first place so that the cost of the
## coercion can be avoided.

dir.create("data")

## Warning in dir.create("data"): 'data' already exists

fwrite(d,"data/exercise_2.csv")
```

## 8.3 Exercise 3

We are given a dataset containing counts of diseases  $y$  from three geographical areas (`fylke`). We want to identify:

- Is there a general yearly trend (i.e. increasing or decreasing from year to year?)
- Does seasonality exist (use the categorical variable “season”)?
- What season has the most cases? (Spring/summer/autumn/winter?)
- Is `numberOfCows` associated with the outcome  $y$ ?

The data for this chapter is available at: [http://rwhite.no/longitudinal\\_analysis/data/exercise\\_3.csv](http://rwhite.no/longitudinal_analysis/data/exercise_3.csv)

```
library(data.table)
set.seed(4)

d <- data.table(date=seq.Date(
  from=as.Date("2010-01-01"),
  to=as.Date("2015-12-31"),
  by=1))

temp <- vector("list",length=3)
for(i in 1:3){
  temp[[i]] <- copy(d)
  temp[[i]][,fylke:=i]
}
d <- rbindlist(temp)

d[,numberOfCows:=rpois(.N,5)]

d[,year:=as.numeric(format.Date(date,"%G"))]
d[,week:=as.numeric(format.Date(date,"%V"))]
d[,month:=as.numeric(format.Date(date,"%m"))]
d[,season:="Winter"]
d[month %in% c(3:5), season:="Spring"]
d[month %in% c(6:8), season:="Summer"]
d[month %in% c(9:11), season:="Autumn"]

d[,seasonIntercept:=0]
d[season=="Spring",seasonIntercept:=1]
d[season=="Summer",seasonIntercept:=2]

d[,yearMinus2000:=year-2000]

d <- d[sample(1:.N,600)]

d[,mu := round(exp(0.1 + yearMinus2000*0.2 + seasonIntercept + 0.0*numberOfCows + 0.1*(fylke-2)))]
d[,y:=rpois(.N,mu)]

dir.create("data")

## Warning in dir.create("data"): 'data' already exists
fwrite(d,"data/exercise_3.csv")
```





## Chapter 9

# Solutions

### 9.1 Exercise 1

```
library(data.table)
d <- fread("data/exercise_1.csv")

fit0 <- glm(y ~ yearMinus2000 + numberOfCows, data=d, family=poisson())
fit1 <- glm(y ~ season + yearMinus2000 + numberOfCows, data=d, family=poisson())

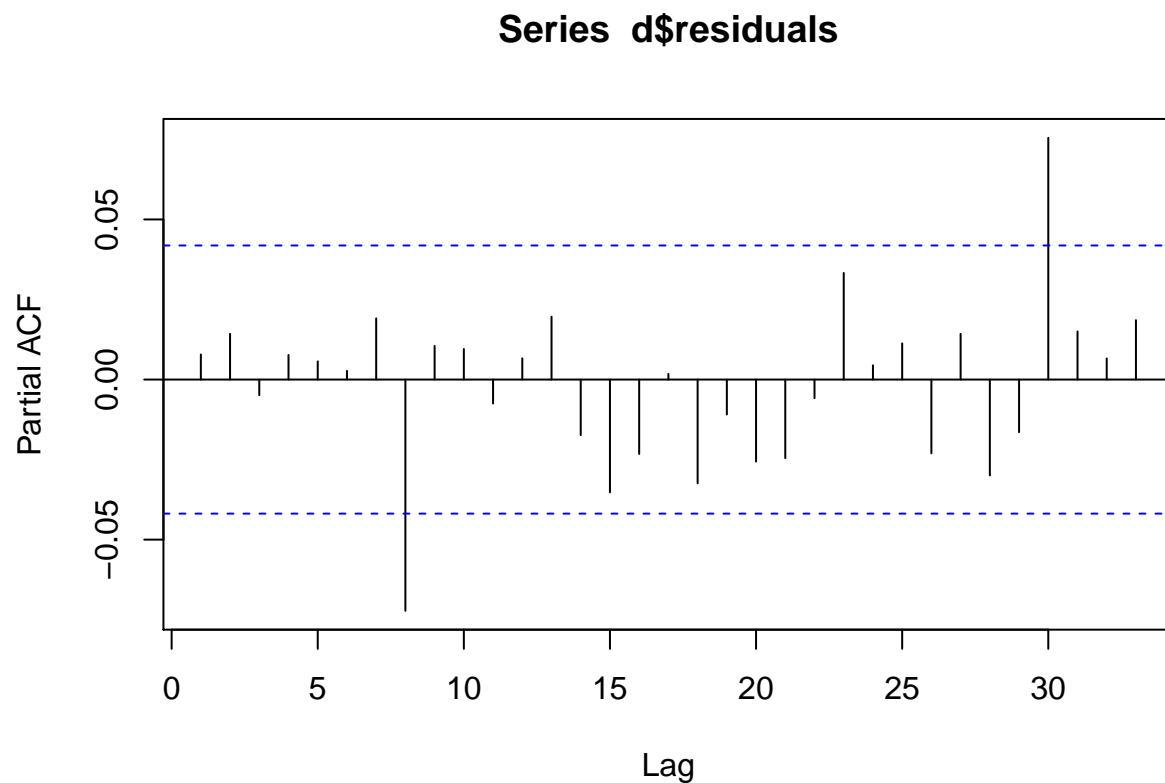
print(lmtest::lrtest(fit0, fit1))

