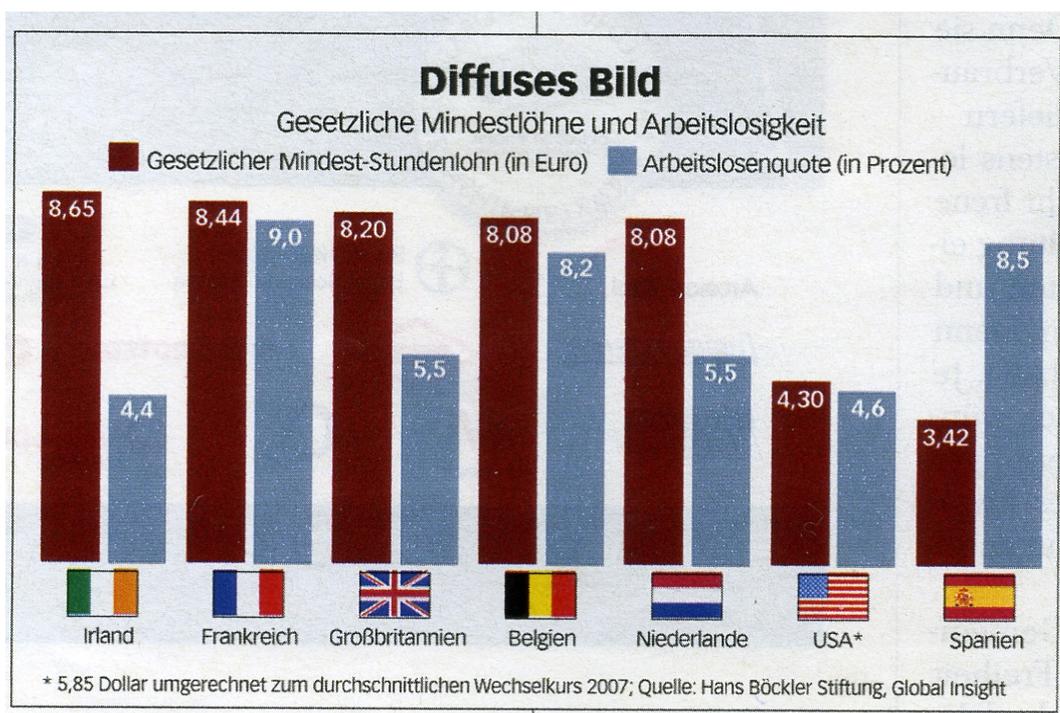


## Zwei Rechenbeispiele für die einfache lineare Regression

### 1. Mindestlöhne – Beispiel<sup>1</sup>

#### 1.1. Daten

Entnommen aus Rolf Ackermann, Spielball des Lobbyisten, Mindestlöhne schaden nicht nur bei Postdiensten sondern in allen Branchen, in: Wirtschaftswoche Nr. 50 (10.12.2007)



Es soll gelten  $x_i$  = Höhe des Mindestlohns ( $x$  ist später bei einer Erweiterung der Aufgabe  $x_1$ ),  $y_i$  = Arbeitslosenquote. Wir haben hier Querschnitts-, nicht Zeitreihendaten, daher der Laufindex  $i = 1, 2, \dots, n$  statt  $t = 1, 2, \dots, T$

#### 1.2. Normalgleichungen

Die Zahlen sind leicht der auf der nächsten Seite wiedergegebenen Excel-Tabelle zu entnehmen

Allgemein	Mit den Daten	In Matrixschreibweise
$\alpha n + \beta \sum x = \sum y$ $\alpha \sum x + \beta \sum x^2 = \sum xy$	$7 \cdot \alpha + 49,17 \cdot \beta = 45,7$ $49,17 \cdot \alpha + 374,06 \cdot \beta = 318,67$	$\begin{bmatrix} n & \sum x \\ \sum x & \sum x^2 \end{bmatrix} \cdot \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} \sum y \\ \sum xy \end{bmatrix}$

#### Exkurs

(kein Muss für Hörer, die diese Darstellungsart nicht mögen)

Die relevanten Matrizen und Vektoren in diesem Beispiel sind die Momentenmatrix

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} n & \sum x \\ \sum x & \sum x^2 \end{bmatrix} \text{ und die Datenmatrix } \mathbf{X} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \text{ sowie der Datenvektor } \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}.$$

<sup>1</sup> Anders als die zweite Aufgabe (Affenaufgabe) wird diese Aufgabe wieder aufgegriffen bei der multiplen Regression.

