

Bus 274: Further Statistics For Business

Harald Schmidbauer



About These Slides

- The present slides are not self-contained; they need to be explained and discussed. This will be done in the lectures.
- Even though being a “work in progress” and subject to revision, the slides constitute copyrighted material.
If you want to reproduce or copy anything from the slides, please ask:

Harald Schmidbauer **harald** at **hs-stat** dot **com**
Angi Rösch **angi** at **angi-stat** dot **com**

- The slides were produced using \LaTeX (www.latex-project.org) and R (www.R-project.org) on a GNU/Linux system, all of which are free and open source software (FOSS).
- R files used for this course are available upon request.



Chapter 16:

Multiple Regression



16.1 Introduction

SLR and Multiple Linear Regression.

- Goal of SLR:

Explain the variability in Y , using a variable X .

- Goal of multiple linear regression:

Explain the variability in Y , using a set of variables X_1, X_2, \dots, X_k .



16.1 Introduction

The problem.

Given are points $(x_{1i}, x_{2i}, \dots, x_{ki}, y_i)$, where:

- y_i : observations from a variable Y , the dependent variable;
- x_{ji} : observations from a variable X_j , which is an independent variable.

Given a $(k+1)$ -dimensional cloud of points, how can we fit a hyperplane?



16.1 Introduction

Outlook on Chapter 16.

- 16.2 An Intuitive Approach
three-dimensional scatterplots and a regression plane
- 16.3 The Regression Plane
the method of least squares
- 16.4 Explanatory Power of the Model
decomposition of variance; coefficient of determination
- 16.5 A Stochastic Model of Multiple Regression
stochastic model and statistical inference
- 16.6 Examples
- 16.7 Prediction Based on Multiple Regression
point prediction and prediction intervals



16.2 An Intuitive Approach

The case of three variables: X_1, X_2, Y .

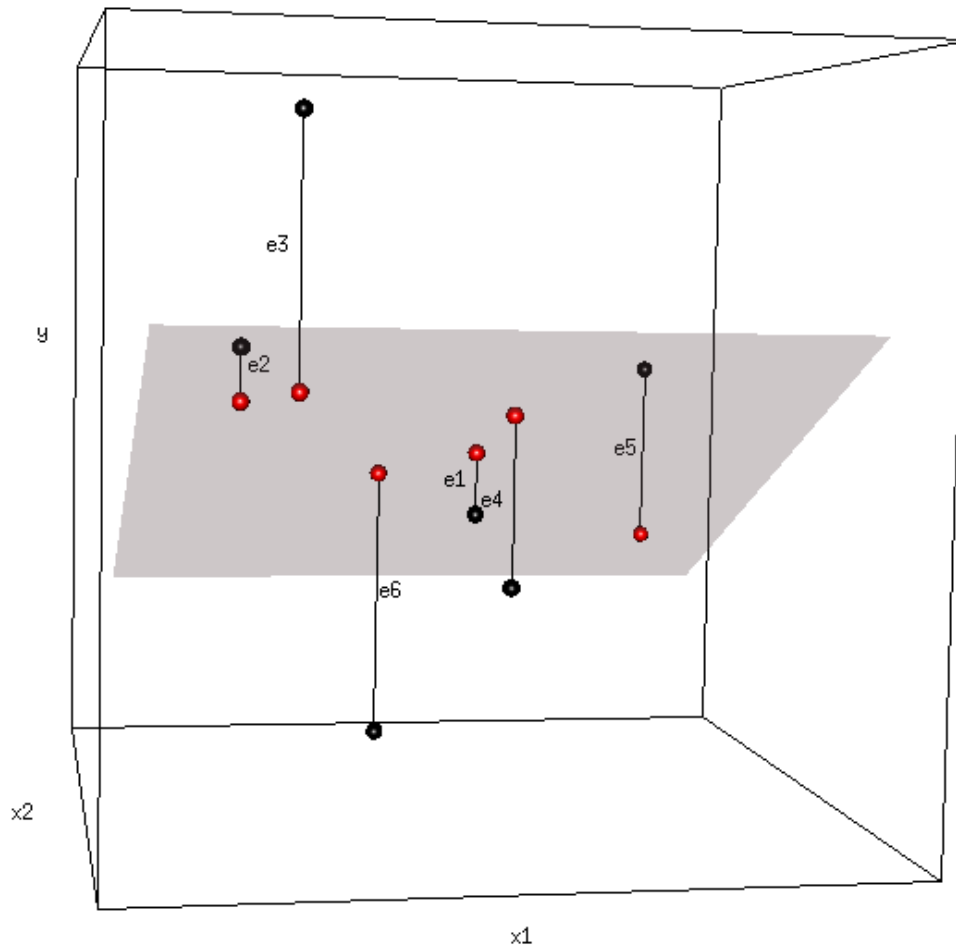
We shall now see a three-dimensional scatterplot in two perspectives with:

- black points, representing the observations,
- a plane, which somehow fits these points,
- red points, the projection of the black points onto the plane,
- the distance between the black and the red points.



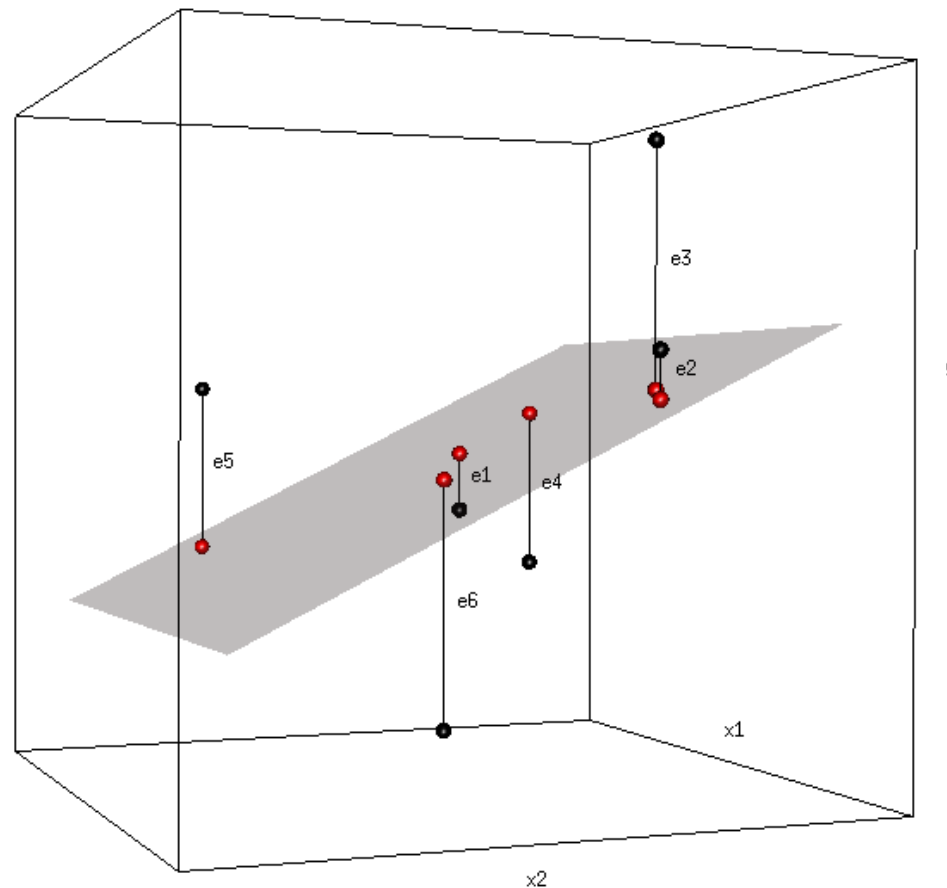
16.2 An Intuitive Approach

Observed points and their projections onto the plane.



16.2 An Intuitive Approach

Observed points and their projections onto the plane.



16.2 An Intuitive Approach

How to find that plane. . . .

in order to find a “good” plane to represent the cloud of points, we need:

- the equation of a plane, depending on parameters,
- a distance function,
- to find the parameter values such that the distance function is minimized.