## Likelihood ratio test
##
## Model 1: y ~ yearMinus2000 + numberOfCows
## Model 2: y ~ season + yearMinus2000 + numberOfCows
##   #Df LogLik Df  Chisq Pr(>Chisq)
## 1    3 -104814
## 2    6  -7847  3 193933 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

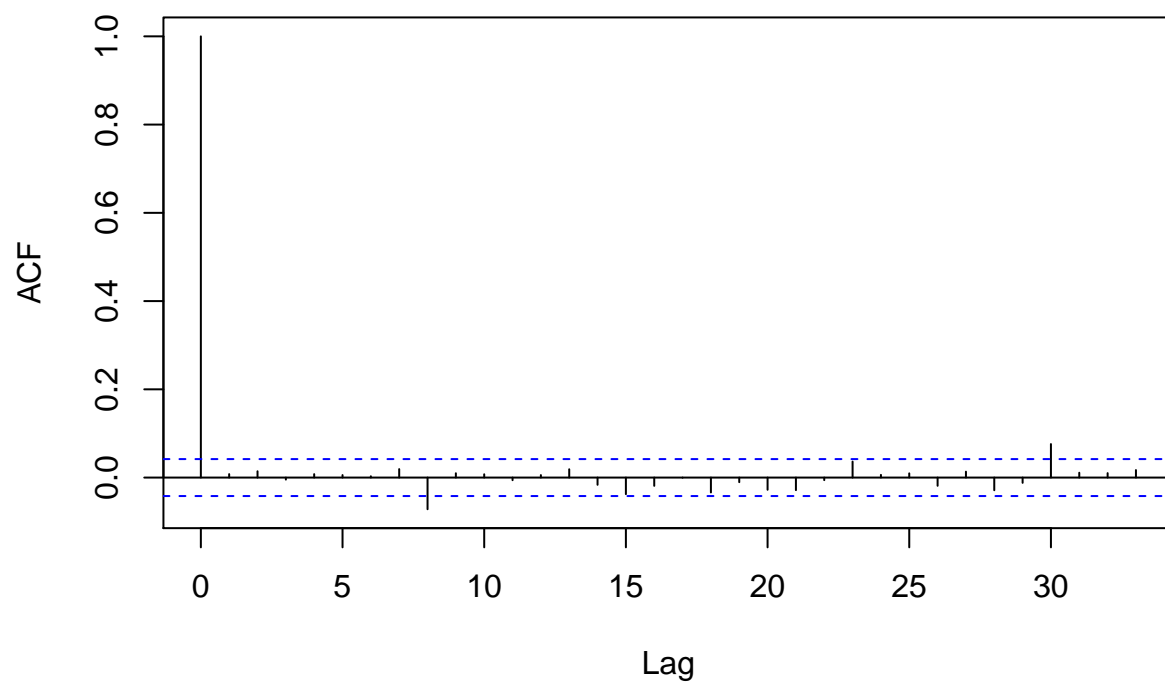
summary(fit1)