Das Modell lautet somit in Matrixschreibweise  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  und die Normalgleichungen als Ergebnis der Methode der kleinsten Quadrate sind  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$  mit  $\hat{\boldsymbol{\beta}}' = [\hat{\alpha} \quad \hat{\beta}]$  und  $\mathbf{u}$  analog zu  $\mathbf{y}$ .

### 1.3. Excel-Tabelle, Berechnung der Regressionskoeffizienten und ihrer Varianzen

A	B	C	D	E	F
	x	y	xy	x <sup>2</sup>	y <sup>2</sup>
Irland	8,65	4,4	38,06	74,8225	19,36
Frankreich	8,44	9	75,96	71,2336	81
Großbritannien	8,2	5,5	45,1	67,24	30,25
Belgien	8,08	8,2	66,256	65,2864	67,24
Niederlande	8,08	5,5	44,44	65,2864	30,25
USA	4,3	4,6	19,78	18,49	21,16
Spanien	3,42	8,5	29,07	11,6964	72,25
Berechnungen					
Spaltensumme S	49,17	45,7	318,666	374,0553	321,5100
Mittelwert S/n	7,0243	6,5286	45,524	53,4365	45,9300

<b>Varianz von x</b>	4,7785*	⇐ Excel	4,0959**	Steigung	-0,08175
<b>Varianz von y</b>	3,8590*	⇐ Excel	3,3078**	Ordinate	7,102804
<b>Kovarianz</b>	-0,3348	⇐ Excel		Korrelation	-0,09097

<b>Die Regressionsfunktion lautet 7,1028 – 0,08175*x</b>	
<b>Die Variablen sind praktisch nicht miteinander korreliert</b>	
Die Bestimmtheit ist nur 0,00828	

\* Diese Werte errechnet Excel als Varianzen (durch n-1 statt durch n geteilt)

\*\* Excel-Werte mit  $(n-1)/n = 6/7$  multipliziert

Die eingegebenen Daten sind in helltürkis markiert. Man kann bestimmte Werte durch Eingeben einer Berechnungsformel berechnen, etwa die Spaltensummen oder die Mittelwerte, um hiermit weiter zu rechnen, etwa um  $s_x^2 = \frac{1}{n} \sum_i x_i^2 - \bar{x}^2 = 53,4365 - (7,0243)^2 = 4,0957$  (Rundungsfehler,<sup>2</sup> vgl. oben 4,0959) zu bestimmen, oder man lässt dies mit der Excel Funktion (mit f<sub>x</sub> wählen!) Mittelwert bzw. Varianz (oder auch Kovarianz) "automatisch" berechnen. Bei den Varianzen wird von Excel jedoch durch n - 1 = 6 geteilt (gelb markierte Felder). Um auf die bekannten Formeln

$s_x^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$  und  $s_y^2$  umzurechnen ist mit 6/7 zu multiplizieren. Hieraus lassen sich Größen

$\hat{\alpha}$  (Ordinatenabschnitt),  $\hat{\beta}$  (Steigung), r und r<sup>2</sup> bestimmen. Ferner gilt (Cramerschen Regel)

$$\hat{\alpha} = \frac{\begin{vmatrix} \Sigma y & \Sigma x \\ \Sigma xy & \Sigma x^2 \end{vmatrix}}{\begin{vmatrix} n & \Sigma x \\ \Sigma x & \Sigma x^2 \end{vmatrix}} = \frac{\begin{vmatrix} 45,7 & 49,17 \\ 318,666 & 374,0553 \end{vmatrix}}{\begin{vmatrix} 7 & 49,17 \\ 49,17 & 374,0553 \end{vmatrix}} = \frac{1425,52}{200,6982} = 7,1028$$