16.3 The Regression Plane

A plane and the observations.

- Plane in 3-dimensional space: $y = a + b_1x_1 + b_2x_2$
- With observations (x_{1i}, x_{2i}, y_i) , $i = 1, \dots, n$:

$$\begin{array}{ll} \hat{y}_1 = a + b_1x_{11} + b_2x_{21}, & e_1 = y_1 - \hat{y}_1 \\ \hat{y}_2 = a + b_1x_{12} + b_2x_{22}, & e_2 = y_2 - \hat{y}_2 \\ \vdots & \vdots \\ \hat{y}_n = a + b_1x_{1n} + b_2x_{2n}, & e_n = y_n - \hat{y}_n \end{array}$$

- The \hat{y}_i are called the fitted values.



16.3 The Regression Plane

Using matrices. — The last relations can be written as:

$$\hat{y} = Xb, \quad e = y - \hat{y} = y - Xb,$$

where

$$\hat{y} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{pmatrix}, \quad X = \begin{pmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} \end{pmatrix}, \quad b = \begin{pmatrix} a \\ b_1 \\ b_2 \end{pmatrix},$$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad e = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}.$$



16.3 The Regression Plane

Definition.

- Define $\hat{y}_i = a + b_1x_{1i} + b_2x_{2i}$ and $e_i = y_i - \hat{y}_i$.
- The regression plane of Y with respect to X_1 and X_2 is the plane $y = a + b_1x_1 + b_2x_2$ with a , b_1 and b_2 such that

$$\begin{aligned} Q(a, b_1, b_2) &= \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - a - b_1x_{1i} - b_2x_{2i})^2 \end{aligned}$$

attains its minimum.

- b_1 and b_2 : regression coefficients.



16.3 The Regression Plane

Regression: some first comments.

- This procedure is asymmetric — like SLR!
- It conforms to the idea: Given X_1 and X_2 , what is Y ?
- X_1, X_2 : “independent variables”,
 Y : “dependent variable”
- This procedure can be easily generalized to $k > 2$ independent variables.
- The case $k > 2$ cannot be easily visualized in terms of a scatterplot.



16.3 The Regression Plane

Example: Used cars.

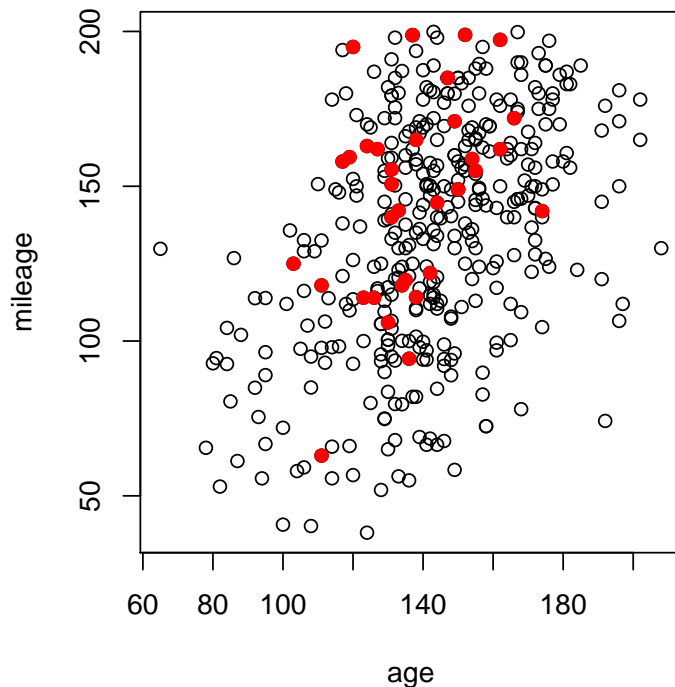
- For a set of used cars, consider these variables:
 - mileage (km)
 - age (months)
 - price (€)
- A natural choice is:
 - dependent variable: price
 - independent variables: mileage, age



16.3 The Regression Plane

Example: Used cars.