##
## Call:
## glm(formula = y ~ season + yearMinus2000 + numberOfCows, family = poisson(),
##      data = d)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.5547  -0.6743  -0.0203   0.6393   3.2527
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.0998769  0.0168980   5.911 3.41e-09 ***
## seasonSpring  0.9996116  0.0077048 129.739 < 2e-16 ***
## seasonSummer  2.0061609  0.0070148 285.990 < 2e-16 ***
## seasonWinter -0.0048955  0.0093124  -0.526  0.599
```

```
## yearMinus2000  0.2001843  0.0011420 175.298 < 2e-16 ***
## numberOfCows  0.1987005  0.0007667 259.153 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
## Null deviance: 296600.3 on 2190 degrees of freedom
## Residual deviance: 2167.4 on 2185 degrees of freedom
## AIC: 15707
##
## Number of Fisher Scoring iterations: 4
d[,residuals:=residuals(fit1, type = "response")]
pacf(d$residuals)
```



```
acf(d$residuals)
```

**Series d\$residuals**

## 9.2 Exercise 2

```

library(data.table)
d <- fread("data/exercise_2.csv")

fit0 <- MASS::glmmPQL(y~yearMinus2000 + numberOfCows, random = ~ 1 | fylke,
                      family = poisson, data = d)

## iteration 1
## iteration 2

fit1 <- MASS::glmmPQL(y~season + yearMinus2000 + numberOfCows, random = ~ 1 | fylke,
                      family = poisson, data = d)

## iteration 1
## iteration 2
## iteration 3

print(lmtest::lrtest(fit0, fit1))

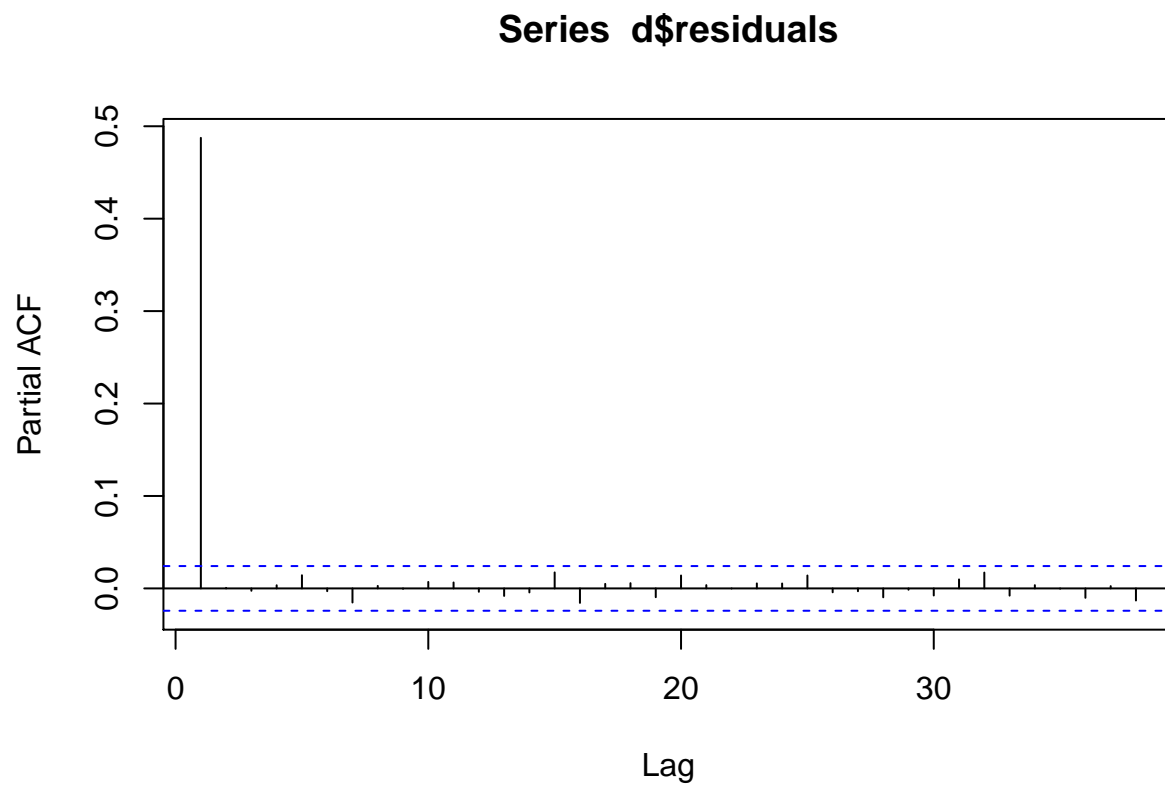
## Likelihood ratio test
##
## Model 1: y ~ yearMinus2000 + numberOfCows
## Model 2: y ~ season + yearMinus2000 + numberOfCows
##   #Df LogLik Df Chisq Pr(>Chisq)
## 1    5
## 2    8      3

