

### 3. Beschreibung von mehrdimensionalen Daten

#### 3.1 Grafische Darstellungen einer Stichprobe von Zahlenpaaren

Werden bei einer Stichprobe zwei verschiedene Merkmale  $x_i$  und  $y_i$  beobachtet, so spricht man von einer bivariaten Stichprobe. Dabei stellt sich die Frage, ob es einen Zusammenhang, eine Abhängigkeit zwischen den beiden Grössen gibt. In diesem Fall sagt man, dass die beiden Variablen korreliert sind.

Beispiele:

Druck und Volumen bei einem physikalischen Experiment, Spraydosen-Treibgasverbrauch und mittlerer Ozongehalt der Atmosphäre in einem Jahr, Grösse des Vaters und des Sohnes bei einem Vater-Sohn-Paar.

Die Daten  $(x_i, y_i)$  können als Punkte in der Koordinatenebene dargestellt werden. Diese Darstellung heisst Streudiagramm (englisch: Scatterplot). Sie erleichtert es, bei Zahlenpaaren Trends, Muster, Ausreisser zu erkennen.

Beispiel:

Das Rauschen von Wasserfällen (Quelle: Science 164,1969, p. 1513-1514)

Bei verschiedenen Wasserfällen wurde die Höhe und die dominierende Frequenz im Spektrum der Bodenvibrationen gemessen. Im Streudiagramm ist ersichtlich, dass hohe Frequenzen mit niedrigen Höhen einhergehen und umgekehrt.

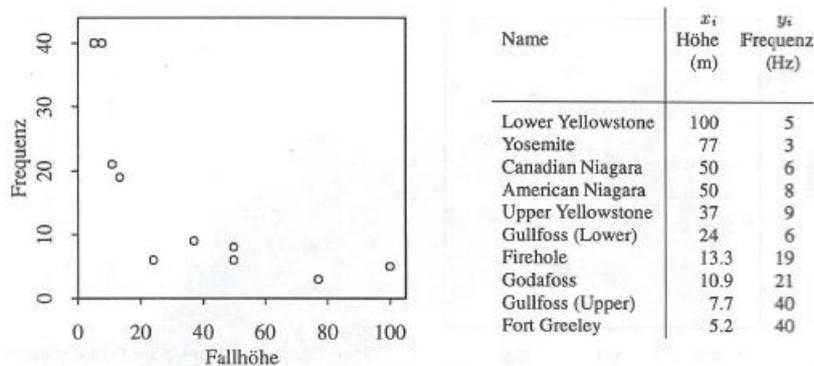


Bild 3.1.d Streudiagramm für Fallhöhe und Schwingungsfrequenz bei Wasserfällen

Sind mehr als zwei Variablen von Interesse, dann wird - wie im folgenden Beispiel - in der so genannten Scatterplot-Matrix jeweils für zwei Merkmale ein Streudiagramm abgebildet.

Beispiel:

Es sind die Streudiagramme der mittleren Monatsrenditen der vier Aktien: Volkswagen; BASF, Siemens, Münchner Rück und zusätzlich der Libor-Zins abgebildet. Der Libor-Zins (London Interbank Offered rate) ist der Zinssatz, zu dem sich Geschäftsbanken gegenseitig Geld mit einer Laufzeit bis zu einem Jahr leihen. In der Matrix ist zu erkennen, dass die Höhe der Renditen für Aktien untereinander deutlich zusammenhängen, während kein deutlicher Zusammenhang zwischen Aktienrendite und Zins erkennbar ist.

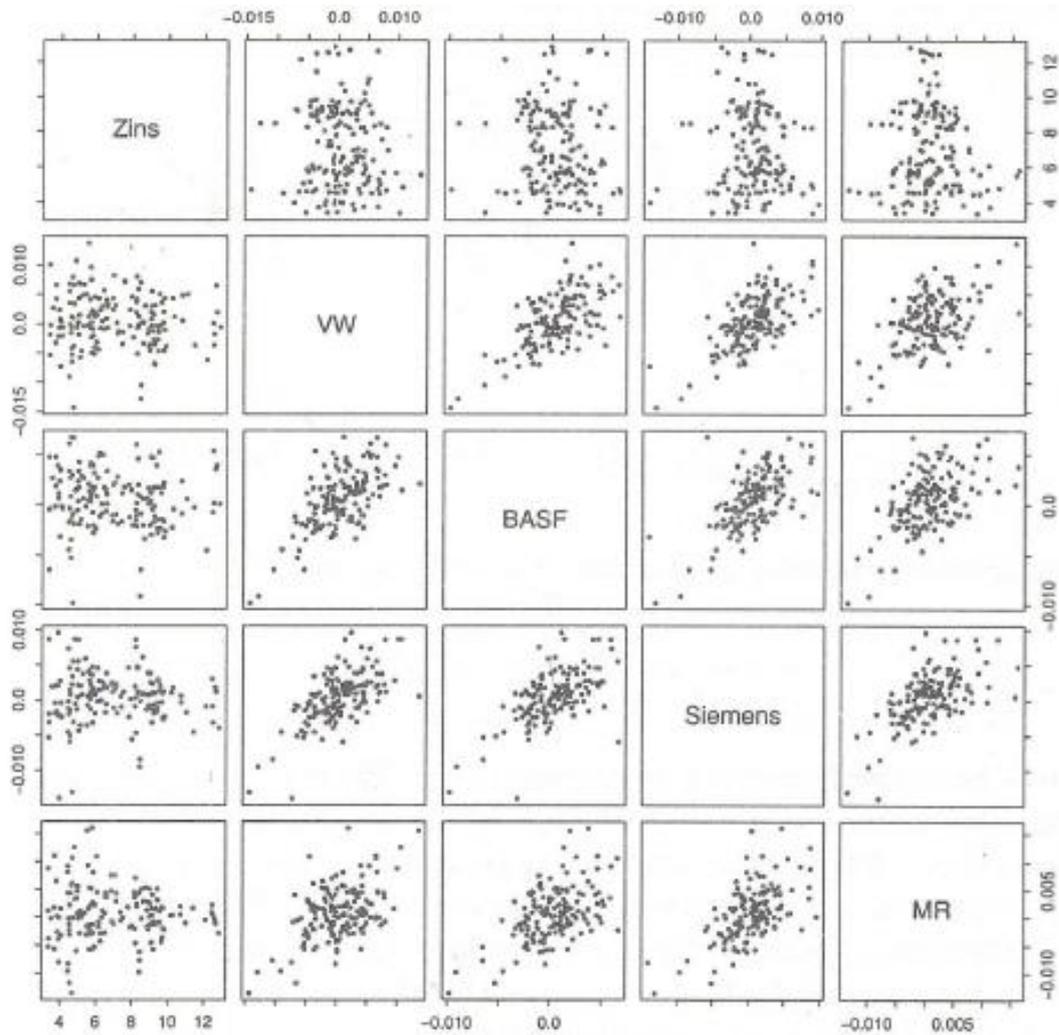


ABBILDUNG 3.10: Streudiagramm-Matrix der mittleren Monatsrenditen für Aktien und des Zinses in Form des Libors

Quelle: Fahrmeier: Statistik, Springer, 2001

### 3.2 Kennzahlen für bivariate Stichproben

Zunächst können die bekannten Kennzahlen für univariate Stichproben bestimmt werden:

Empirischer Mittelwert

$$\boxed{\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i} \quad \boxed{\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i} \quad 1)$$

Varianz  $s_x^2$  bzw. Standardabweichung  $s_x$

$$\boxed{\begin{aligned} s_x^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 \right) \\ &= \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) \end{aligned}} \quad 2)$$

bzw.

$$\boxed{\begin{aligned} s_y^2 &= \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n y_i^2 - \frac{1}{n} \left( \sum_{i=1}^n y_i \right)^2 \right) \\ &= \frac{1}{n-1} \left( \sum_{i=1}^n y_i^2 - n\bar{y}^2 \right) \end{aligned}} \quad 2)$$

Als Gesamtmaß wird die folgende sogenannte Produkt-Momentenkorrelation  $r_{xy}$  (englisch: Pearson-Correlation) verwendet:

$$\boxed{r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}} \quad 5)$$

Definiert man als Empirische Kovarianz  $s_{xy}$

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n-1} \left( \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \right) \quad 6)$$

$$= \frac{1}{n-1} \left( \sum_{i=1}^n x_i y_i - \frac{1}{n} \left( \sum_{i=1}^n x_i \right) \left( \sum_{i=1}^n y_i \right) \right)$$

dann kann die Produkt-Momentenkorrelation  $r_{xy}$  mit 6) kurz folgendermassen geschrieben werden:

$$r_{xy} = \frac{s_{xy}}{s_x s_y} \quad \dots\dots\dots 5)$$

Schreibt man die Definition 5) in der folgenden Form um, so ist zu erkennen, dass die Korrelation nicht vom Nullpunkt der Skalen von x bzw. y und von den entsprechenden Masseneinheiten abhängt:

$$r_{xy} = \sum_{i=1}^n \frac{(x_i - \bar{x})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \cdot \frac{(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Führt man analog zum Vorgehen bei der Normalverteilung die folgenden neuen standardisierten Variablen ein

$$u_i = \frac{x_i - \bar{x}}{s_x}$$

bzw.

$$v_i = \frac{y_i - \bar{y}}{s_y}$$

und bildet man mit diesen Werten die Vektoren  $\vec{u}$  und  $\vec{v}$ , so vereinfacht sich wegen

$|\vec{u}| = |\vec{v}| = \sqrt{n-1}$  die Formel und die Summe kann als Skalarprodukt aufgefasst werden.

$$r_{xy} = \frac{1}{n-1} \cdot \sum_{i=1}^n u_i v_i = \frac{1}{n-1} \cdot (\vec{u} \cdot \vec{v}) \quad 7)$$

Da die Definition der Korrelation nicht von der Wahl der linearen Skala abhängt, gilt zudem:

$$r_{xy} = r_{uv}$$

Wenn beide Koordinaten das gleiche Vorzeichen haben (die zugehörigen Punkte liegen dann im verschobenen Koordinatensystem mit Zentrum im Schwerpunkt S  $(\bar{x}, \bar{y})$  im ersten und dritten Quadranten), so deutet dies auf einen „positiven“ Zusammenhang zwischen den beiden Variablen, und entsprechend bedeuten Koordinaten mit ungleichem Vorzeichen auf einen „negativen“ Zusammenhang hin (die zugehörigen Punkte liegen im zweiten und vierten Quadranten).

