

Regression und Korrelation (Gurtner)

Oft muss man Trends verlängern. Dazu nimmt man die Daten einer Entwicklung, legt eine passende Kurve durch und verlängert die Kurve in die Zukunft. Hier verwenden wir nur die Gerade als Kurve – das nennt sich dann **Regressionsgerade** oder Ausgleichsgerade (in EXCEL: Trendkurve)

Beispiel:

Folgender (fiktiver) Zusammenhang zwischen den Geburtszahlen und der Anzahl der Störche in Rust möge für verschiedene Jahre gegeben sein:

Störche	8	7	5	3	4
Geburten	12	10	8	5	6

Untersuchen Sie, ob ein linearer Zusammenhang zwischen der Anzahl der Geburten und der Anzahl der Störche besteht, um die These zu belegen, dass die Störche die Babys bringen (!)
 Berechnen Sie die Korrelation und die Regressionsgrade und bestimmen Sie damit näherungsweise für 10 Störche die Anzahl der Geburten.
 Zeichnen Sie für die Daten und die Gerade ein Diagramm!

Lösung:

Dazu müssen wir eine Tabelle erstellen, welche die Daten und deren Quadrate und Produkte enthält und zum Schluss die Summen und Mittelwerte:

	xi	yi	xi ²	xi*yi	yi ²
	8	12	64	96	144
	7	10	49	70	100
	5	8	25	40	64
	3	5	9	15	25
	4	6	16	24	36
Summen	27	41	163	245	369
Mittelwert (:5)	5,4	8,2	32,6	49	73,8

Dadurch lassen sich die **Varianzen** von x und y und die **Kovarianz** berechnen:

$$V(x) = \sigma_x^2 = \text{Mittelwert der x-Quadrate minus Quadrat des x-Mittelwerts} = 32,6 - 5,4^2 = \mathbf{3,44}$$

$$V(y) = \sigma_y^2 = \text{Mittelwert der y-Quadrate minus Quadrat des y-Mittelwerts} = 73,8 - 8,2^2 = \mathbf{6,56}$$

$$\text{coV}(xy) = \sigma_{xy} = \text{Mittelwert der Produkte minus Produkt der Mittelwerte} = 49 - 5,4 \cdot 8,2 = \mathbf{4,72}$$

Damit kann man die Werte k und d der Regressionsgeraden bestimmen:

$$k = \text{coV}(xy) / V(x) = 4,72 / 3,44 = \mathbf{1,37}$$

$$d = \bar{y} - k \cdot \bar{x} = 8,2 - 1,372 \cdot 5,4 = \mathbf{0,79}$$

→ $y = \mathbf{1,37 \cdot x} + \mathbf{0,79}$ ist die Regressionsgerade.

Um die **Gerade einzeichnen** zu können, bestimmen wir die Werte am Beginn und Ende der Tabelle:

$$\text{für } x = 3 \text{ gilt: } y = 1,37 \cdot 3 + 0,79 = 4,9$$

$$\text{für } x = 8 \text{ gilt: } y = 1,37 \cdot 8 + 0,79 = 11,7$$

Für die **Korrelation** gilt:

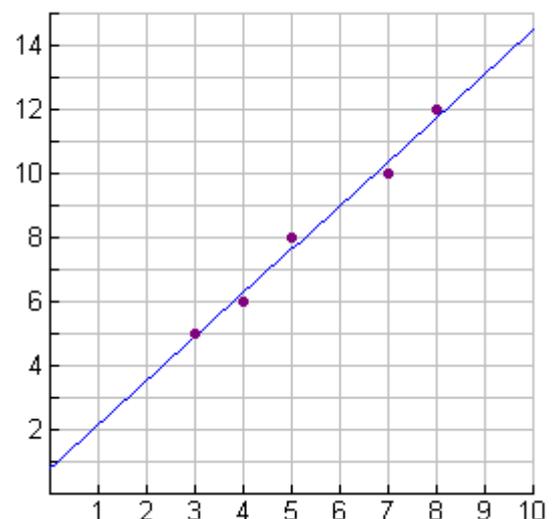
$$r = \text{coV}(xy) / \sqrt{V(x) \cdot V(y)} = \frac{4,72}{\sqrt{3,44 \cdot 6,56}} = 0,9936$$

und das Quadrat der Korrelation gibt die **Güte der Regressionsgeraden** an: $r^2 = 0,897 = 89,7\%$.

Das ist gut, die Güte ist erst unter 50% so schlecht, dass man keine sichere Aussage mehr machen kann, ob die

Daten korreliert sind oder nicht. Die Gerade liefert gute Prognosen!