summary(fit1)

## Linear mixed-effects model fit by maximum likelihood
## Data: d
##   AIC BIC logLik
##   NA  NA   NA
##
## Random effects:
## Formula: ~1 | fylke
##      (Intercept) Residual
## StdDev:  0.08342256 1.298934
##
## Variance function:
## Structure: fixed weights
## Formula: ~invwt
## Fixed effects: y ~ season + yearMinus2000 + numberOfCows
##              Value Std.Error DF t-value p-value
## (Intercept)  0.8483946 0.05053613 6565  16.78788  0.0000
## seasonSpring  0.9334080 0.00685147 6565 136.23480  0.0000
## seasonSummer  1.9312703 0.00621739 6565 310.62400  0.0000
## seasonWinter -0.0822382 0.00841368 6565  -9.77434  0.0000
## yearMinus2000 0.2004222 0.00104237 6565 192.27503  0.0000
## numberOfCows  0.0005788 0.00077223 6565   0.74954  0.4536
## Correlation:
##              (Intr) ssnSpr ssnSmm ssnWnt yM2000
## seasonSpring -0.097

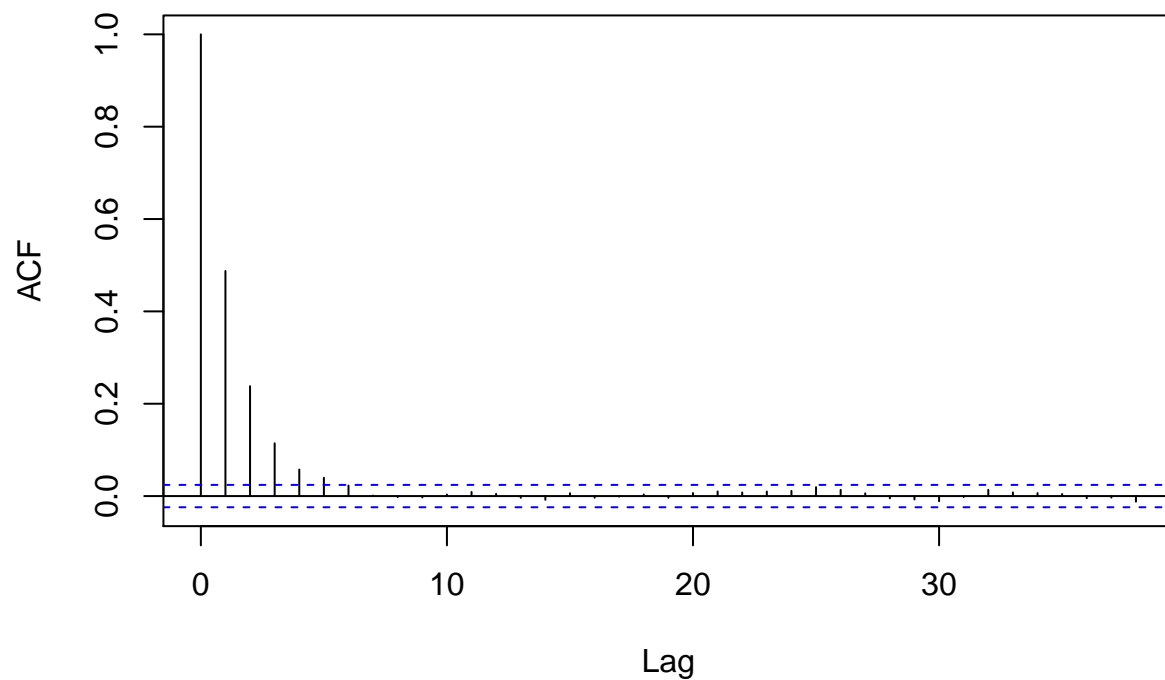
```

```
## seasonSummer -0.106 0.793
## seasonWinter -0.079 0.586 0.646
## yearMinus2000 -0.268 0.000 0.000 -0.002
## numberOfCows -0.070 -0.002 -0.018 0.004 -0.020
##
## Standardized Within-Group Residuals:
##      Min      Q1      Med      Q3      Max
## -5.44473503 -0.49365704 -0.05256441 0.39697143 16.32534219
##
## Number of Observations: 6573
## Number of Groups: 3
d[,residuals:=residuals(fit1, type = "normalized")]