und entsprechend für die Steigung

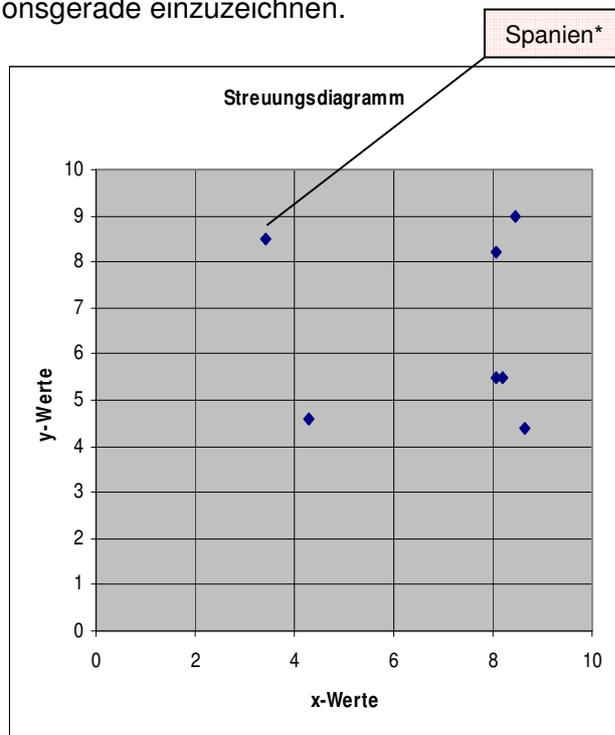
<sup>2</sup> Die Berechnung ist offenbar sehr fehleranfällig (Rundungsfehler!) zumal n sehr klein ist und sie erfolgt am besten ausgehend von den Normalgleichungen mit der Cramerschen Regel. In dieser Hinsicht ist das zweite Rechenbeispiel (Aufgabenstellung) sehr viel angenehmer, weil hier zwar mit noch weniger, dafür aber "glatteren" Zahlen gearbeitet wird.

$$\hat{\beta} = \frac{\begin{vmatrix} n & \Sigma y \\ \Sigma x & \Sigma xx \end{vmatrix}}{\begin{vmatrix} n & \Sigma x \\ \Sigma x & \Sigma x^2 \end{vmatrix}} = \frac{\begin{vmatrix} 7 & 45,7 \\ 49,17 & 318,666 \end{vmatrix}}{\begin{vmatrix} 7 & 49,17 \\ 49,17 & 374,0553 \end{vmatrix}} = \frac{-16407}{200,6982} = -0,08175.$$

Die Parameter  $\hat{\alpha}$ ,  $\hat{\beta}$ ,  $r$  und  $r^2$  kann man auch direkt mit den Excelfunktionen  $f_x$  bestimmen. Man erhält dann:

Parameter	Excel Funktion $f_x$	Ergebnis
Ordinatenabschnitt $\hat{\alpha}$	"Achsenabschnitt"	7,102804
Steigung $\hat{\beta}$	"Steigung"	-0,08174961
Korrelation $r$	"Pearson" oder "Korrel"	-0,09096885
Bestimmtheit $r^2$	"Bestimmtheitsmaß"	0,008275332

Das Streuungsdiagramm erhält man als Graphik vom Typ "Punkt (XY)"<sup>3</sup>. Anders als beim nächsten Beispiel (Affenaufgabe) ist hier darauf verzichtet worden, die von Excel bestimmte Regressionsgerade einzuzeichnen.



Man kann nun auch die Größen bestimmen, die wichtig sind für das Schätzen und Testen von Regressionskoeffizienten. Man erhält (in der Symbolik des Buches von v. Auer) die folgenden Werte:

$$S_{yy} = 45,93 - (6,5286)^2 = 3,3077551$$

$$S_{\hat{y}\hat{y}} = r^2 S_{yy} = 0,027372771$$

$$S_{\hat{u}\hat{u}} = S_{yy} - S_{\hat{y}\hat{y}} = 3,28038233$$

Die geschätzte Varianz der Störgröße ist

$$\text{danach } \hat{\sigma}^2 = \frac{S_{\hat{u}\hat{u}}}{n-2} = \frac{3,2804}{5} = 0,656076$$

$$\text{Ferner ist } S_{xx} = 53,4365 - (7,0243)^2 = 4,09571$$

$$\text{und mit } \hat{\sigma}_{\hat{\beta}}^2 = \frac{\hat{\sigma}^2}{S_{xx}} = \frac{0,656}{4,096} = 0,16018 \text{ erhält}$$

man die geschätzte Varianz von  $\hat{\beta}$  und somit für die Standardabweichung von  $\hat{\beta}$  den Wert 0,400223.