Die aufgeführten Umformungen bei 1) bis 7) ermöglichen die Berechnung aller Kennzahlen aus geeigneten Summen. Allerdings besteht bei dieser Berechnungsart die Gefahr von Rundungsfehlern.

Beispiel Wasserfälle:

Höhe [m]	Frequenz [Hz]							
$x_i$	$y_i$	$x_i^2$	$y_i^2$	$x_i y_i$	$u_i$	$v_i$	$u_i v_i$	
100	5	10000.00	25	500.0	1.9559956	-0.75911772	-1.48483092	
77	3	5929.00	9	231.0	1.23607403	-0.90100889	-1.11371369	
50	6	2500.00	36	300.0	0.39094871	-0.68817214	-0.26904001	
50	8	2500.00	64	400.0	0.39094871	-0.54628098	-0.21356784	
37	9	1369.00	81	333.0	-0.01596348	-0.4753354	0.00758801	
24	6	576.00	36	144.0	-0.42287567	-0.68817214	0.29101125	
13.3	19	176.89	361	252.7	-0.7577957	0.23412042	-0.17741545	
10.9	21	118.81	441	228.9	-0.83291795	0.37601158	-0.3131868	
7.7	40	59.29	1600	308.0	-0.93308095	1.72397763	-1.60861069	
5.2	40	27.04	1600	208.0	-1.0113333	1.72397763	-1.74351598	
375.1	157	23256.03	4253	2905.6	Summe	0	0	-6.62528213
37.51					$\bar{x}$			
	15.7				$\bar{y}$			
		1020.67			$s_x$			
		31.95			$s_y$			
			198.68		$s_{xy}$			
			14.10		$r_{xy}$			
				-331.50				
				-0.736				

Ein weiteres Beispiel:

In der Klasse 1C der Kantonsschule Zofingen wurden am 21.3.2003 die Körpergrößen der Paare Tochter- Mutter bzw. Sohn-Vater gemessen.

### Körpergröße 1C

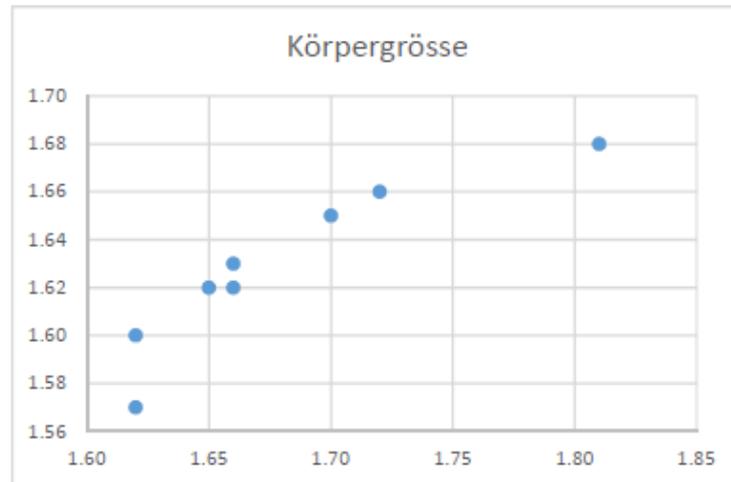
21.03.2003

Tochter x	Mutter y
1.62	1.57
1.62	1.60
1.65	1.62
1.66	1.62
1.66	1.63
1.70	1.65
1.72	1.66
1.81	1.68

Korrelation 0.911

#### Teilresultate

$\bar{x}$	1.680
$sd_x$	0.063
$\bar{y}$	1.629
$sd_y$	0.035
Kovarianz	0.002
Korrelation	0.911

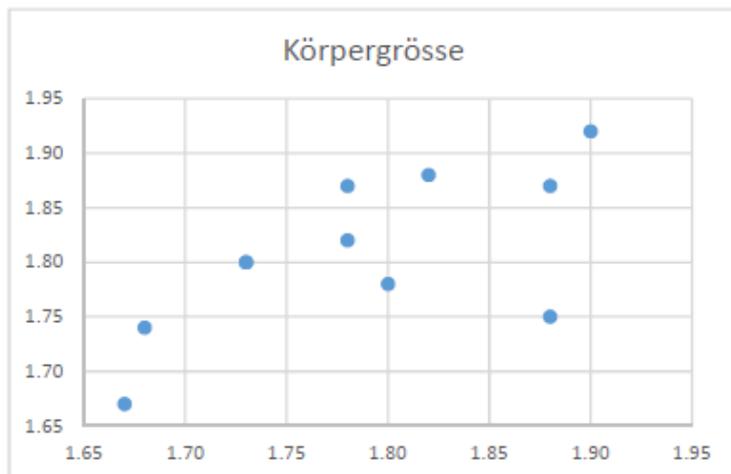


Sohn x	Vater y
1.67	1.67
1.68	1.74
1.73	1.80
1.73	1.80
1.78	1.87
1.78	1.82
1.80	1.78
1.82	1.88
1.88	1.75
1.88	1.87
1.90	1.92

Korrelation 0.669

#### Teilresultate

$\bar{x}$	1.786
$sd_x$	0.079
$\bar{y}$	1.809
$sd_y$	0.073
Kovarianz	0.004
Korrelation	0.669



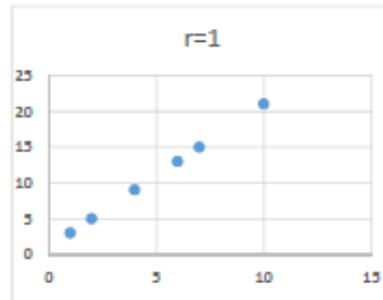
Hypothese:

Die Körpergrößen des Paares Vater-Sohn korrelieren schwächer als beim Paar Mutter-Tochter. Im Kapitel „Schliessende Statistik“ wird die Frage behandelt, wie eine solche Hypothese überprüft werden kann.

In den folgenden Abbildungen sind die Scatterplots von acht Stichproben mit verschiedenen Korrelationskoeffizienten dargestellt. Zusätzlich sind die Werte der Standardabweichungen von x bzw. y und die Kovarianz angegeben, aus denen sich die Korrelation gemäss 5) berechnen lässt.

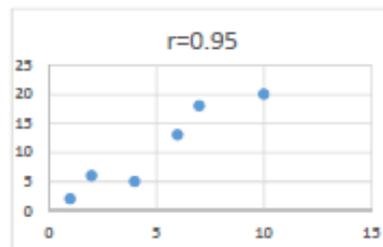
#### Scatterplot

	$x_i$	$y_i$
i		
1	1	3
2	2	5
3	4	9
4	6	13
5	7	15
6	10	21



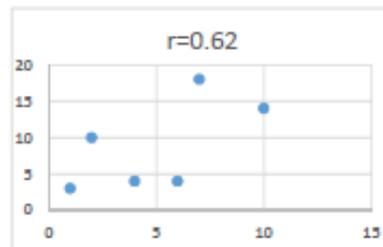
$\bar{x}$	5	$\bar{y}$	11.00
$s_x$	3.35	$s_y$	6.69
Kovarianz			22.40
Korrelation			1.00

	$x_i$	$y_i$
i		
1	1	2
2	2	6
3	4	5
4	6	13
5	7	18
6	10	20



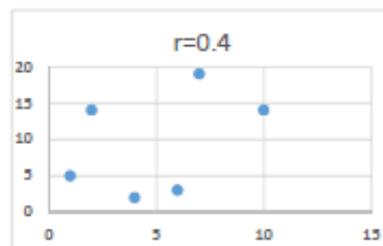
$\bar{x}$	5	$\bar{y}$	10.67
$s_x$	3.35	$s_y$	7.42
Kovarianz			23.60
Korrelation			0.95

	$x_i$	$y_i$
i		
1	1	3
2	2	10
3	4	4
4	6	4
5	7	18
6	10	14



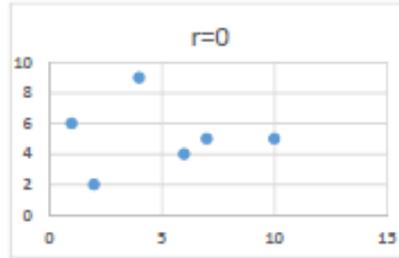
$\bar{x}$	5	$\bar{y}$	8.83
$s_x$	3.35	$s_y$	6.21
Kovarianz			12.80
Korrelation			0.62