pacf(d$residuals)
```



```
acf(d$residuals)
```

### Series d\$residuals



```
fit1 <- MASS::glmPQL(y~season + yearMinus2000 + numberOfCows, random = ~ 1 | fylke,
  family = poisson, data = d,
  correlation=nlme::corAR1(form=~dayOfSeries|fylke))
```

```
## iteration 1
```

```
## iteration 2
```

```
## iteration 3
```

```
summary(fit1)
```

```
## Linear mixed-effects model fit by maximum likelihood
```

```
## Data: d
```

```
## AIC BIC logLik
```

```
## NA NA NA
```

```
##
```

```
## Random effects:
```

```
## Formula: ~1 | fylke
```

```
## (Intercept) Residual
```

```
## StdDev: 0.08328798 1.319938
```

```
##
```

```
## Correlation Structure: AR(1)
```

```
## Formula: ~dayOfSeries | fylke
```

```
## Parameter estimate(s):
```

```
## Phi
```

```
## 0.5525116
```

```
## Variance function:
```

```
## Structure: fixed weights
## Formula: ~invwt
## Fixed effects: y ~ season + yearMinus2000 + numberOfCows
##
```

	Value	Std.Error	DF	t-value	p-value
## (Intercept)	0.9283222	0.05561940	6565	16.69062	0.000
## seasonSpring	0.8631442	0.01224757	6565	70.47476	0.000
## seasonSummer	1.8166993	0.01098229	6565	165.42086	0.000
## seasonWinter	-0.1394364	0.01488823	6565	-9.36554	0.000
## yearMinus2000	0.2001812	0.00197415	6565	101.40142	0.000
## numberOfCows	0.0004206	0.00057695	6565	0.72909	0.466