Die Prognose für 10 Störche: $y(10) = 1,37 \cdot 10 + 0,79 = \mathbf{14,5}$ Geburten



Im Taschenrechner TI-30 II ist das einfach + ohne Formeln einzutippen → siehe Anhang!

Was bedeuten nun Korrelation und Bestimmtheitsmaß der Güte?

Die Korrelation gibt an, wie gut die Daten zusammenhängen.

Positive Korrelationen heißen: Je mehr x desto mehr y → Zeichnung A,C

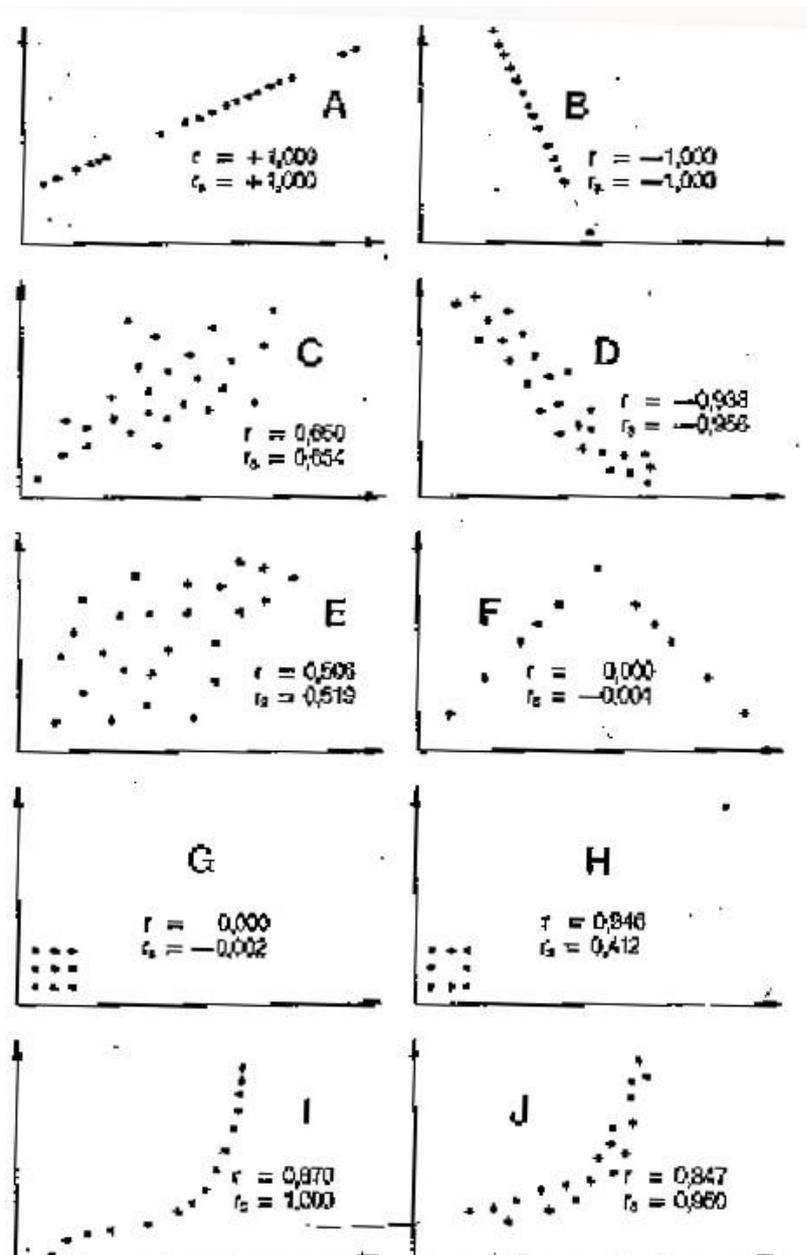
Negative Korrelationen heißen: Je mehr x desto weniger y → Zeichnung B,D

Korrelationen im Bereich zwischen $-0,4$ und $0,4$ heißen, dass die Daten nur wenig korrelieren.

→ Zeichnung E,F,G

Es könnten auch ganz andere Trägerfunktionen sinnvoll sein, wie Beispiel I und J zeigen. Da wäre eine quadratische Funktion weitaus besser angepasst.

Das **Bestimmtheitsmaß** gibt an, wie viel Prozent der Daten durch die Regressionsgerade gut dargestellt werden. Bei $r = 0,5$ ergibt sich $r^2 = 0,25$, also werden nur 25 % der Daten durch die Regressionsgerade richtig dargestellt.