\* Spanien scheint ein Ausreißer zu sein. Rechnet man ohne Spanien, so ist der Korrelationskoeffizient 0,37482 statt -0,09097 (allerdings ist n dann auch nur noch 6).

Für die geschätzte Varianz und die (geschätzte) Standardabweichung (standard deviation S.D. oder "Std. Error") von  $\hat{\alpha}$  ergibt sich daraus  $\hat{\sigma}_{\hat{\alpha}}^2 = \overline{x^2} \hat{\sigma}_{\hat{\beta}}^2$  mit dem quadrierten quadratischen Mittel  $\overline{x^2} = \sum x_i^2 / n = 53,4365$ , so dass die Varianz  $\hat{\sigma}_{\hat{\alpha}}^2$  den Wert  $53,4365 \cdot 0,16018 = 8,5595$  annimmt ( $\hat{\sigma}_{\hat{\alpha}} = 2,9256$ ). Die t-Werte sind demnach  $t = 7,1028 / 2,9256 = 2,4283$  bei der Hypothesen  $H_0: \alpha = 0$  und  $t = -0,08175 / 0,400223 = -0,20426$  bei der  $H_0: \beta = 0$ . Somit ist zwar  $\alpha$ , nicht aber  $\beta$  signifikant von 0 verschieden. Berechnungen dieser Art (das Schätzen und Testen betreffend) und vor allem eine multiple Regression lassen sich besser mit **EViews**, statt mit Excel durchführen. Die Erweiterung des Beispiels Mindestlöhne zu einer Aufgabe der multiplen Regression mit den ent-

<sup>3</sup> Wenn man den Bereich markiert, auf den sich die Grafik beziehen soll, dann sollte man auch die Felder x und y mit markieren.

sprechenden Berechnungen findet sich in einem weiteren Download. Das dort mit EViews ermittelte Ergebnis (y in Abhängigkeit von  $x = x_1$ ) sei hier jedoch bereits (verkürzt) wiedergegeben (Ergebnisse, die mit den oben [z.T. mit Excel] berechneten Ergebnissen verglichen werden können sind gelb unterlegt):

Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	7.102804	2.925650	2.427769	0.0595
X1	-0.081750	0.400224	-0.204259	0.8462
R-squared	0.008275	Mean dependent var		6.528571
Adjusted R-squared	-0.190070	S.D. dependent var		1.964446*
S.E. of regression	2.143020	Akaike info criterion		4.597266
Sum squared resid	22.96268	Schwarz criterion		4.581811

\* das ist die Wurzel aus 3,8590 in der Excel-Tabelle auf Seite 2 oben.

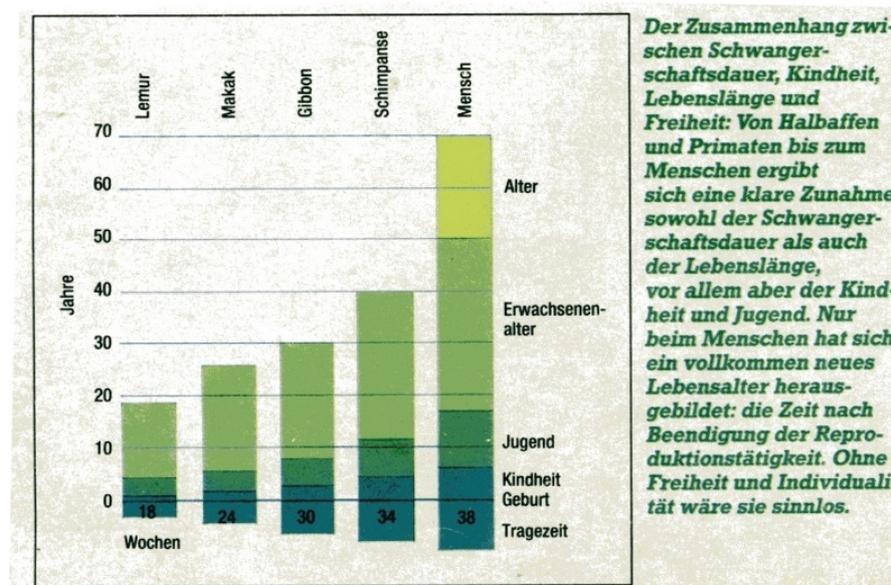
## 2. "Affenaufgabe"<sup>4</sup>

### 2.1. Daten

Die folgenden Daten über den Zusammenhang zwischen Dauer der Schwangerschaft und Lebenserwartung (der Mensch als "Ausreißer") sind entnommen aus der Zeitschrift Focus