- Important: The so-called “independent variables” need not be uncorrelated.
- For our sample of 400 cars (VW Golf 1.8):



– correlation: 0.43

– red points: cars with ac



16.3 The Regression Plane

Computing the regression plane.

- Minimizing Q leads to the following vector equation:

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- The fitted values are:

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{b} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- These formulas apply to any number k of independent variables.
- For $k = 1$, the formulas of SLR are obtained.



16.3 The Regression Plane

Multiple regression — some properties in the context of descriptive statistics.

- The vector of arithmetic means $(\bar{x}_1, \bar{x}_2, \bar{y})$ is on the regression plane.
- The average error \bar{e} equals zero.
- The matrix $X(X'X)^{-1}X'$ in $\hat{y} = Xb = X(X'X)^{-1}X'y$ is a projection matrix: y is projected onto a sub-space of \mathbb{R}^n .



16.3 The Regression Plane

Example: Used cars.

- Data from 400 used cars (VW Golf 1.8, age at least 5 years, mileage at most 200000 km).
- The fitted regression plane is:

$$\text{price} = 14146.2 - 24.61 \cdot \text{mileage} - 49.13 \cdot \text{age}$$

(Price in €, mileage in 1000 km, age in months.)

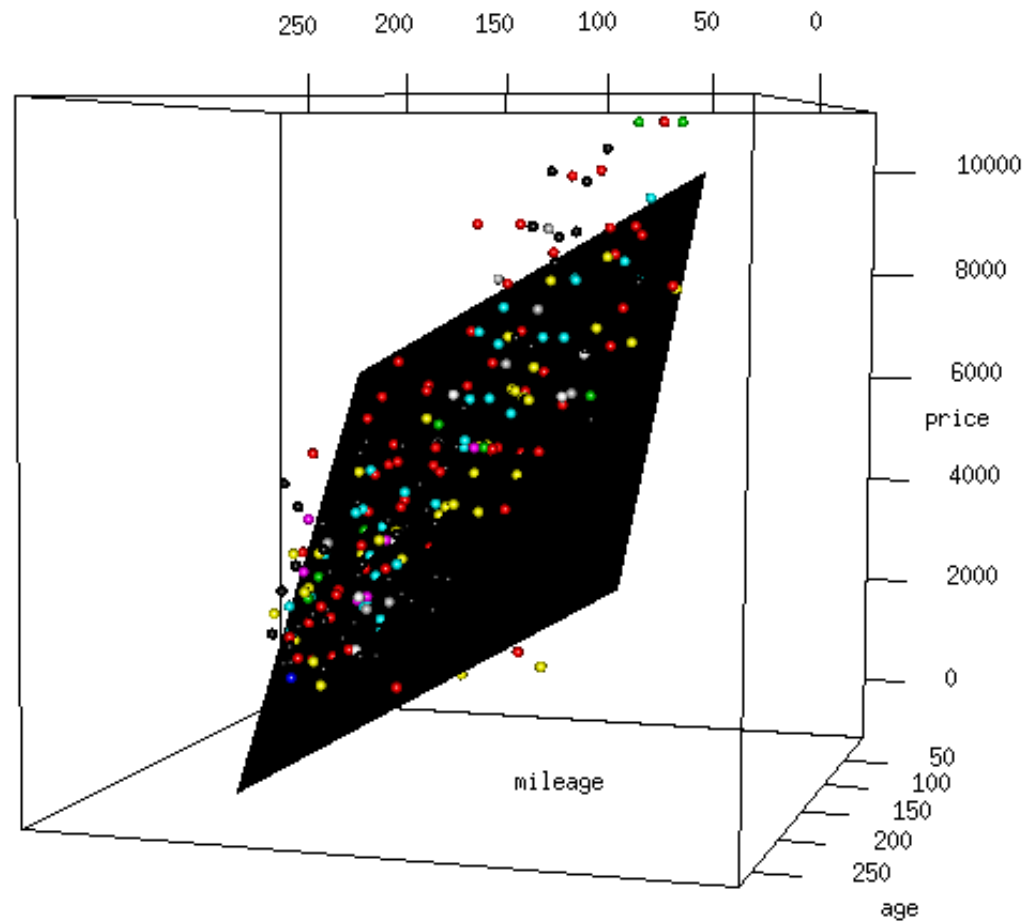
- According to this result: What is the average price of a car with mileage 100000 km, age 10 years?
- How much will this decrease if the car is used for another year, for another 12000 km?



