

# Time series regression model

## Multiple linear regression: potential problems

When we fit a linear regression model to a particular data set, many problems may occur.

Most common among these are the following:

- ▶ Non-linearity of the response-predictor relationships
- ▶ Correlation of error terms
- ▶ Non-constant variance of error terms
- ▶ Outliers

we will discuss some of these problems in more detail . . .

## Multiple linear regression with time series

Many business and economic problems involve the use of time series data.

The linear regression model may be usefully employed to model monthly, quarterly or yearly data.

- ▶ A linear trend may be easily included through a predictor  $X_{1,t} = t$ .
- ▶ Seasonality may modeled with seasonal dummy variables. As a general rule, we use  $s - 1$  dummy variables to describe  $s$  periods (to avoid perfect multicollinearity).

## Multiple linear regression with time series

For instance, a model for quarterly data with trend and seasonality may be

$$Y_t = \beta_0 + \beta_1 t + \beta_2 S_2 + \beta_3 S_3 + \beta_4 S_4 + \varepsilon_t$$

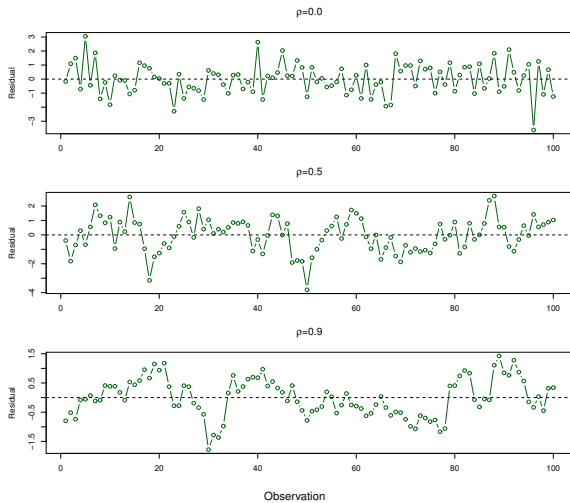
Trend and seasonality are modelled as a series of straight lines with different intercept and same slope. The first quarter is described with the model  $Y_t = \beta_0 + \beta_1 t$ .

Parameters  $\beta_2, \beta_3, \beta_4$  describe the variation with respect to  $\beta_0$  due to seasonality.

## Multiple linear regression with time series

- ▶ Time series data tend to be **autocorrelated**
- ▶ Autocorrelation occurs when the effect of a variable is spread over time. For example, a change in prices may have an effect on both current and future sales
- ▶ Autocorrelation may be detected through a **graphical inspection of residuals**
- ▶ Specific tests on residuals

# Autocorrelated residuals



# Autocorrelated residuals

A typical example of autocorrelation is defined as

$$Y_t = \beta_0 + \beta_1 X_t + \varepsilon_t$$

with

$$\varepsilon_t = \rho\varepsilon_{t-1} + \nu_t$$

where  $\rho$  is the correlation between sequential errors and  $\nu_t$  is an erratic component with mean zero and constant variance.

If  $\rho = 0$  then  $\varepsilon_t = \nu_t$ .

The **Durbin-Watson test** is typically used to diagnose this kind of autocorrelation.

The system of hypothesis is

$$H_0 : \rho = 0 \quad H_1 : \rho > 0$$

## Durbin-Watson test

The Durbin-Watson test is defined as

$$DW = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}$$

The values of DW range between 0 and 4 with a central value of 2. For large samples, the following holds

$$DW = 2(1 - r_1(e))$$

where  $r_1(e)$  is the residual autocorrelation at lag 1. Since  $-1 < r_1(e) < 1$ , then  $0 < DW < 4$ .