$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \cdot \bar{y}$$

Die **Kovarianz**:

("Mittelwert der Produkte minus Produkt der Mittelwerte")

gibt die Größe des linearen Zusammenhangs zweier Größe x und y an.

Bei der Methode der *linearen Regression* nimmt man an, dass zwischen den beiden Werten ein linearer Zusammenhang besteht, das heißt:

$$y = k \cdot x + d + \text{ein zufälliger Fehler}$$

Die Konstanten k und d werden so bestimmt, dass die Summe der Quadrate der Fehler möglichst klein wird (Methode der kleinsten Fehlerquadrate von C.F. Gauß). Anschaulich können wir uns das so vorstellen, dass wir x und y als Koordinaten von Punkten auffassen und in ein Koordinatensystem einzeichnen. Wir suchen dann die Gerade, die diese Punktwolke am besten annähert. Diese Aufgabe kann man mit Hilfe der Differentialrechnung lösen und erhält als Gleichung der Regressionsgeraden:

$y = kx + d$, wobei

$$a = \frac{\text{Cov}(x, y)}{V(x)}$$

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

Die zweite Formel ergibt sich daraus, dass die Regressionsgerade durch den "Schwerpunkt" (\bar{x} / \bar{y}) der Punktwolke geht. Der *Korrelationskoeffizient* r liefert ein Maß dafür, wie gut die gegebenen Werte durch diese lineare Funktion angenähert werden. Er ist definiert durch

$$r = \frac{\text{Cov}(x, y)}{\sqrt{V(x) \cdot V(y)}}$$

Beispiel 2

soll zeigen, wie es mit negativen Korrelationen aussieht: Der Zusammenhang von (fiktiven) Mathe- und Deutschnoten soll gefunden werden

x=Mathenote	y=Englischnote	x*y	x ²	y ²	Ausgleichs- Näherung	Differenz zur Englischnote
1	5	5	1	25	4,9	-0,1
4	1	4	16	1	2,5	1,5
3	4	12	9	16	3,3	-0,7
3	5	15	9	25	3,3	-1,7
2	4	8	4	16	4,1	0,1
3	3	9	9	9	3,3	0,3
5	1	5	25	1	1,7	0,7
5	4	20	25	16	1,7	-2,3
4	1	4	16	1	2,5	1,5
3	3	9	9	9	3,3	0,3

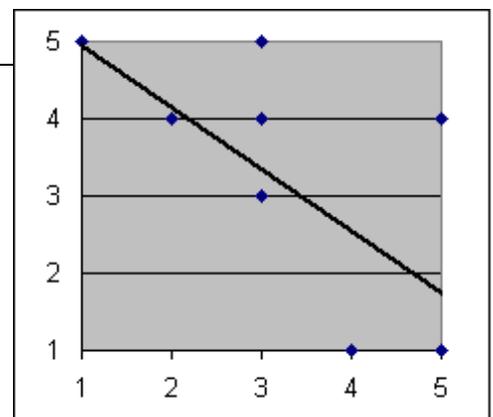
Summe	33	31	91	123	119
Mittelwerte (:10)	3,3	3,1	9,1	12,3	11,9

Varianz von x	$12,3 - 3,3^2 =$	1,41
Varianz von y	$11,9 - 3,1^2 =$	2,29
Kovarianz von xy	$9,1 - 3,3 \cdot 3,1 =$	-1,13

$k = \text{Cov}(xy) / V(x)$	$-1,13 : 1,41 =$	-0,80
$d = \bar{y} - k \cdot \bar{x}$	$3,1 - (-0,8) \cdot 3,3 =$	5,74

$r = -1,13 : \sqrt{1,41 \cdot 2,29}$	-0,63
$r^2 = (-0,63)^2$	0,40

Ausgleichsgerade: **$y = -0,80 \cdot x + 5,74$**



Wie man sieht ist die Korrelation negativ, je besser die Note in Mathematik ist, desto schlechter ist sie in Englisch. Diese Aussage trifft auf 40% der Daten zu. Das ist nicht sehr aussagekräftig, wie man an der Zeichnung sieht – die Punkte sind recht weit weg von der Geraden!