16.3 The Regression Plane

Example: Used cars.

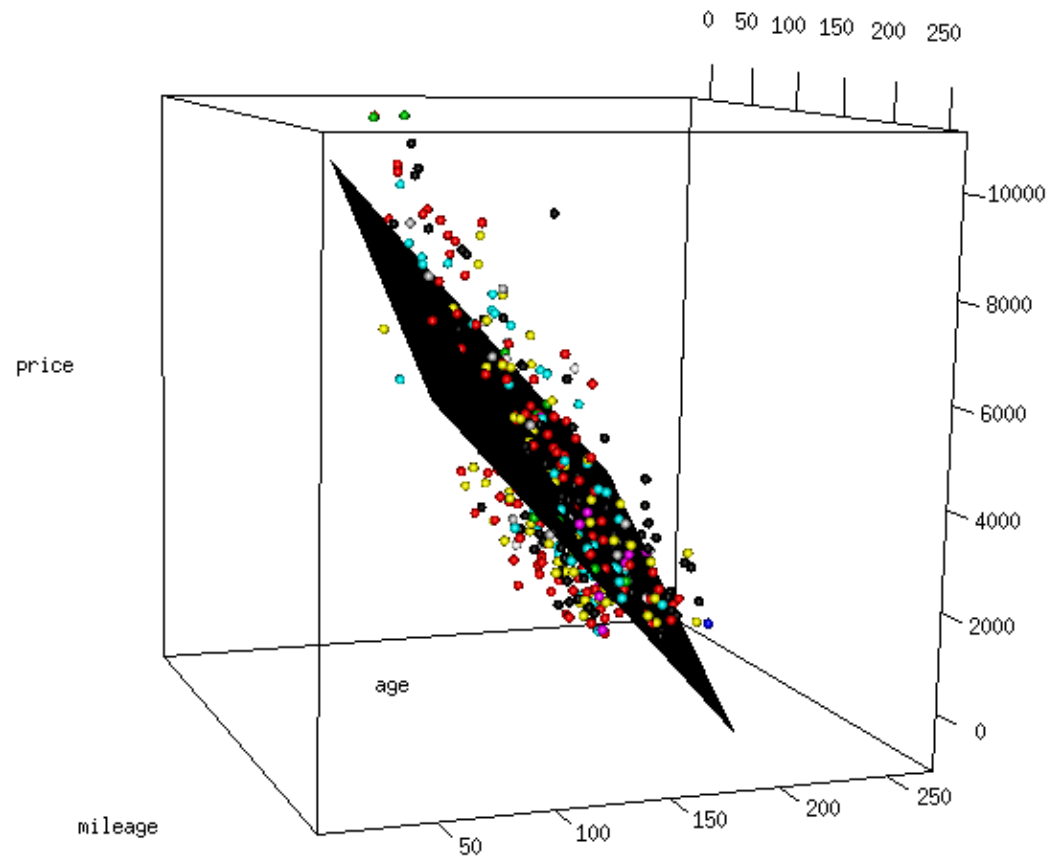
Scatterplot:



16.3 The Regression Plane

Example: Used cars.

Scatterplot:



16.4 Explanatory Power of the Model

Decomposition of variance.

As in SLR, it holds that:

$$\begin{aligned} \sum (y_i - \bar{y})^2 &= \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2, \\ \text{SST} &= \text{SSR} + \text{SSE} \end{aligned}$$

where

SST: total sum of squares

SSR: regression sum of squares

SSE: error sum of squares



16.4 Explanatory Power of the Model

The coefficient of determination.