```
## Correlation:
##
```

	(Intr)	ssnSpr	ssnSmm	ssnWnt	yM2000
## seasonSpring	-0.155				
## seasonSummer	-0.171	0.784			
## seasonWinter	-0.123	0.574	0.621		
## yearMinus2000	-0.464	0.000	0.000	-0.002	
## numberOfCows	-0.049	0.001	-0.006	0.004	-0.007

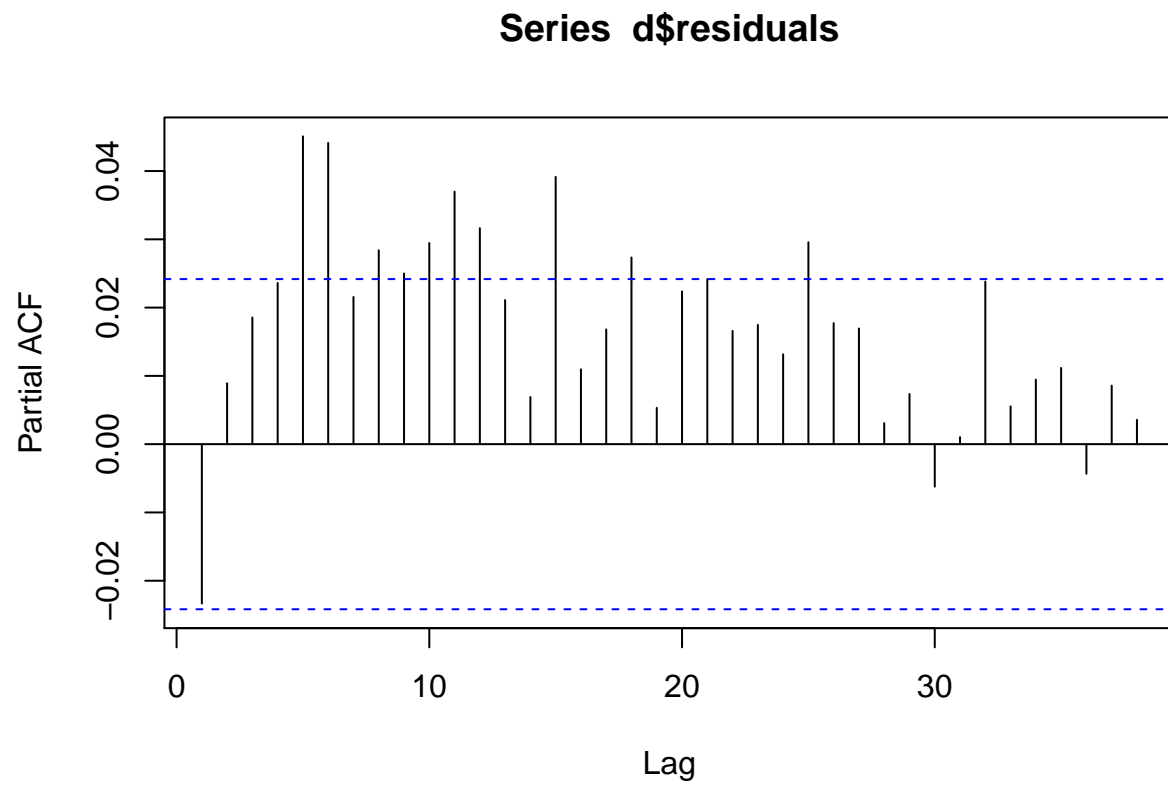
```
##
## Standardized Within-Group Residuals:
##
```

	Min	Q1	Med	Q3	Max
##	-5.03056012	-0.54478730	-0.04721577	0.46011628	15.05853958