	$x_i$	$y_i$
i		
1	1	5
2	2	14
3	4	2
4	6	3
5	7	19
6	10	14



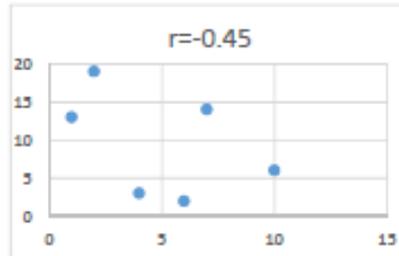
$\bar{x}$	5	$\bar{y}$	9.50
$s_x$	3.35	$s_y$	7.06
Kovarianz			9.40
Korrelation			0.40

i	$x_i$	$y_i$
1	1	6
2	2	2
3	4	9
4	6	4
5	7	5
6	10	5



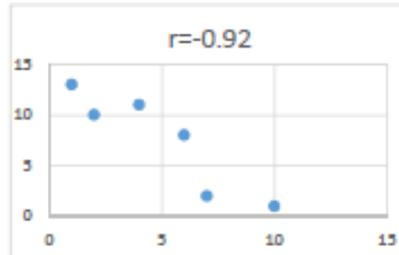
$\bar{x}$	5	$\bar{y}$	5.17
$s_x$	3.35	$s_y$	2.32
Kovarianz			0.00
Korrelation			0.00

i	$x_i$	$y_i$
1	1	13
2	2	19
3	4	3
4	6	2
5	7	14
6	10	6



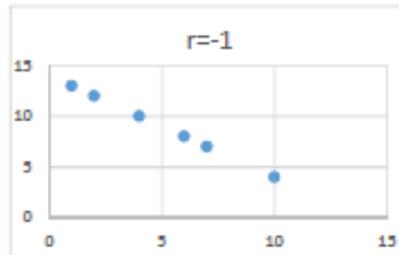
$\bar{x}$	5	$\bar{y}$	9.50
$s_x$	3.35	$s_y$	6.83
Kovarianz			-10.40
Korrelation			-0.45

i	$x_i$	$y_i$
1	1	13
2	2	10
3	4	11
4	6	8
5	7	2
6	10	1



$\bar{x}$	5	$\bar{y}$	7.50
$s_x$	3.35	$s_y$	4.93
Kovarianz			-15.20
Korrelation			-0.92

i	$x_i$	$y_i$
1	1	13
2	2	12
3	4	10
4	6	8
5	7	7
6	10	4



$\bar{x}$	5	$\bar{y}$	9.00
$s_x$	3.35	$s_y$	3.35
Kovarianz			-11.20
Korrelation			-1.00

## Eigenschaften der Produkt-Momenten-Korrelation

a)

Die Korrelation  $r_{xy}$  ist +1 für eine Gerade mit positiver Steigung, -1 für eine mit negativer Steigung. Für Werte von  $r_{xy}$  nahe bei +1 oder -1 liegen die Punkte eng um eine Gerade. Die Korrelation ist also ein Mass für die Stärke und Richtung eines linearen Zusammenhangs. Ist  $r_{xy} = 0$ , so besteht kein linearer Zusammenhang.

b)

Es gilt  $|r_{xy}| \leq 1$

Ein einfacher Beweis ergibt sich nach der sogenannten Cauchy-Schwarzschen Ungleichung, die von Ebene und Raum her bereits bekannt ist:

$$(\vec{u} \cdot \vec{v})^2 \leq (|\vec{u}| |\vec{v}|)^2$$

Zum Beweis wählt man für die Berechnung der Korrelation wie in 3.2 die standardisierten Variablen

$$u_i = \frac{x_i - \bar{x}}{s_x} \quad \text{bzw.} \quad v_i = \frac{y_i - \bar{y}}{s_y}$$

und bildet mit diesen Werten die Vektoren  $\vec{u}$  und  $\vec{v}$ . gilt

Wegen der Standardisierung folgt daraus für die Varianz nach 2)

$$s_u^2 = \frac{1}{n-1} |\vec{u}|^2$$

$$s_v^2 = \frac{1}{n-1} |\vec{v}|^2$$

und wegen

$$\left( \sum_{i=1}^n u_i \right) = \left( \sum_{i=1}^n v_i \right) = 0$$

nach 6) für die Kovarianz

$$s_{uv} = \frac{1}{n-1} \cdot (\vec{u} \cdot \vec{v})$$

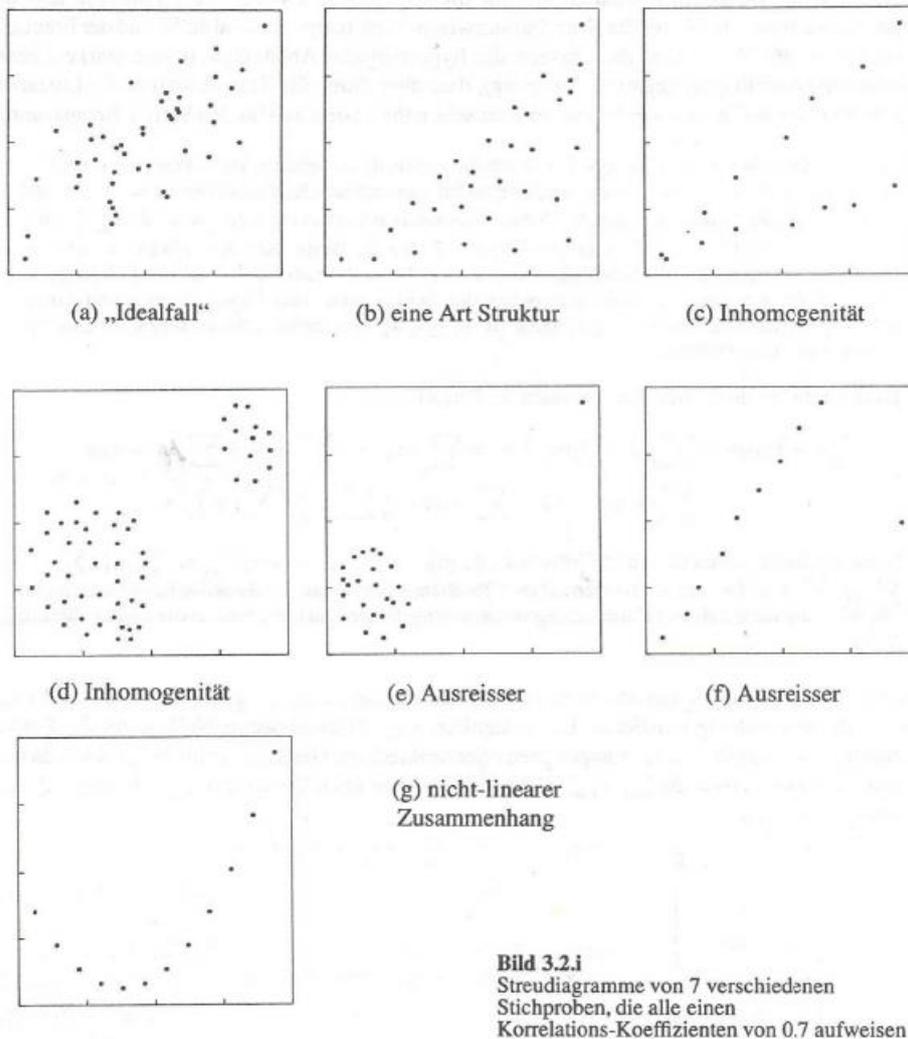
und schliesslich für die Produkt-Momenten-Korrelation nach 5) und Cauchy-Schwarz

$$r_{xy}^2 = r_{uv}^2 = \left( \frac{s_{uv}}{s_u s_v} \right)^2 = \frac{\frac{(\vec{u} \cdot \vec{v})^2}{(n-1)^2}}{\frac{1}{n-1} |\vec{u}|^2 \cdot \frac{1}{n-1} |\vec{v}|^2} \frac{(\vec{u} \cdot \vec{v})^2}{(|\vec{u}| \cdot |\vec{v}|)^2} \leq 1$$

c)

Die folgenden Streudiagramme stellen Stichproben mit gleichem Korrelationskoeffizienten 0.7 dar (Quelle: Chambers, Cleveland, Kleiner and Tukey 1983).

Es ist der grosse Einfluss von allfälligen Ausreissern zu erkennen. In Anwendungen sollte deshalb nicht nur die Korrelation allein, sondern auch das zugehörige Streudiagramm beachtet werden.



d)

Korrelationen sollten vorsichtig interpretiert werden. Selbst ein offensichtlicher Zusammenhang braucht keinen ursächlichen Zusammenhang zu bedeuten. Es kann nämlich sein, dass beide Größen von einer dritten Größe abhängen.

Dass die Anzahl der bei einem Brand eingesetzten Feuerwehrleute hoch korreliert ist mit dem angerichteten Schaden, bedeutet nicht, dass die Feuerwehrleute besser daheim bleiben würden.

Heikel sind insbesondere Zeitreihen mit linearer Komponente. Im bekannten Beispiel ist die hohe Korrelation zwischen der Zahl der Störche und der Zahl der Geburten zwischen 1900 und 1970 kein Beweis, dass der Klapperstorch die Babys bringt.

Es besteht z.B. eine positive Korrelation zwischen Körpergewicht und manueller Geschicklichkeit bei Schulkindern, denn Kinder werden mit steigendem Alter im Mittel schwerer und gleichzeitig geschickter.

Bei Männern besteht zwischen dem Einkommen und der Zahl der Haare auf dem Kopf eine negative Korrelation. Die erklärende Variable ist auch in diesem Fall das Alter.

### 3.3 Die lineare Regression

#### 3.3.1 Berechnung der Ausgleichsgeraden

Bei der Korrelation werden die beiden Grössen  $x$  und  $y$  gleich behandelt. In vielen Anwendungen ist allerdings die eine Grösse die Zielgrösse  $y$ , die andere die Ausgangsgrösse, erklärende Variable  $x$ .

Beispiele:

Erklärende Variable	Zielgrösse
Höhe des Wasserfalls	Schwingungsfrequenz
Spraydosen-Verbrauch	Ozongehalt
Geschwindigkeit	Bremsweg
Jahresnote	Maturanote

In vielen Fällen ist „im wesentlichen“  $y$  eine lineare Funktion von  $x$ . Damit stellt sich die Frage, wie eine Gerade rechnerisch möglichst gut an eine Stichprobe von Wertepaaren angepasst werden kann.