It is defined as:

$$\frac{SSR}{SST}$$

- The coefficient of determination is the share of variability in the data which is explained by the regression.
- In contrast to SLR, the coefficient of determination cannot be computed as the square of a coefficient of correlation.
- $R^2 = 100\%$ if and only if all observed points are on the regression plane.
- $R^2 = 0\%$ means that no linear combination of independent variables contributes to explaining Y .



16.4 Explanatory Power of the Model

Example: Used cars.

Compare the following fitted models and their R^2 s:

- Model 1 ($R^2 = 0.434$):
price = $8984.41 - 38.20 \cdot \text{mileage}$
- Model 2 ($R^2 = 0.528$):
price = $13160.68 - 65.61 \cdot \text{age}$
- Model 3 ($R^2 = 0.675$):
price = $14146.2 - 24.61 \cdot \text{mileage} - 49.13 \cdot \text{age}$
- According to each model: What is the average price of a car with mileage 100000 km, age 10 years?



16.5 A Stochastic MLR Model

Multiple regression in descriptive and inductive statistics.

- So far, we have seen multiple regression from a purely *descriptive* point of view.
(There were no probabilities, no stochastic models.)
- A stochastic model is needed to
 - obtain insight into the mechanism which created the data,
 - make reliable statements about out-of-sample cases.
- We shall now see this model, written out for $k = 2$ independent variables.



16.5 A Stochastic MLR Model

A *stochastic* multiple linear regression model.

$$Y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i, \quad i = 1, \dots, n$$

- The random variable Y_i represents the observation belonging to x_{1i} and x_{2i} .
- α , β_1 and β_2 are unknown parameters (to be estimated).
- x_{ji} is the observation of the independent variable X_j .
- ϵ_i is a random variable; it contains everything not accounted for in the equation $y = \alpha + \beta_1 x_1 + \beta_2 x_2$.



16.5 A Stochastic MLR Model

Matrix form of the stochastic model.

The system

$$Y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i, \quad i = 1, \dots, n,$$

can be written as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{2n} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}.$$

The generalization to k independent variables is straightforward.



16.5 A Stochastic MLR Model

Assumptions in the stochastic multiple linear regression model.

For statistical inference, we assume:

- The matrix X has full rank.
- The matrix X is considered fixed (non-stochastic).
- $\epsilon_i \sim N(0, \sigma_\epsilon^2)$ iid for $i = 1, \dots, n$.

With the last assumption, it holds that

$$E(Y_i | x_1, x_2) = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i}, \quad i = 1, \dots, n.$$



16.5 A Stochastic MLR Model

Computing estimators.

- The method of least squares leads to the following estimator for β :

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

- As a random vector, $\hat{\beta}$ has a covariance matrix. It is given by

$$\text{var}(\hat{\beta}) = \sigma_{\epsilon}^2 \cdot (\mathbf{X}'\mathbf{X})^{-1}.$$

- The residual error variance can be estimated as

$$s_{\epsilon}^2 = \frac{\text{SSE}}{n - k - 1}$$



16.5 A Stochastic MLR Model

Statistical inference about the parameters.

- Statistical inference about β_j is based on the following property:

$$\frac{\hat{\beta}_j - \beta_j}{s_{\beta_j}} \sim t_{n-k-1},$$

where s_{β_j} is the standard error of $\hat{\beta}_j$.

- The standard error s_{β_j} can be obtained from

$$\hat{\text{var}}(\hat{\beta}) = s_{\epsilon}^2 \cdot (\mathbf{X}'\mathbf{X})^{-1}.$$

(This may be tedious to compute, but it is standard output in statistical software packages.)



16.5 A Stochastic MLR Model

Which variables to include?

- We prefer models with large R^2 and small s_ϵ^2 .
- Should an additional variable be included as independent variable in the model?
- Including an additional variable will *always*
 - increase R^2 ,
 - reduce SSE,
 - decrease the degrees of freedom.
- This is why including an additional variable need not reduce s_ϵ^2 — care needs to be taken!



16.6 Examples

Example: Returns on OSG stock.

Overseas Shipholding Group, Inc. (“OSG”), is a marine transportation company whose stock is listed at New York Stock Exchange (NYSE).