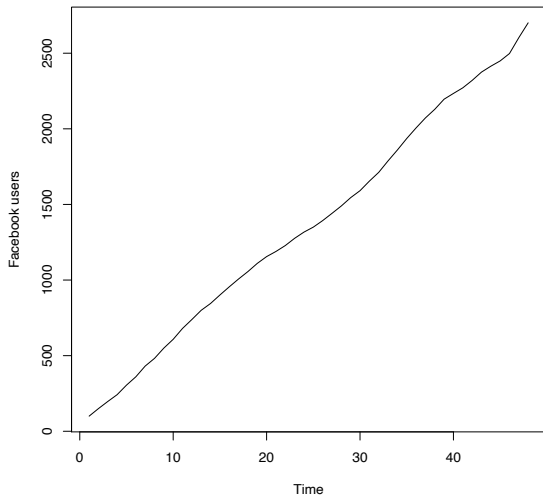
## Autocorrelation: solutions

To solve the problem of autocorrelation we need to examine the model:

- ▶ is the functional form correct?
- ▶ are there any omitted variables?

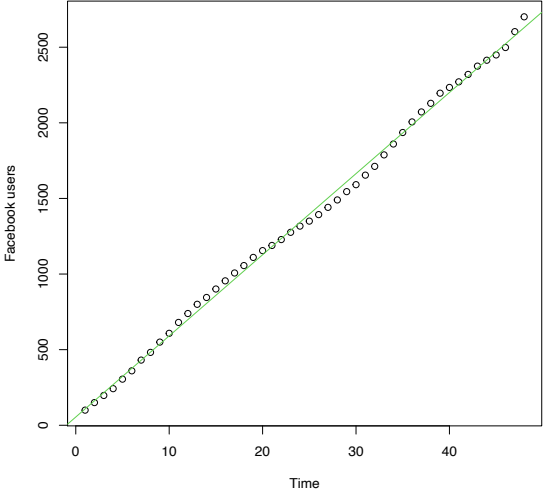
## Example

Facebook users: quarterly data 2008-2020



# Example

## Facebook users: simple linear regression



# Example

## Facebook users: simple linear regression

```
lm(formula = fb ~ time)
```

Coefficients:

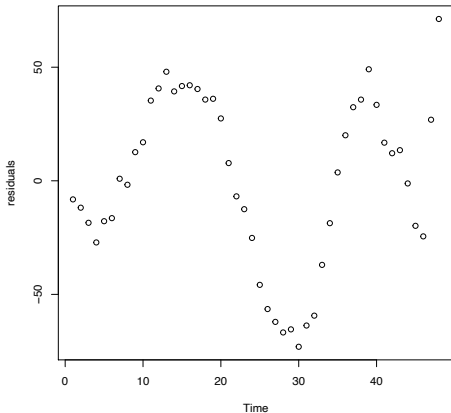
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	54.5363	10.9917	4.962	1e-05 ***
time	53.6507	0.3905	137.378	<2e-16 ***

---

Residual standard error: 37.48 on 46 degrees of freedom  
Multiple R-squared: 0.9976, Adjusted R-squared: 0.9975  
F-statistic: 1.887e+04 on 1 and 46 DF, p-value: < 2.2e-16