Die Steigung  $b$  (englisch: slope) und der  $y$ -Achsenabschnitt (englisch: intercept) sind so zu bestimmen, dass die Abweichungen der Punkte von der Geraden möglichst klein sind. Als Kriterium wird nach Gauss und Legendre (um 1800) die Methode der Kleinsten Quadrate (least squares) verwendet (Maximum likelihood):

Die Steigung  $b$  und der  $y$ -Achsenabschnitt  $a$  der Geraden sind so zu bestimmen, dass die folgende Quadratsumme minimal wird:

$$Q(a, b) = \sum_{i=1}^n e_i^2$$

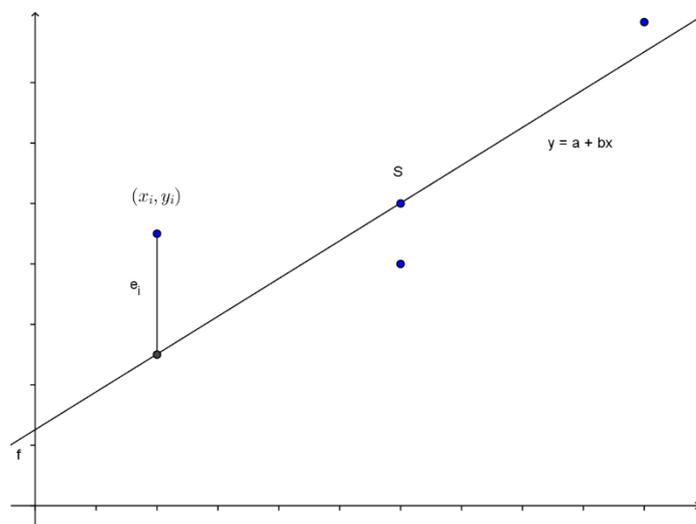
Es kann gezeigt werden, dass die Regressionsgerade durch den „Schwerpunkt“  $S(\bar{x}, \bar{y})$  der Punkte geht.

Die Lösung des Problems vereinfacht sich deshalb, wenn man wie bereits im Abschnitt Korrelation erläutert von den ursprünglichen Daten zu standardisierten Daten  $(u, v)$  übergeht. d.h. die Abweichungen der  $x$  bzw.  $y$ -Werte vom entsprechenden Mittelwert werden in Einheiten der Standardabweichung gemessen.

$$u_i = \frac{x_i - \bar{x}}{s_x}$$

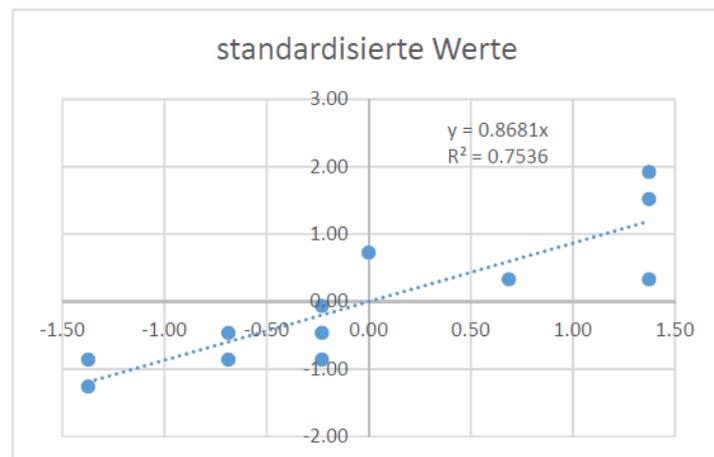
und

$$v_i = \frac{y_i - \bar{y}}{s_y}$$



Das folgende Beispiel dient der Illustration der Idee

$x_i$	$y_i$	standardisiert	
		$u_i$	$v_i$
4.7	3	-1.37	-0.86
5	3	-0.69	-0.86
5.2	4	-0.23	-0.46
5.2	5	-0.23	-0.07
5.9	10	1.37	1.92
4.7	2	-1.37	-1.26
5.9	9	1.37	1.52
5.2	3	-0.23	-0.86
5.3	7	0.00	0.73
5.9	6	1.37	0.33
5.6	6	0.69	0.33
5	4	-0.69	-0.46
$\bar{x}$	$\bar{y}$	$\bar{u}_i$	$\bar{v}_i$
5.30	5.17	0.00	0.00
$s_x$	$s_y$	$s_u$	$s_v$
0.44	2.52	1.00	1.00



Korrelation

Wertepaare	standardisiert
0.868	0.868

Steigung der Geraden

5.00	0.868
------	-------

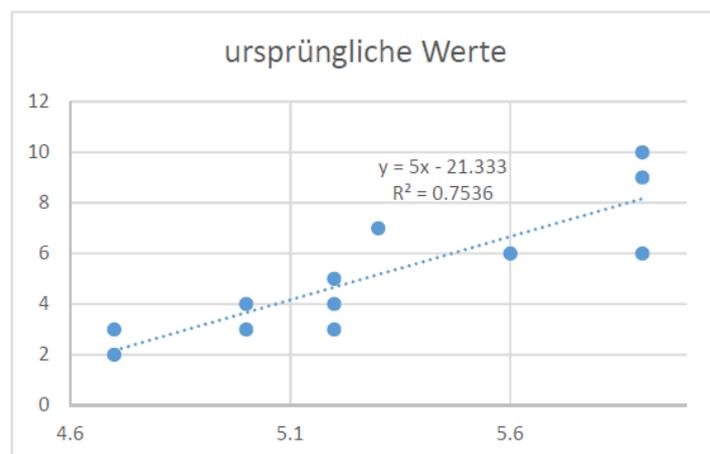
Achsenabschnitt

-21.33	0
--------	---

Bestimmtheitsmass

$r^2$	0.7536
-------	--------

$r^2$	0.7536
-------	--------



Bemerkung:

In der Abbildung ist ausser dem Quadrat des Korrelationskoeffizienten  $r_{xy}^2$ , das erst später definierte Bestimmtheitsmass  $R^2$  angegeben. Im Beispiel stimmt das Quadrat des Korrelationskoeffizienten  $r_{xy}^2$  mit dem Bestimmtheitsmass  $R^2$  überein. In 16) wird gezeigt, dass diese Aussage allgemein gilt.

Wie das Beispiel zeigt, bringt die Standardisierung die folgenden Vereinfachungen:

Die Mittelwerte sind 0:

$$\sum_{i=1}^n u_i = n \cdot \bar{u} = 0$$

und

$$\sum_{i=1}^n v_i = n \cdot \bar{v} = 0$$

Die Standardabweichungen sind 1 und wegen 2) gilt:

$$(n-1)s_u^2 = n-1 = \sum_{i=1}^n u_i^2 - n \cdot \bar{u}^2$$

oder also

$$\boxed{\sum_{i=1}^n u_i^2 = n-1} \quad 8)$$

und entsprechend

$$(n-1)s_v^2 = n-1 = \sum_{i=1}^n v_i^2 - n \cdot \bar{v}^2$$

oder also

$$\boxed{\sum_{i=1}^n v_i^2 = n-1} \quad 8')$$

und nach Definition der Korrelation zwischen  $u_i$  und  $v_i$  nach 7)

$$\boxed{r_{uv} = \frac{1}{n-1} \sum_{i=1}^n u_i \cdot v_i} \quad 9)$$

Da im  $(u, v)$ -System die Regressionsgerade durch den Ursprung, den Schwerpunkt geht, kann sie in der Form  $v = mu$  angesetzt werden. Die Steigung  $m$  ist so zu bestimmen, dass die folgende Quadratsumme minimal wird:

$$Q(m) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (v_i - m \cdot u_i)^2$$

Nach der Kettenregel gilt dann für die 1. Ableitung von  $Q$ :

$$Q'(m) = -2 \sum_{i=1}^n (v_i - m \cdot u_i) \cdot u_i = 0$$

und nach Ausmultiplizieren und Vereinfachen

$$\sum_{i=1}^n u_i v_i = m \sum_{i=1}^n u_i^2$$

und wegen der Standardisierung nach 9) und 8)

$$(n-1)r_{uv} = (n-1) \cdot m$$

und nach Division durch  $(n-1)$

$$m = r_{uv}$$

Für die standardisierten Daten stimmt also die Steigung  $m$  der Regressionsgeraden mit der Korrelation  $r_{uv}$  überein.

Die Steigung der Regressionsgeraden für die Originaldaten ergibt sich, indem man die Standardisierung rückgängig macht.

$$r_{uv} = m = \frac{v_2 - v_1}{u_2 - u_1} = \frac{\frac{y_2 - \bar{y}}{s_y} - \frac{y_1 - \bar{y}}{s_y}}{\frac{x_2 - \bar{x}}{s_x} - \frac{x_1 - \bar{x}}{s_x}} = \frac{y_2 - y_1}{x_2 - x_1} \cdot \frac{s_x}{s_y} = b \cdot \frac{s_x}{s_y}$$

Löst man diese Gleichung nach b auf, so erhält man für die Steigung der Regressionsgeraden zu den Originaldaten

$$b = r_{uv} \cdot \frac{s_y}{s_x} \quad \text{Steigung der Regressionsgeraden}$$

Da sich bei einer linearen Transformation die Korrelation nicht ändert, gilt  $r_{uv} = r_{xy}$  und damit

$$b = r_{xy} \cdot \frac{s_y}{s_x} \quad \text{Steigung der Regressionsgeraden} \quad 10)$$

Da für die Korrelation gemäss 5') gilt:

$$r_{xy} = \frac{s_{xy}}{s_x \cdot s_y}$$

erhält man

$$b = r_{xy} \cdot \frac{s_y}{s_x} = \frac{s_{xy}}{s_x \cdot s_y} \cdot \frac{s_y}{s_x} = \frac{s_{xy}}{s_x^2}$$

oder

$$b = \frac{s_{xy}}{s_x^2} \quad \text{oder} \quad s_{xy} = b \cdot s_x^2 \quad 11)$$

Damit ist auch der y-Achsenabschnitt bestimmt, denn da der „Schwerpunkt“ S ( $\bar{x}, \bar{y}$ ) auf der Regressionsgeraden liegt gilt:

$$\bar{y} = a + b\bar{x} \quad \text{oder} \\ a = \bar{y} - b\bar{x}$$

Zusammenfassung:

Die Regressionsgerade hat die Steigung

$$b = \frac{s_{xy}}{s_x^2} = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad 12)$$

und den y-Achsenabschnitt

$$a = \bar{y} - b\bar{x} \quad 13)$$

Bemerkung:

Ist die Kettenregel nicht bekannt, so kann die Steigung  $m$  auch folgendermassen bestimmt werden:

Multipliziert man bei  $Q(m)$  die Summanden nach der binomischen Formel aus, so erhält man eine in  $m$  quadratische Funktion der Form

$$Q(m) = a \cdot m^2 + bm + c$$

Der Graph von  $Q$  ist eine nach oben geöffnete Parabel, deren Scheitelabszisse gleich dem ersten Summanden in der quadratischen Auflösungsformel ist. Für das Minimum von  $Q$  ergibt sich damit

$$m = -\frac{b}{2a}$$

in Übereinstimmung mit 11)

### 3.3.2 Das Bestimmtheitsmass

Die Abweichungen  $e_i = y_i - \hat{y}_i$  zwischen den beobachteten Werten  $y_i$  und den durch die Regressionsgerade vorhergesagten Werten  $\hat{y}_i$  heissen Residuen. Mit ihnen kann punktweise überprüft werden, wie gut das Modell mit den Daten übereinstimmt. Es fehlt aber ein Mass um die Güte des Regressionsmodells insgesamt beurteilen zu können. Ein solches Mass kann über die sogenannte Streuungszersetzung erhalten werden.