```
##
## Number of Observations: 6573
## Number of Groups: 3

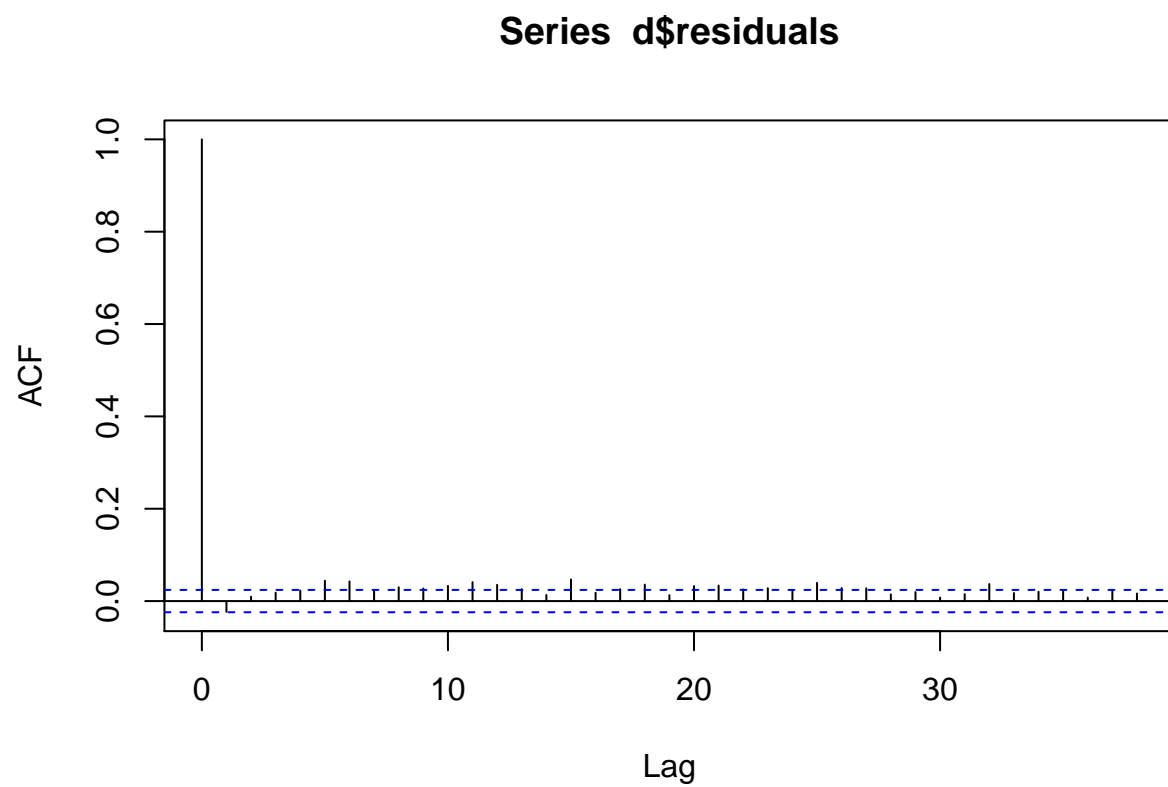
d[,residuals:=residuals(fit1, type = "normalized")]

pacf(d$residuals)
```



```
acf(d$residuals)
```





### 9.3 Exercise 3

```

library(data.table)
d <- fread("data/exercise_3.csv")

fit0 <- lme4::glmer(y ~ yearMinus2000 + numberOfCows + (1|fylke), family = poisson, data = d)
fit1 <- lme4::glmer(y ~ season + yearMinus2000 + numberOfCows + (1|fylke), family = poisson, data = d)

## Warning in checkConv(attr(opt, "derivs"), opt$par, ctrl = control
## $checkConv, : Model failed to converge with max|grad| = 0.0013139 (tol =
## 0.001, component 1)

print(lmtest::lrtest(fit0, fit1))

## Likelihood ratio test
##
## Model 1: y ~ yearMinus2000 + numberOfCows + (1 | fylke)
## Model 2: y ~ season + yearMinus2000 + numberOfCows + (1 | fylke)
##   #Df   LogLik Df Chisq Pr(>Chisq)
## 1    4 -10144.4
## 2    7 -1794.9  3 16699  < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

summary(fit1)

## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: poisson ( log )
## Formula: y ~ season + yearMinus2000 + numberOfCows + (1 | fylke)
## Data: d
##
##      AIC      BIC   logLik deviance df.resid
## 3603.8   3634.6  -1794.9   3589.8     593
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -3.4683 -0.6176 -0.0064  0.5895  3.1431
##
## Random effects:
## Groups Name          Variance Std.Dev.
## fylke (Intercept) 0.005643 0.07512
## Number of obs: 600, groups: fylke, 3
##
## Fixed effects:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.0817819  0.0707498   1.156   0.248
## seasonSpring  1.0213789  0.0246508  41.434 <2e-16 ***
## seasonSummer  2.0118660  0.0220035  91.434 <2e-16 ***
## seasonWinter -0.0001082  0.0294244  -0.004   0.997
## yearMinus2000 0.2019749  0.0038745  52.129 <2e-16 ***
## numberOfCows -0.0045764  0.0028115  -1.628   0.104
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##

```

```
## Correlation of Fixed Effects:
##          (Intr) ssnSpr ssnSmm ssnWnt yM2000
## seasonSprng -0.225
## seasonSummr -0.276  0.756
## seasonWintr -0.200  0.565  0.633
## yearMns2000 -0.708 -0.020  0.008 -0.011
## numberOfCws -0.200  0.025  0.038  0.053 -0.005
## convergence code: 0
## Model failed to converge with max|grad| = 0.0013139 (tol = 0.001, component 1)
```