# Example

## Facebook users: residuals

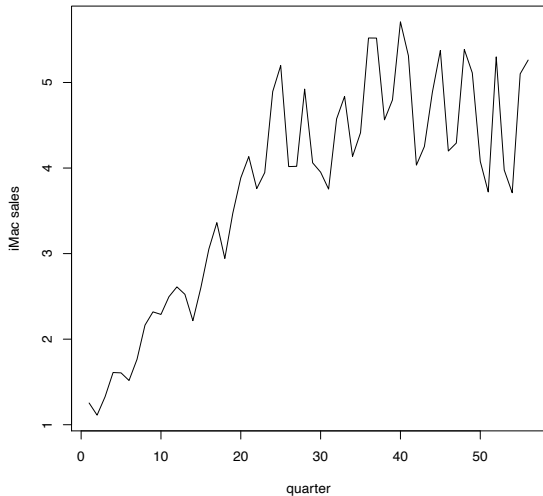


Durbin-Watson test:  $DW = 0.16378$ ,  $p\text{-value} < 2.2e-16$

Positive autocorrelation in residuals

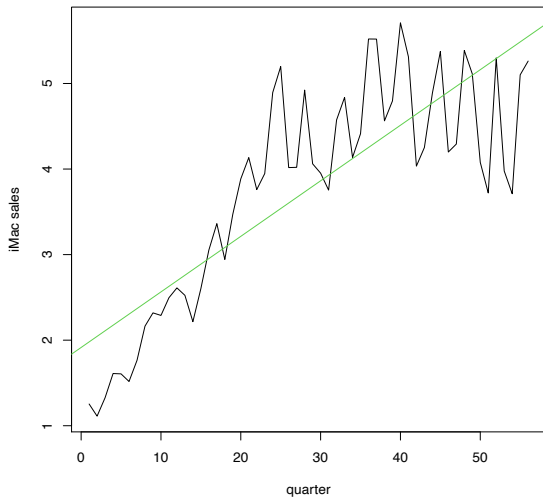
# Example

iMac sales: quarterly data 2006-2019



# Example

iMac sales: simple linear regression



## Example

iMac sales: linear regression with trend and seasonality

Call:

```
tslm(formula = mac.ts ~ trend + season)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.60158	-0.42293	-0.00687	0.54972	1.42797

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	2.155255	0.236078	9.129	2.62e-12	***
trend	0.064591	0.005613	11.507	8.68e-16	***
season2	-0.640448	0.256052	-2.501	0.0156	*
season3	-0.460039	0.256237	-1.795	0.0785	.
season4	0.176727	0.256544	0.689	0.4940	
---					

Residual standard error: 0.6773 on 51 degrees of freedom

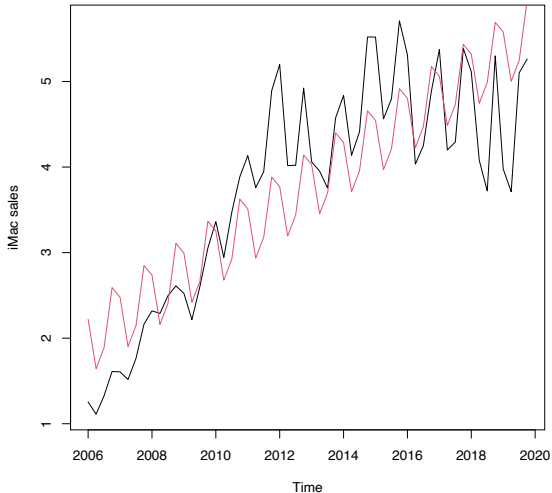
Multiple R-squared: 0.7436, Adjusted R-squared: 0.7235

F-statistic: 36.97 on 4 and 51 DF, p-value: 1.695e-14