Dazu führt man die folgenden Summen ein:

$SQT$  („Sum of squares Total“) ist die Gesamtstreuung,

$$SQT = \sum_{i=1}^n (y_i - \bar{y})^2$$

$SQE$  („Sum of squares Explained“) ist durch die Regression erklärte Streuung

$$SQE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$SQR$  (Sum of squares Residual) ist durch die Rest- oder Residualstreuung erklärte Streuung.

$$SQR = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Für diese drei Summen gilt die folgende

Streuungszerlegung

$$SQT = SQE + SQR \text{ oder } SQR = SQT - SQE$$

also

$$\boxed{\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (\hat{y}_i - \bar{y})^2}$$

14)

wobei für die  $y$ -Werte der Regressionsgeraden gilt:  $\hat{y}_i = a + b \cdot x_i$

Beweis:

Aus

$$y_i - \hat{y}_i = y_i - (a + bx_i)$$

und

$$a = \bar{y} - b\bar{x}$$

folgt

$$y_i - \hat{y}_i = y_i - (\bar{y} - b\bar{x} + bx_i) = (y_i - \bar{y}) - b(x_i - \bar{x})$$

und damit nach dem binomischen Lehrsatz

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 - 2b \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) + b^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

Da wegen 11) gilt:

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = b \sum_{i=1}^n (x_i - \bar{x})^2$$

erhält man schliesslich

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 - b^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

SQR = SQT - SQE

oder auch

$$b^2 \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{SQE} = \text{SQT} - \text{SQR}$$

Wegen

$$(\hat{y}_i - \bar{y})^2 = (bx_i + a - \bar{y})^2 = (bx_i + a - (b\bar{x} + a))^2 = b^2(x_i - \bar{x})^2$$

gilt also wie behauptet

$$\text{SQR} = \text{SQT} - \text{SQE}$$

Als Masszahl für die Güte des Modells wird das Bestimmtheitsmass  $R^2$  (auch Determinationskoeffizient) verwendet. Das Bestimmtheitsmass gibt an, welchen Anteil der Gesamtstreuung die durch die Regression erklärte Streuung ausmacht.

$$R^2 = \frac{\text{SQE}}{\text{SQT}} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad 15)$$

Aus der Definition folgt:  $0 \leq R^2 \leq 1$

Spezialfälle:

Für  $R = 0$  ist das Modell denkbar schlecht

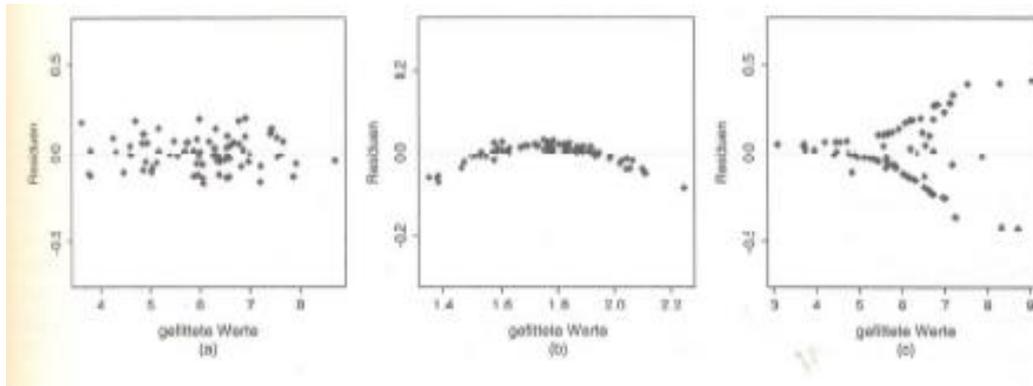
Für  $R = 1$  liegen die Datenpunkte auf einer Geraden.



Bei der Beurteilung der Güte des Modells sollte aber nicht allein das Bestimmtheitsmass, sondern auch das Muster der Residuen betrachtet werden.

Das Verhalten in der Abbildung a) ist ideal, denn die Residuen schwanken unsystematisch und sie sind nahezu 0.

Im Gegensatz dazu verändern sich die Residuen bei b) und c) je nach der Einflussgrösse. Eine nicht lineare Abhängigkeit wird durch das Modell nicht beschrieben.



### Beispiele aus verschiedenen Anwendungsgebieten

#### a) Finanzwirtschaft: Das CAP-Modell (Quelle: Fahrmeier)

Bei Anwendungen im Finanzbereich wird die Steigung  $b$  der Regressionsgeraden mit  $\beta$  und der y-Achsenabschnitt  $a$  mit  $\alpha$  bezeichnet.

Das „Capital Asset Pricing Model“ wird verwendet, um das Risiko von verschiedenen Aktien zu vergleichen. Dazu wird auf der Basis von Finanzmarktdaten über einen längeren Zeitraum für jede Aktie der sogenannte Beta-Faktor bestimmt. Mit diesem Faktor kann das Risiko dieser Aktie gemessen am Risiko des Gesamtmarkts beurteilt werden. Beispielsweise bedeutet ein Beta-Faktor grösser als eins, dass das Risiko der Aktie grösser als das Marktrisiko ist. Es interessiert der Zusammenhang zwischen der Einflussgrösse  $x$  (Marktrisiko minus Zins) und dem abhängigen Merkmal  $y$  (Rendite der Aktie minus Zins).

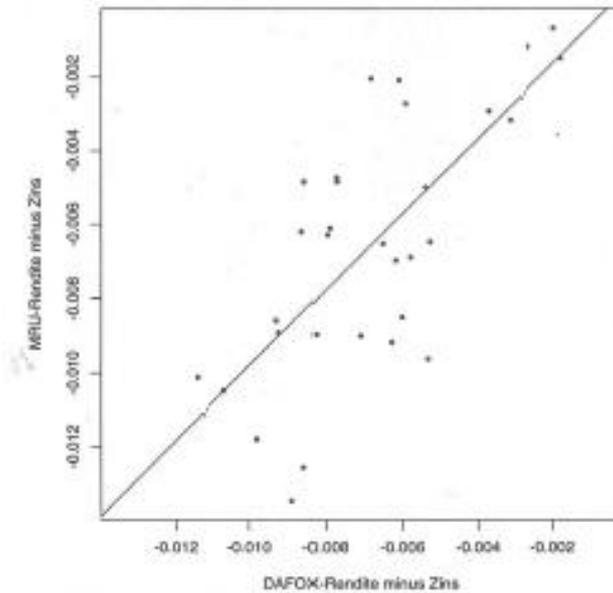


ABBILDUNG 3.20: Streudiagramm und Ausgleichsgerade für das CAP-Modell

Im abgebildeten Beispiel ergibt sich aufgrund der Daten die Regressionsgerade  
 $y = 0.0004 + 1.0216x$

Das Ergebnis bedeutet, dass im Mittel

bei einer Änderung der Marktrendite  $x$  um ein Prozent sich die Rendite  $y$  der Aktie um 1.02% also wegen  $\beta \approx 1.02 \approx 1$  etwa auch um ein Prozent verändert.

Die Residuen enthalten die Renditekomponenten der Aktie, die sich nicht durch den Markt erklären lassen. Diese Aktie hat also ein mit dem Markt vergleichbares Risiko. Das Bestimmtheitsmass  $R^2 = 0.5$  bedeutet, dass 50% der Renditeschwankungen der Aktie durch Schwankungen des Marktes bedingt sind und  $1 - 0.5$  d.h. 50% titelspezifischer Natur sind. Der Korrelationskoeffizient  $r_{xy}$  zwischen den Marktrendite und Rendite der Aktie beträgt  $\sqrt{0.5} \approx 71\%$

Das Residuenmuster zeigt keine besonderen Auffälligkeiten

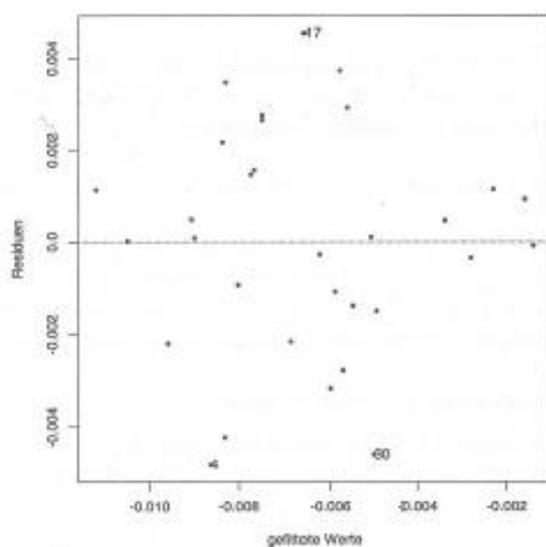


ABBILDUNG 3.21: Residualplot für dasCAP-Modell

## b) Physik: Die Zustandsgleichung für ideale Gase

Bei konstantem Druck  $p$  ist das Volumen  $V$  eines Gases nach dem Gesetz von Gay-Lussac proportional zur absoluten Temperatur  $T$  d.h. es gilt:

$$V = cT$$

Die Proportionalitätskonstante  $c$  ist vom konstanten Druck und von der vorliegenden Gasmenge abhängig. Misst man die Temperatur in  $^{\circ}\text{C}$ , dann ist das Volumen also eine lineare Funktion der Temperatur

$$\vartheta = T - T_n$$

wobei  $T_n$  die sogenannte Normtemperatur 273.15 bedeutet.