Let variables be defined as

osg.ret = monthly return on OSG stock;

nyse.ret = monthly return on the NYSE Composite Index;

sop.ret = monthly change in spot oil price (WTI);

export = exported goods (from USA), in million USD

Question: Which variables can explain returns on OSG stock?



16.6 Examples

Example: Returns on OSG stock.

Model 1:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.4989	1.1801	1.270	0.209
nyse.ret	1.4737	0.3067	4.805	1.2e-05 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.962 on 56 degrees of freedom

Multiple R-Squared: 0.2919, Adjusted R-squared: 0.2793

F-statistic: 23.09 on 1 and 56 DF, p-value: 1.200e-05



16.6 Examples

Example: Returns on OSG stock.

Model 2:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	3.592e+00	1.167e+01	0.308	0.759
nyse.ret	1.478e+00	3.101e-01	4.764	1.43e-05 ***
export	-3.319e-05	1.841e-04	-0.180	0.858

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.041 on 55 degrees of freedom

Multiple R-Squared: 0.2923, Adjusted R-squared: 0.2666

F-statistic: 11.36 on 2 and 55 DF, p-value: 7.419e-05



16.6 Examples

Example: Returns on OSG stock.

Model 3:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.9753	1.1812	0.826	0.4125
nyse.ret	1.5615	0.3024	5.163	3.45e-06 ***
sop.ret	0.3025	0.1536	1.970	0.0539 .

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8.74 on 55 degrees of freedom

Multiple R-Squared: 0.3386, Adjusted R-squared: 0.3145

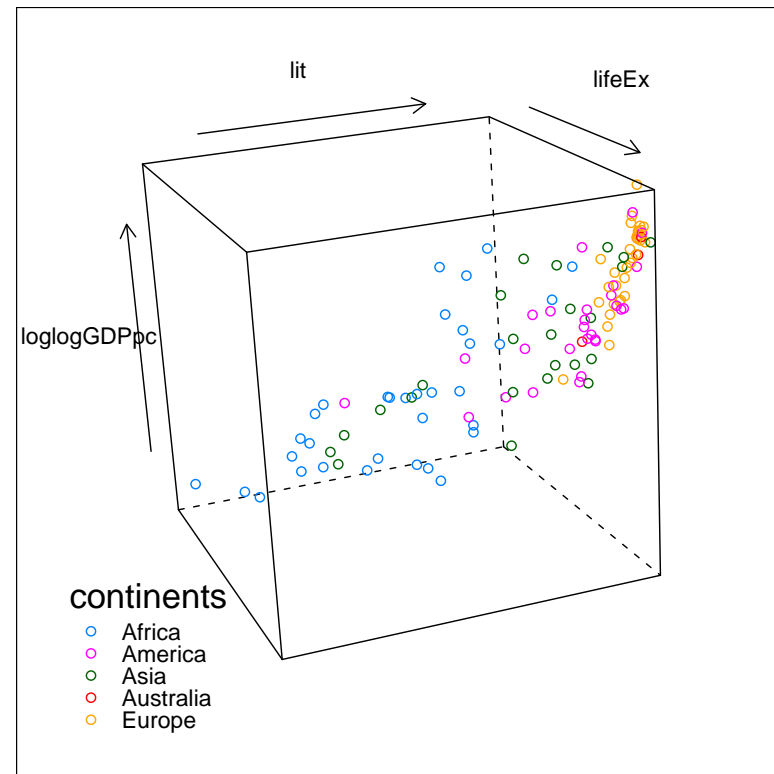
F-statistic: 14.08 on 2 and 55 DF, p-value: 1.156e-05



16.6 Examples

Example: Life expectancy, literacy, GDP.

What is the relation between literacy, the expectation of life, and (doubly logged) GDP per capita?



16.6 Examples

Example: Life expectancy, literacy, GDP.