	X = Dauer der Schwangerschaft	Y = Lebenserwartung
Lemur	18	18
Makak	24	26
Gibbon	30	30
Schimpanse	34	40
Mensch	38	70
Summe	144	184

Die hier wiedergegebene Abbildung aus FOCUS zeigt, wie schwierig es ist, die Daten verständlich graphisch darzustellen, wenn man glaubt, bei den statistisch nicht vorgebildeten Lesern nicht Gebrauch machen zu können von der Möglichkeit eines Streudiagramms. Man muss dann wohl mit den verschiedensten Farben operieren und es ist sehr fraglich, ob die Dinge so klarer und leichter verständlich werden als mit einem Streudiagramm.

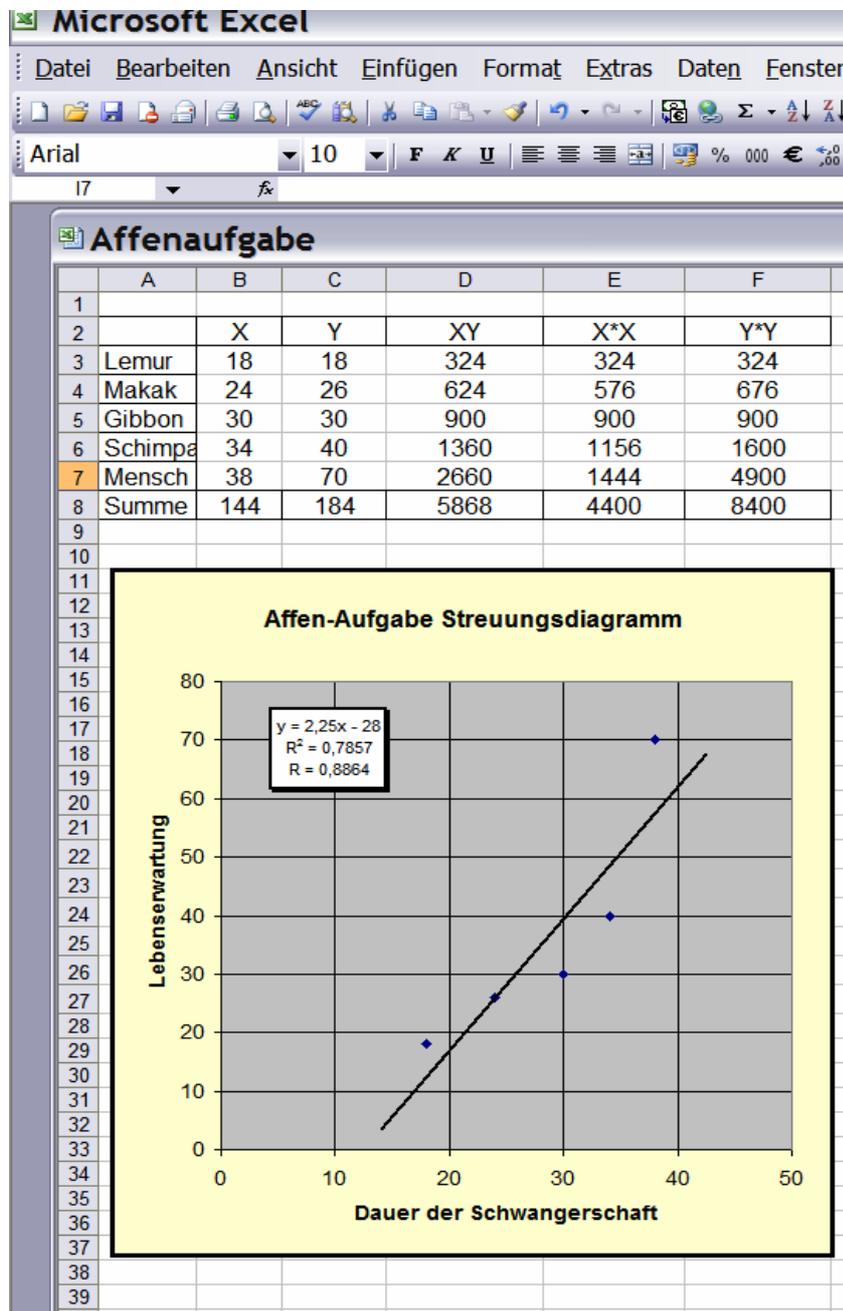


<sup>4</sup> Vorteil dieses Beispiels:

sehr wenige und zudem glatte Zahlen als Daten, so dass es leicht möglich ist, alles mit dem Taschenrechner nachzurechnen.

## 2.2. Berechnungen zur Deskriptiven Statistik, Excel Tabelle/Grafik und Normalgleichungen

Man sieht hier den Bildschirm und die Eingabe der Daten, wobei in den Spalten D, E und F einige einfacher Berechnungen durchgeführt werden, die zur Bestimmung der Normalgleichungen notwendig sind:



Man kann mit diesen Angaben leicht die Normalgleichungen zusammenstellen und erhält so

$$5 \cdot \alpha + 144 \cdot \beta = 184$$

$$144 \cdot \alpha + 4400 \cdot \beta = 5868$$

Die Schätzwerte  $\hat{\alpha}$  und  $\hat{\beta}$  erhält man nach der Cramerschen Regel mit drei Determinanten wie folgt