Für das Volumen gilt damit

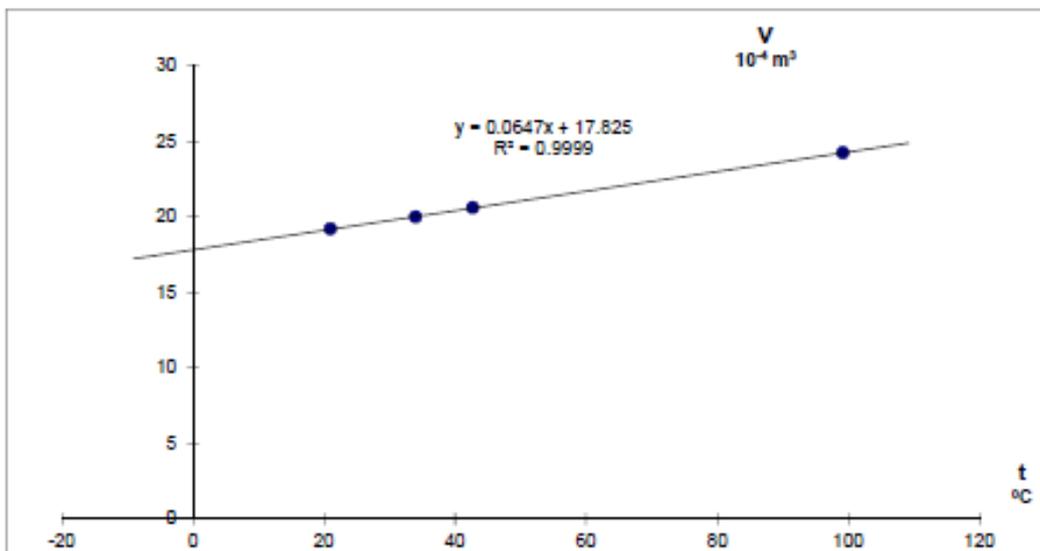
$$V(\vartheta) = V_0 + b \cdot \vartheta$$

Dabei ist  $V_0$  das Volumen bei  $0^{\circ}\text{C}$  (der y-Achsenabschnitt der Regressionsgeraden) und  $b$  die Steigung der Regressionsgeraden. Die folgenden Daten stammen von einem Versuch im Physikunterricht an der KSZ (ac)

**Volumenausdehnung**

2C 26.6.1994

Nr.	Temp. T $^{\circ}\text{C}$	V $\text{cm}^3$	$V(t) = V_0 + \gamma V_0 t$			
1	21.0	19.20	$m = \gamma V_0$	$V_0$	$\gamma$	$T_0$
2	34.0	20.00	0.06474905	17.8252282	0.00363	-275.3
3	42.7	20.60				
4	99.2	24.25				



Die bekannte Rechnung ergibt für  $b = 0.06474905$  und für  $a = V_0 = 17.8252282$ .

Die Ausgleichsgerade schneidet die Temperaturachse bei  $-275.3^{\circ}\text{C}$ , dem absoluten Nullpunkt der Temperatur. Ein genauere Wert ist  $-273.15^{\circ}\text{C}$ .

Besteht zwischen Einflussgrösse und Zielgrösse kein linearer Zusammenhang, dann kann dies eine geeignete Transformation der Daten ermöglichen.

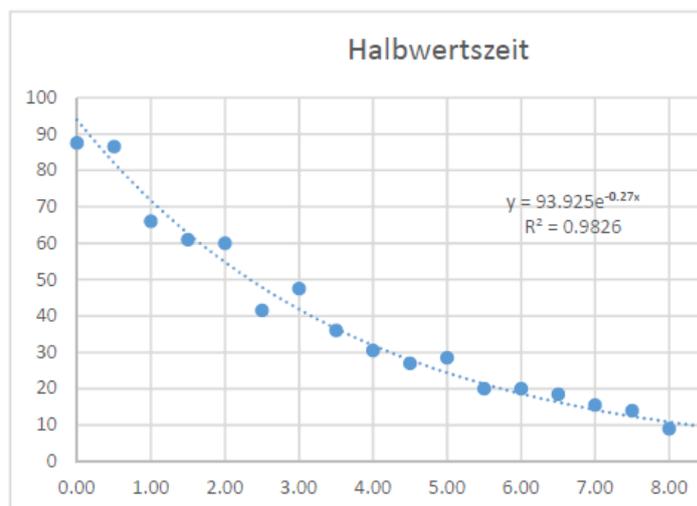
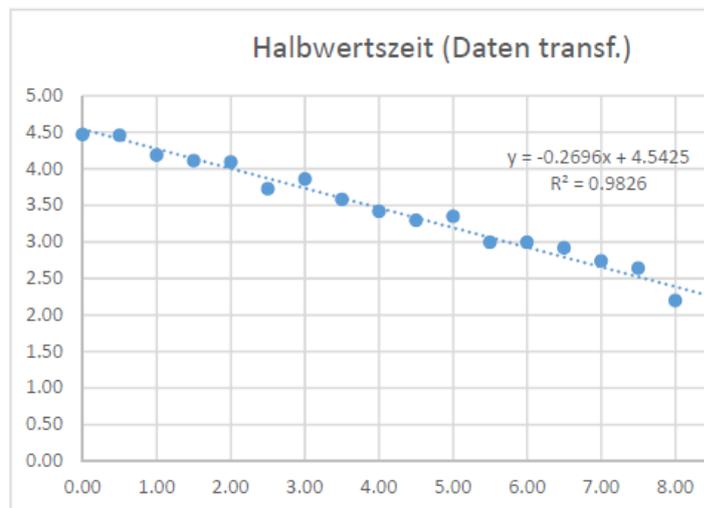
Wachstumsgesetz	Einflussgrösse	Zielgrösse
Exponentialfunktion	$\ln(\text{Einflussgrösse})$	Zielgrösse
Potenzfunktion	$\ln(\text{Einflussgrösse})$	$\ln(\text{Zielgrösse})$

c) Halbwertszeit von Ba-137 (SPAM-Maturaaufgabe an der AKSA 2007)

Zur Bestimmung der Halbwertszeit von Ba-137 wurde im Abstand von 30 Sekunden die Impulsrate pro Sekunde gemessen. Dabei ergaben sich die in der Tabelle ermittelten Werte:

Halbwertszeit

Zeit [min]	Impulsrate	$\ln(r)$
t	r [s <sup>-1</sup> ]	
0.00	87.5	4.47
0.50	86.5	4.46
1.00	66	4.19
1.50	61	4.11
2.00	60	4.09
2.50	41.5	3.73
3.00	47.5	3.86
3.50	36	3.58
4.00	30.5	3.42
4.50	27	3.30
5.00	28.5	3.35
5.50	20	3.00
6.00	20	3.00
6.50	18.5	2.92
7.00	15.5	2.74
7.50	14	2.64
8.00	9	2.20
8.50	9	2.20
9.00	7.5	2.01



Für den radioaktiven Zerfall gilt das Gesetz:

$$N(t) = N(0) \cdot e^{-kt}$$

Logarithmiert man diese Gleichung, dann gilt:

$$\ln(N(t)) = \ln(N(0)) - kt$$

d.h. die logarithmierten Funktionswerte sind eine lineare Funktion der Zeit. Die Ausgleichsgerade für die Punkte  $(t, \ln(N(t)))$  hat die Steigung

$b = -0.2695772$  und den y-Achsenabschnitt  $a = 4.5424927$

Wegen  $\ln(N(0)) = a$  und  $-k = b$  folgt

$$N(0) = e^{4.5425} \approx 93.925$$

und damit das Zerfallsgesetz

$$N(t) = N(0) \cdot e^{-kt} \approx 93.925 \cdot e^{-0.26958 \cdot t}$$

Die Halbwertszeit  $\tau$ , die auch grafisch bestimmt werden kann, ergibt sich wegen

$$N(\tau) = \frac{1}{2} \cdot N(0) = N(0) \cdot e^{-k\tau}$$

nach Division durch  $N(0)$

$$\frac{1}{2} = e^{-k\tau}$$

zu

$$\tau = \frac{\ln 2}{k} \approx 2.57 \text{ min}$$

d) Barometerdruck in verschiedenen Höhenlagen

Quelle: Lambert, Beiträge zum Gebrauch der Mathematik, Theil 1, Berlin 1765

Es wurden die Höhen über Meer verschiedener Berge im Languedoc, der Auvergne und der Provence gemessen.

Längenmass: 1 Toise entspricht 1.494 m

Barometer: Linie

Lambert interessierte sich für die Form der Abhängigkeit des Drucks von der Höhe. Im folgenden Beispiel interessiert umgekehrt die für einen Höhenmesser wichtige Abhängigkeit der Höhe vom Druck. In den ersten drei Spalten der Tabelle sind die von Lambert erfassten Werte aufgeführt.