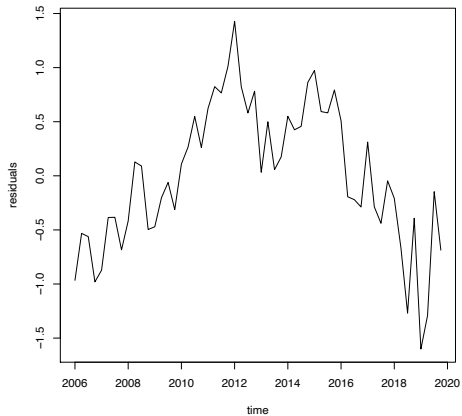
# Example

iMac sales: linear regression with trend and seasonality



# Example

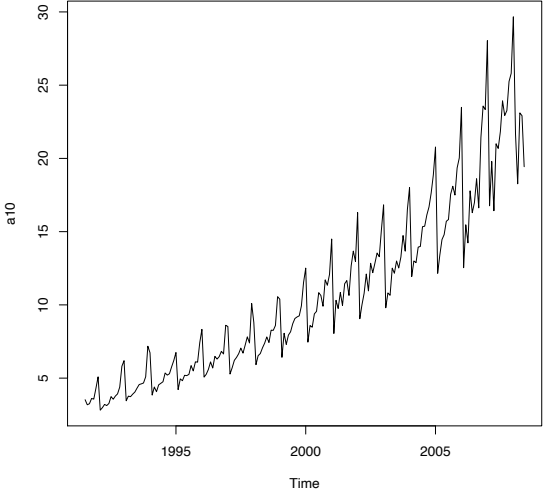
## iMac sales: residuals



Residuals clearly show a nonlinear behaviour

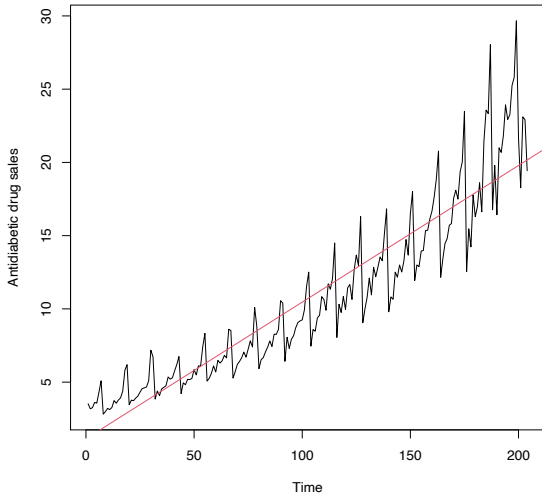
# Example

## Monthly sales of a drug



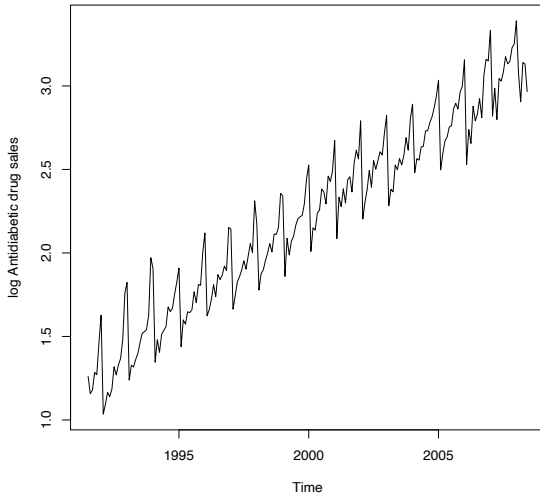
# Example

Monthly sales of a drug: simple linear regression



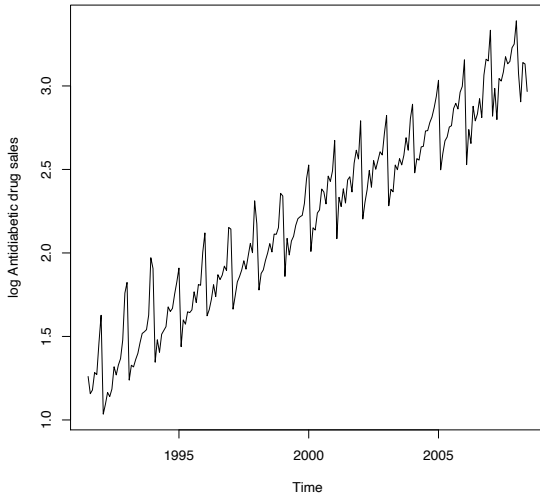
# Example

Monthly sales of a drug: log transformation



# Example

Monthly sales of a drug: log transformation



## Example

Monthly sales of a drug: simple linear regression with log transformation

Call:

```
lm(formula = la10 ~ t)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-0.36954	-0.09621	-0.00889	0.07139	0.43395

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.2577135	0.0216920	57.98	<2e-16 ***
t	0.0093211	0.0001835	50.80	<2e-16 ***

---

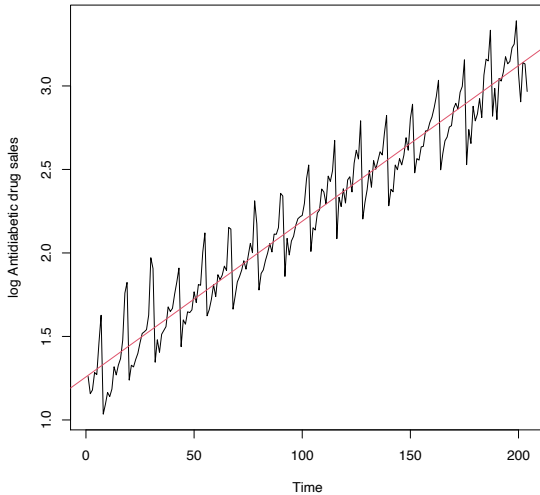
Residual standard error: 0.1543 on 202 degrees of freedom