Model 1:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-103.386	9.158	-11.29	<2e-16	***
log(log(GDPpc))	78.875	4.253	18.55	<2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.538 on 119 degrees of freedom

Multiple R-Squared: 0.743, Adjusted R-squared: 0.7408

F-statistic: 344 on 1 and 119 DF, p-value: < 2.2e-16



16.6 Examples

Example: Life expectancy, literacy, GDP.

Model 2:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	27.66047	3.55972	7.77	3.08e-12	***
lit	0.46619	0.04199	11.10	< 2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 9.038 on 119 degrees of freedom

Multiple R-Squared: 0.5088, Adjusted R-squared: 0.5046

F-statistic: 123.2 on 1 and 119 DF, p-value: < 2.2e-16



16.6 Examples

Example: Life expectancy, literacy, GDP.

Model 3:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-90.64350	11.36348	-7.977	1.09e-12	***
log(log(GDPpc))	69.62269	6.51710	10.683	< 2e-16	***
lit	0.08656	0.04655	1.860	0.0654	.

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.471 on 118 degrees of freedom

Multiple R-Squared: 0.7503, Adjusted R-squared: 0.7461

F-statistic: 177.3 on 2 and 118 DF, p-value: < 2.2e-16



16.7 Prediction Based on MLR

Point prediction vs. interval prediction. (Case $k = 2$.)

Let x_1, x_2 be given. The outcome of the random variable $Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \epsilon$ can be predicted in terms of. . .

- a single point: $\hat{Y} = \hat{\alpha} + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$
 - This has disadvantages similar to those of a point estimate.
- a prediction interval.
It has to cope with two sources of uncertainty:
 - The parameters α, β_1, β_2 are unknown.
 - There is a random error ϵ , which has an unknown variance σ_ϵ^2 .



16.7 Prediction Based on MLR

Prediction intervals. (Case $k = 2$.)

Given a vector $x_0 = (1, x_{1,n+1}, x_{2,n+1})'$ with out-of-sample values $x_{1,n+1}$ and $x_{2,n+1}$, a 95% prediction interval for the corresponding Y_{n+1} has bounds

$$\hat{Y}_{n+1} \pm t_{n-k-1, 0.975} \cdot s_\epsilon \cdot \sqrt{1 + x_0'(X'X)^{-1}x_0}$$

These are the bounds of an interval which will contain the random variable $Y_{n+1} = \alpha + \beta_1 x_{1,n+1} + \beta_2 x_{2,n+1} + \epsilon$ with probability 95%.

Here, \hat{Y}_{n+1} is a point prediction, obtained as

$$\hat{Y}_{n+1} = \hat{\alpha} + \hat{\beta}_1 x_{1,n+1} + \hat{\beta}_2 x_{2,n+1}.$$



16.7 Prediction Based on MLR

Prediction intervals. (Case $k = 2$.)

An approximation formula for the interval bounds is

$$\hat{Y}_{n+1} \pm t_{n-k-1, 0.975} \cdot s_{\epsilon} \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_{1,n+1} - \bar{x}_1)^2}{\sum (x_{1i} - \bar{x}_1)^2} + \frac{(x_{2,n+1} - \bar{x}_2)^2}{\sum (x_{2i} - \bar{x}_2)^2}}$$

- This formula may be used if the independent variables are uncorrelated and n is large.
- The generalization to $k > 2$ is straightforward.



16.7 Prediction Based on MLR

Example: Used cars.

- Based on a sample of size $n = 400$, the fitted model is:

$$\text{price} = 14146.2 - 24.61 \cdot \text{mileage} - 49.13 \cdot \text{age}$$

- Point forecast of the price of a car with mileage 100000 km, age 10 years:

$$14146.2 - 24.61 \cdot 100 - 49.13 \cdot 10 = 5789.6$$



16.7 Prediction Based on MLR

Example: Used cars.

- Bounds of a 95% prediction interval:

exact formula: $5789.6 \pm 1.966 \cdot 1240 \cdot 1.002807$

approximate formula: $5789.6 \pm 1.966 \cdot 1240 \cdot 1.003476$

- Corresponding 95% prediction intervals:

exact formula: $[3345.0, 8234.3]$

approximate formula: $[3343.4, 8235.9]$