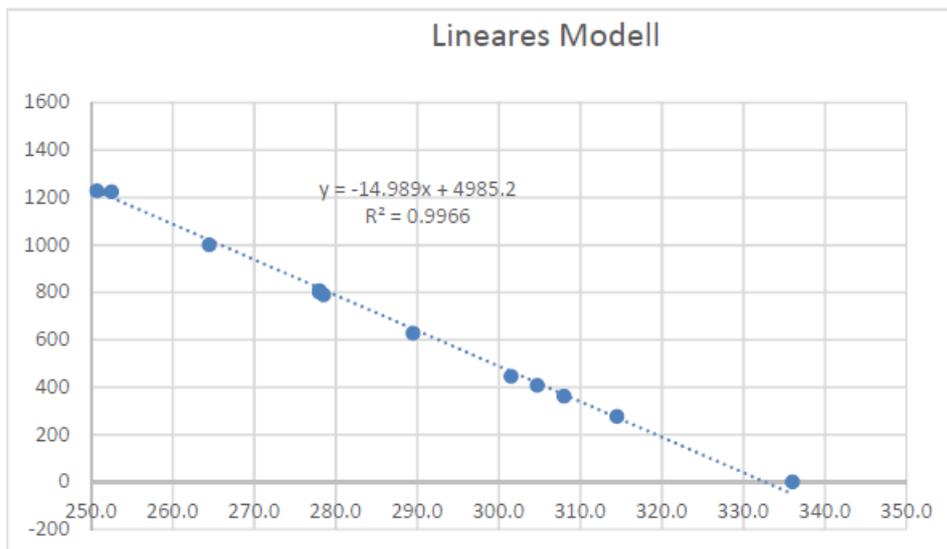
Wählt man für die Daten ein lineares Modell, dann ist die Höhe des Berges eine lineare Funktion des Barometerdrucks  $p$  oder

$$h = b \cdot p + a$$

### Barometerdruck in verschiedenen Höhenlagen

Lineares Modell			a	b
			4985.2	-14.989
Ort	Druck p	Höhe h	erwartet	Residuen
Meereshöhe	336.0	0	-50.99	50.99
Clairret	314.5	277	271.27	5.73
Rodez	308	362	368.69	-6.69
Massane	304.7	408	418.16	-10.16
Rupeyroux	301.5	446	466.12	-20.12
Bugarac	289.5	628	645.98	-17.98
Puy du Dôme	278.5	789	810.86	-21.86
La Coste	278	807	818.35	-11.35
La Courlande	278	801	818.35	-17.35
Mont d'Or	264.5	1001	1020.70	-19.70
St. Barthélémy	252.5	1225	1200.56	24.44
Mousset	250.7	1228	1227.54	0.46
Canigou	240.5	1424	1380.42	43.58

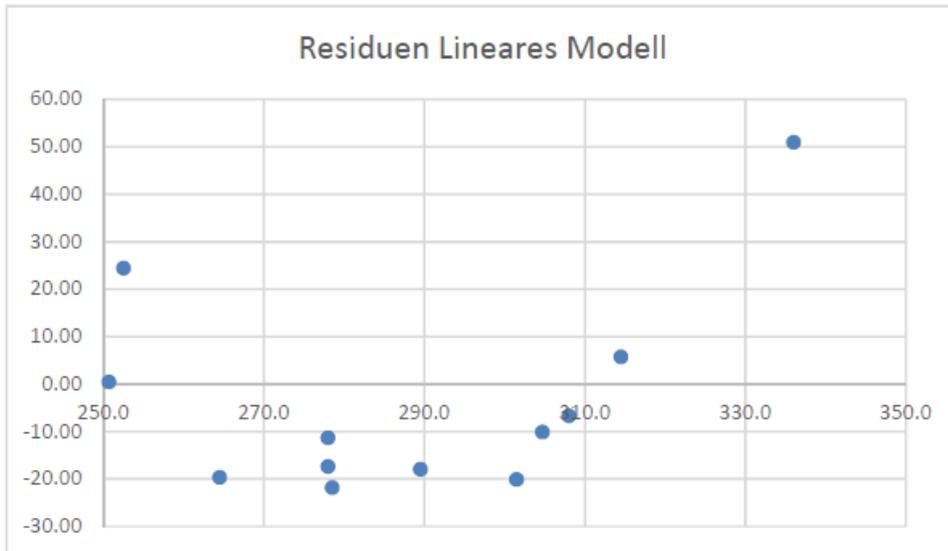
Die zugehörigen Punkte sind in der folgenden Abbildung dargestellt:



Wie in den bisherigen Beispielen können der y-Achsenabschnitt a und die Steigung b bestimmt werden mit dem folgenden Ergebnis:  $a = -14.989$  und  $b = 4985.2$ .

Mit diesen Werten können die nach dem Modell geschätzten Höhen bestimmt werden, die in der vierten Spalte der Tabelle dargestellt sind. Der Skizze entnimmt man, dass das Modell die Daten scheinbar gut beschreibt. Zudem ist auch das Bestimmtheitsmass sehr hoch.

Berücksichtigt man hingegen die in der fünften Spalte aufgeführten Residuen, so zeigen diese ein auffälliges Muster:



Masszahlen oder grafische Darstellungen sollten nicht überbewertet werden. Es ist auch zu überprüfen, ob die Residuen - wie in diesem Beispiel- allenfalls von der Einflussgrösse abhängen.

Nun wird die die Abhängigkeit des Drucks von der Höhe ein exponentielles Modell gewählt:

$$p(h) = p_0 \cdot e^{-kh}$$

Logarithmiert man diese Gleichung so erhält man:

$$\ln(p(h)) = \ln(p_0) - kh$$

Oder nach h aufgelöst:

$$\ln(p(h)) - \ln(p_0) = -kh$$

bzw.

$$\ln(p_0) - \ln(p(h)) = kh$$

oder schliesslich nach Division durch k

$$h = \frac{1}{k} \ln(p_0) - \frac{1}{k} \ln(p(h))$$

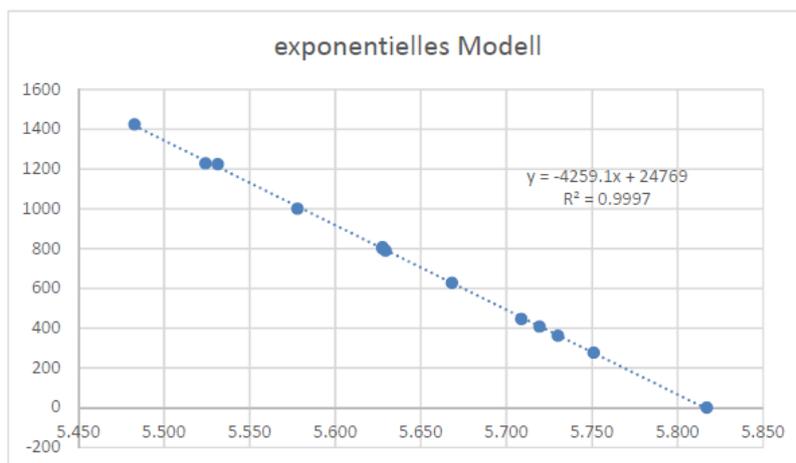
Bei dieser Transformation der Daten ist also die Höhe eine lineare Funktion von  $\ln(p(h))$

## Barometerdruck in verschiedenen Höhenlagen

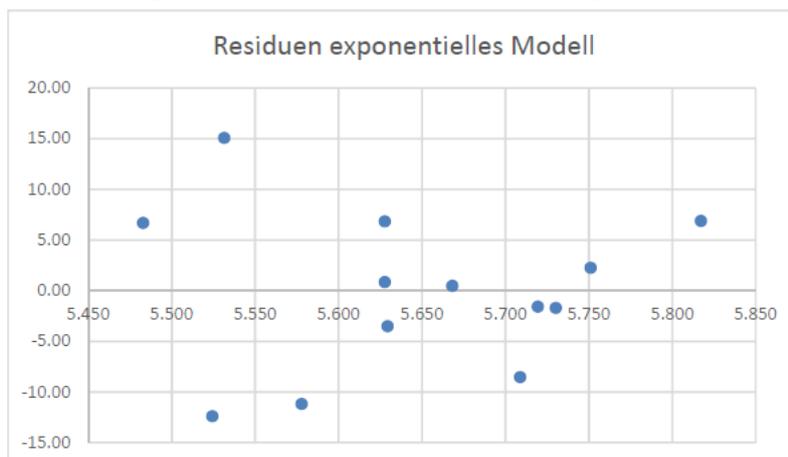
### Exponentielles Modell

Ort	Druck p	$\ln(p)$	Höhe h	erwartet	Residuen
Meereshöhe	336.0	5.817	0	-6.88	6.88
Clairret	314.5	5.751	277	274.76	2.24
Rodez	308	5.730	362	363.71	-1.71
Massane	304.7	5.719	408	409.59	-1.59
Rupeyroux	301.5	5.709	446	454.55	-8.55
Bugarac	289.5	5.668	628	627.54	0.46
Puy du Dôme	278.5	5.629	789	792.52	-3.52
La Coste	278	5.628	807	800.17	6.83
La Courlande	278	5.628	801	800.17	0.83
Mont d'Or	264.5	5.578	1001	1012.19	-11.19
St. Barthélemy	252.5	5.531	1225	1209.94	15.06
Mousset	250.7	5.524	1228	1240.41	-12.41
Canigou	240.5	5.483	1424	1417.32	6.68

In der Abbildung sind die Punkte  $(\ln(p(h)), h)$  dargestellt:



Für die Ausgleichsgerade ergeben sich die folgenden Werte:  $a = 24768.7$  und  $b = -4259.1$



Das Muster der Residuen zeigt keine Auffälligkeiten.