Multiple R-squared: 0.9274, Adjusted R-squared: 0.927

F-statistic: 2580 on 1 and 202 DF, p-value: < 2.2e-16

# Example

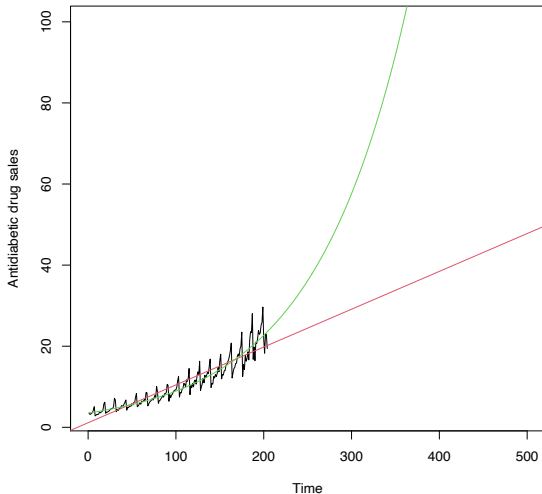
## Monthly sales of a drug: log transformation





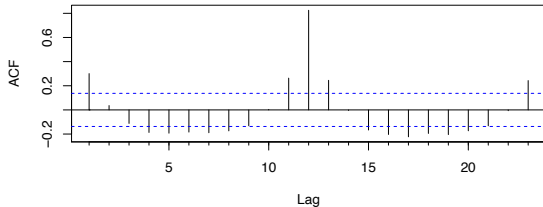
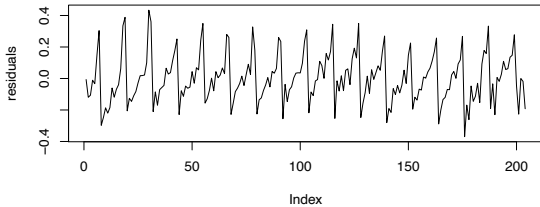
# Example

## Monthly sales of a drug: model comparison



# Example

## Monthly sales of a drug: residuals



## Selecting predictors

- ▶ When there are many possible predictors, we need some strategy for **selecting the best predictors** to use in a regression model
- ▶ We may use different approaches for model selection

## Selecting predictors

- ▶ **Best subset regression**: suitable when possible
- ▶ **Stepwise regression**: backward and forward, or hybrid approach
- ▶ **Akaike's Information Criterion**

$$\text{AIC} = T \log \left( \frac{\text{SSE}}{T} \right) + 2(k + 2)$$

The idea is to penalize the fit of the model (SSE) with the number of parameters that need to be estimated.

The model with the minimum AIC is often the best model for forecasting.

## Forecasting with regression

Predictions for  $Y_t$  can be obtained using

$$\hat{Y}_t = \hat{\beta}_0 + \hat{\beta}_1 X_{1,t} + \hat{\beta}_2 X_{2,t} + \cdots + \hat{\beta}_k X_{k,t}$$

However, we are interested in **forecasting future values** of  $Y$ .

## Ex-ante forecasts and ex-post forecasts

- ▶ **Ex-ante forecasts** are those made using only the information available in advance: genuine forecasts
- ▶ **Ex-post forecasts** are those that are made using later information on the predictors, i.e. once these have been observed.
- ▶ **Building a predictive regression model**: obtaining forecasts of the predictors can be very challenging. An alternative formulation is to use as predictors their **lagged values**.

$$Y_{t+1} = \beta_0 + \beta_1 X_{1,t} + \beta_2 X_{2,t} + \cdots + \beta_k X_{k,t} + \varepsilon_{t+1}$$