$$\hat{\alpha} = \frac{\begin{vmatrix} 184 & 144 \\ 5868 & 4400 \end{vmatrix}}{\begin{vmatrix} 5 & 144 \\ 144 & 4400 \end{vmatrix}} = -28$$

$$\hat{\beta} = \frac{\begin{vmatrix} 5 & 184 \\ 144 & 5868 \end{vmatrix}}{\begin{vmatrix} 5 & 144 \\ 144 & 4400 \end{vmatrix}} = 2,25$$

Die Regressionsgerade lautet mithin

$$- 28 + 2,25 \cdot x.$$

Für den Korrelationskoeffizienten erhält man mit Excel den Wert + **0,7857**, so dass die Bestimmtheit  $r^2 = 0,8864$ , also 88,64%.

Die Regressionsgerade erhält man als "Trend" wenn man in der Graphik einen y Punkt anklickt und die rechte Maustaste drückt. Dann kommt ein Menü, mit dem man die Trendlinie (mit denen Typen linear, logarithmisch, gleitende Mittelwerte etc.) wählen kann. Wenn keine Verlängerung "vorwärts" oder "rückwärts" gewählt wird, zeichnet Excel nur die Gerade im Bereich zwischen  $x_{\min}$  und  $x_{\max}$ . Den eingezeichneten Trend anklicken → jetzt kann man die Trendlinie formatieren (dicker, farbig etc). Man kann auch – wie hier geschehen – die Option "Parameter anzeigen" wählen und das dazu gehörige Textfeld mit der Regressionsgleichung, der Bestimmtheit  $R^2$  (entspricht  $r^2$ ) bearbeiten und auch nachträglich den Korrelationskoeffizienten (Funktion "Pearson") eintragen.

Ohne Excel kann man die Parameter  $\hat{\alpha}$  und  $\hat{\beta}$  sowie  $r$  und  $\hat{\sigma}^2$  auch wie folgt berechnen:

Mittelwerte	$\bar{x} = 144/5 = 28,8$	$\bar{y} = 184/5 = 36,8$
Varianzen (und Summen der Abweichungsquadrate)	$s_x^2 = 4400/5 - (28,8)^2 = 50,56$ $S_{xx} = T s_x^2 = 5 \cdot 50,56 = 252,8$	$s_y^2 = 8400/5 - (36,8)^2 = 325,76$ $S_{yy} = 1628,8$
Kovarianz	$s_{xy} = S_{xy}/T - \bar{x} \cdot \bar{y} = 5868/5 - 28,8 \cdot 36,8 = 113,76$	
Steigung	$\hat{\beta} = s_{xy}/s_x^2 = 113,76/50,56 = 2,25$	
Ordinatenabschnitt	$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} = 36,8 - 2,25 \cdot 28,8 = -28$	
Sum squared resid.	$S_{\hat{u}\hat{u}} = 349$ (verschiedene Berechnungsmöglichkeiten siehe unten)	
S.E. of regression $\hat{\sigma}$	$\hat{\sigma}^2 = \hat{\sigma}_u^2 = S_{\hat{u}\hat{u}}/(T-2) = 349/3 = 116,33 \rightarrow \hat{\sigma} = \sqrt{116,33} = 10,786$	
Korrelation	$r = s_{xy}/\sqrt{s_x^2 s_y^2} = 0,8864 \rightarrow \text{Bestimmtheit } r^2 = 0,78573$	
Zweites Anfangsmoment	$\overline{x^2} = \sum x_t^2/T = 4400/5 = 880$	

Es gilt  $S_{xx} = T \cdot s_x^2 = \frac{1}{T} \cdot \left| \begin{matrix} T & \sum x_t \\ \sum x_t & \sum x_t^2 \end{matrix} \right|$ , ferner da  $\hat{y}_t - \bar{y} = \hat{\beta}(x_t - \bar{x})$  ist auch  $S_{\hat{y}\hat{y}} = \hat{\beta}^2 S_{xx} = \hat{\beta}^2 T s_x^2$

und  $S_{\hat{u}\hat{u}} = S_{yy} - S_{\hat{y}\hat{y}} = 1628,8 - 1279,8 = 349$  (oder  $S_{\hat{u}\hat{u}} = (1-r^2) \cdot T \cdot s_y^2$ ).

### 2.3. Berechnungen zur Induktiven Statistik

#### a) Konfidenzintervall für die Streuung

Mit dem  $\alpha/2$  sowie dem  $1-\alpha/2$  Quantil (d.h. mit der unteren und oberen Signifikanzschranke) der  $\chi^2$  Verteilung bei  $\alpha = 0,05$  und  $5-2 = 3$  Freiheitsgraden  $z_u = 0,216$  und  $z_o = 9,348$  erhält man die folgenden Grenzen des 95% Konfidenzintervalls für die unbekannte Varianz  $\sigma^2$  in der Grundgesamtheit: Untergrenze  $\frac{S_{\hat{u}\hat{u}}}{z_o} = \frac{349}{9,348} = 37,33$  und Obergrenze  $\frac{S_{\hat{u}\hat{u}}}{z_u} = \frac{349}{0,216} = 1615,74$ .

Es wird gerne vergessen, dass<sup>5</sup>  $\hat{\sigma}_u^2 = \hat{\sigma}^2 = 349/3 = 116,333$  genauso ein zu schätzender Parameter des Modells der einfachen Regression ist wie die Parameter  $\hat{\alpha}$  und  $\hat{\beta}$ .

#### b) Regressionskoeffizienten, $H_0: \rho = 0$ (Varianzanalyse)

Die weiteren Größen sind für die Bestimmung von Konfidenzintervallen und zur Durchführung von Tests wichtig

$\hat{\sigma}_\beta^2 = \frac{\hat{\sigma}^2}{S_{xx}} = \frac{116,33}{252,8} = 0,4600475$	$\hat{\sigma}_\alpha^2 = \overline{x^2} \hat{\sigma}_\beta^2 = 880 \cdot 0,46 = 404,842$	$\hat{\sigma}_{\alpha\beta} = -\bar{x} \cdot \hat{\sigma}_\beta^2 = -13,249$
---	--	--

Die Matrix der (geschätzten) Varianzen und Kovarianzen der Regressionskoeffizienten ist demnach

$$\mathbf{V} = \hat{\sigma}^2 \mathbf{X}'\mathbf{X} = \begin{bmatrix} 404,842 & -13,249 \\ -13,249 & 0,46005 \end{bmatrix}$$

Alle weiteren Berechnungen (Konfidenzintervalle für  $\alpha$  und  $\beta$  Prognoseintervall für  $y_0$  und sowie Varianzanalyse) vgl. Vorlesung.

<sup>5</sup> Man beachte, dass 116,333 nicht in der Mitte liegt zwischen 37,33 und 1615,74. Das liegt daran, dass die  $\chi^2$  Verteilung nicht symmetrisch ist.