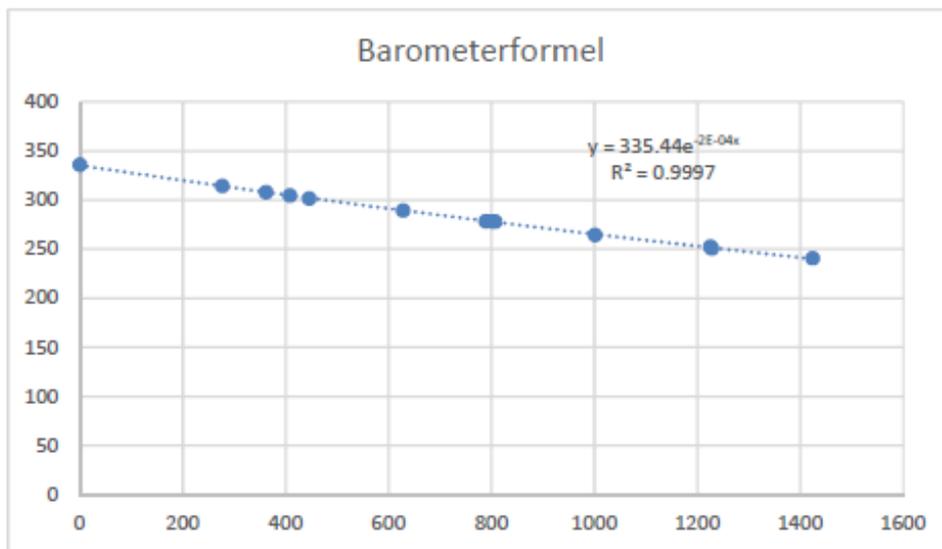
Die Parameter der Barometerformel ergeben sich aus der Steigung und dem y-Achsenabschnitt der Ausgleichsgeraden.

$$-\frac{1}{k} = b \text{ oder } k = -\frac{1}{b} \approx 0.00023479$$

$$\frac{1}{k} \ln(p_0) = a \text{ oder } \ln(p_0) = k \cdot a \approx 5.8155 \text{ und } p_0 \approx 335.458$$

Damit gilt mit den erwähnten Einheiten die folgende Barometerformel:

$$p(h) = p_0 \cdot e^{-kh} = 335.46 \cdot e^{-0.00023479 \cdot h}$$





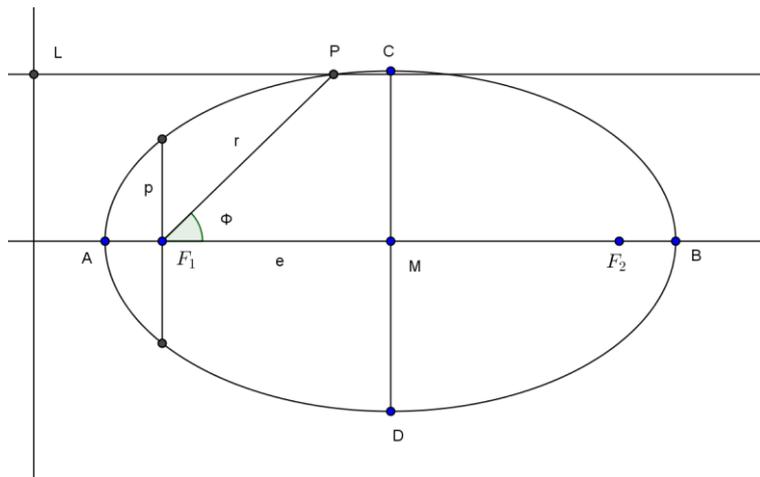
## f) Der Komet Tentax

Der 1968 entdeckte Komet Tentax bewegt sich im Sonnensystem. In einem geeigneten Polarkoordinatensystem  $(r, \varphi)$  wurden folgende Positionen des Kometen beobachtet:

$\varphi^\circ$	$\varphi$ (rad)	$r$
48.0	0.8	2.70
67.0	1.2	2.00
83.0	1.4	1.61
108.0	1.9	1.20
126.0	2.2	1.02

Falls die Einflüsse der Planeten vernachlässigt werden, bewegt sich der Komet nach dem Keplerschen Gesetz auf einer elliptischen oder hyperbolischen Bahn mit der Gleichung

$$r = \frac{p}{1 - e \cdot \cos(\varphi)}$$



Der Parameter  $p$  und die lineare Exzentrizität  $e$  können aus den Daten geschätzt werden. Dazu wird die Kegelschnittgleichung als lineare Gleichung umgeschrieben.

Wird

$$r \cdot (1 - e \cdot \cos(\varphi)) = r - e \cdot r \cdot \cos(\varphi) = p$$

nach  $r$  aufgelöst

$$r = p + r \cdot \cos(\varphi) \cdot e$$

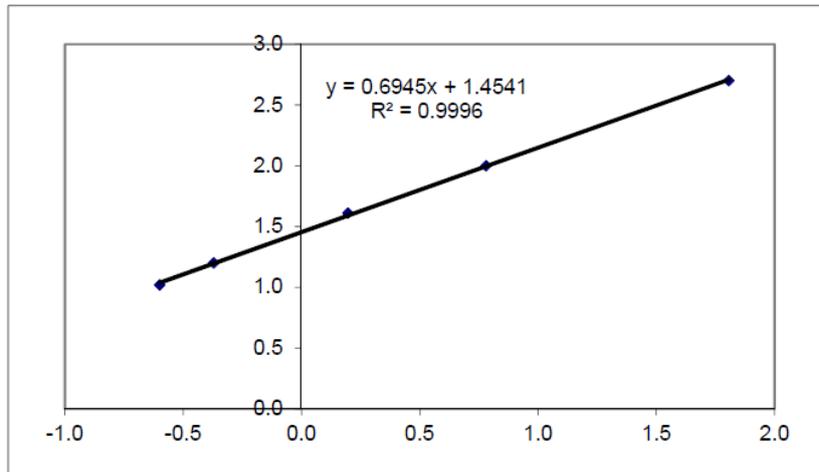
so folgt mit  $x = r \cdot \cos(\varphi)$  und  $y = r$  die lineare Gleichung

$$y = p + e \cdot x$$

Die lineare Exzentrizität beträgt  $e \approx 0.6945$  und der Parameter  $p \approx 1.4541$

Wegen  $e < 1$  bewegt sich also Tentax auf einer elliptischen Bahn mit dem Parameter  $p$

$\varphi^\circ$	$\varphi$ (rad)	r	$x=r \cdot \cos(\varphi)$	$y=r$
48.0	0.8	2.70	1.807	2.70
67.0	1.2	2.00	0.781	2.00
83.0	1.4	1.61	0.196	1.61
108.0	1.9	1.20	-0.371	1.20
126.0	2.2	1.02	-0.600	1.02



### Übungsaufgabe:

J. Kepler (1571 – 1630) hat zu Beginn des 17. Jahrhunderts die Bewegungen der Planetenbahnen untersucht und seine berühmten „Gesetze“ entdeckt. Die Planetenbahnen sind Ellipsen (mit der Sonne im Brennpunkt) und der Verbindungsstrahl Sonne – Erde überstreicht in gleichen Zeiten gleiche Flächen. Das dritte Keplersche Gesetz stellt einen Zusammenhang her zwischen der Halbachse  $a$  der Ellipse und der Umlaufzeit  $T$  her.  $T$  ist eine Potenzfunktion von  $a$

$$T = \alpha \cdot a^\beta$$

$a$  in astronomischen Einheiten,  $T$  in Erdjahren.

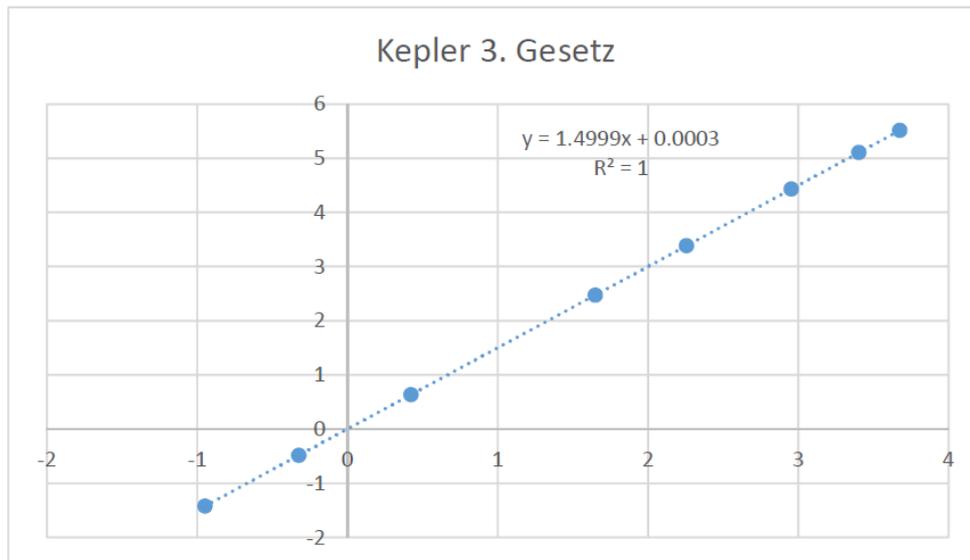
Es wurden folgende Daten gemessen (Datenquelle nicht bekannt):

Planet	$a$	$T$
Merkur	0.387	0.241
Venus	0.723	0.615
Mars	1.524	1.881
Jupiter	5.203	11.862
Saturn	9.539	29.458
Uranus	19.182	84.014
Neptun	30.058	164.793
Pluto	39.44	247.7

Wegen

$$\ln(T) = \ln(\alpha) + \beta \cdot \ln(a)$$

besteht zwischen  $\ln(T)$  und  $\ln(a)$  wie die folgende Abbildung zeigt eine lineare Beziehung:



Die Ausgleichsgerade hat die Steigung  $\beta \approx 3/2$  und den y-Achsenabschnitt  $0.0003 \approx \ln(\alpha)$ . Es gilt also näherungsweise wegen  $\alpha \approx 1$

$$T \approx a^{3/2}$$

oder

$$\frac{T^2}{a^3} \approx \textit{konstant}$$

Heute ist bekannt, dass die Keplerschen Gesetze die Bewegungen der Planeten nur näherungsweise beschreiben.