Übungen:

<p>1) In einem Unternehmen werden für verschiedene Produktionsmengen die Kosten festgestellt:</p> <table border="1"> <thead> <tr> <th>Stückmenge</th> <th>Kosten</th> </tr> </thead> <tbody> <tr><td>30</td><td>128</td></tr> <tr><td>35</td><td>165</td></tr> <tr><td>40</td><td>175</td></tr> <tr><td>55</td><td>240</td></tr> <tr><td>65</td><td>300</td></tr> </tbody> </table> <p>Ermitteln Sie die Regressionsgerade und die Korrelation. Wie groß ist die Differenz zwischen den realen Kosten bei 35 Stück und den berechneten? Bei wie viel Stück sind die Kosten auf 400 gestiegen?</p>	Stückmenge	Kosten	30	128	35	165	40	175	55	240	65	300	<p>2) Die durchschnittlichen Rohölpreise am Spotmarkt Amsterdam in den letzten 8 Monaten waren (in \$/Barrel)</p> <table border="1"> <thead> <tr> <th>Monat</th> <th>Preis</th> </tr> </thead> <tbody> <tr><td>3</td><td>28,4</td></tr> <tr><td>4</td><td>29,0</td></tr> <tr><td>5</td><td>30,3</td></tr> <tr><td>6</td><td>30,6</td></tr> <tr><td>7</td><td>29,8</td></tr> </tbody> </table> <p>Bestimmen Sie die Regressionsgerade und prognostizieren Sie den Ölpreis für die nächsten 2 Monate. Wie viel Prozent der Daten erklärt die Regressionsgerade? In welchem Monat übersteigt der Preis 33 \$?</p>	Monat	Preis	3	28,4	4	29,0	5	30,3	6	30,6	7	29,8	<p>3) Die folgenden Daten sollen den Zusammenhang zwischen dem Gewicht (Masse) und der Körpergröße belegen:</p> <table border="1"> <thead> <tr> <th>Körpergröße</th> <th>Masse</th> </tr> </thead> <tbody> <tr><td>164</td><td>48</td></tr> <tr><td>169</td><td>68</td></tr> <tr><td>160</td><td>51</td></tr> <tr><td>171</td><td>54</td></tr> <tr><td>165</td><td>53</td></tr> <tr><td>165</td><td>66</td></tr> </tbody> </table> <p>Bestimmen Sie die Regressionsgerade und die Korrelation. Wie gut ist die Aussage: „Je größer desto mehr Gewicht“ dokumentiert? Welchen Wert gibt die Regressionsgerade für $x = 165$ cm? Welche Körpergröße sollte ein 100 kg-Mensch haben?</p>	Körpergröße	Masse	164	48	169	68	160	51	171	54	165	53	165	66																																														
Stückmenge	Kosten																																																																																					
30	128																																																																																					
35	165																																																																																					
40	175																																																																																					
55	240																																																																																					
65	300																																																																																					
Monat	Preis																																																																																					
3	28,4																																																																																					
4	29,0																																																																																					
5	30,3																																																																																					
6	30,6																																																																																					
7	29,8																																																																																					
Körpergröße	Masse																																																																																					
164	48																																																																																					
169	68																																																																																					
160	51																																																																																					
171	54																																																																																					
165	53																																																																																					
165	66																																																																																					
<p>4) Das Alter von Ehepaaren ist gegeben:</p> <table border="1"> <thead> <tr> <th>Ehemann</th> <th>Ehefrau</th> </tr> </thead> <tbody> <tr><td>24</td><td>32</td></tr> <tr><td>18</td><td>20</td></tr> <tr><td>37</td><td>35</td></tr> <tr><td>26</td><td>30</td></tr> <tr><td>19</td><td>22</td></tr> <tr><td>21</td><td>26</td></tr> <tr><td>20</td><td>24</td></tr> <tr><td>31</td><td>38</td></tr> </tbody> </table> <p>Kann man auf Grund der Regressionsgeraden von einem Zusammenhang: „Je älter der Mann, desto älter die Frau“ reden? Zeichnen Sie die Punkte und die Gerade in eine Zeichnung. Wie groß sind die Mittelwerte und die Standardabweichungen des Alters von Mann und Frau?</p>	Ehemann	Ehefrau	24	32	18	20	37	35	26	30	19	22	21	26	20	24	31	38	<p>5) Folgende Tabelle zeigt die Weizenproduktion im Jahre 1982</p> <table border="1"> <thead> <tr> <th>Land</th> <th>Fläche in Mio. ha</th> <th>Produktion in Mio. t</th> </tr> </thead> <tbody> <tr><td>Sowjetunion</td><td>57</td><td>87</td></tr> <tr><td>USA</td><td>32</td><td>76</td></tr> <tr><td>China</td><td>28</td><td>63</td></tr> <tr><td>Indien</td><td>22</td><td>38</td></tr> <tr><td>Kanada</td><td>13</td><td>28</td></tr> <tr><td>Australien</td><td>12</td><td>9</td></tr> <tr><td>Türkei</td><td>9</td><td>18</td></tr> <tr><td>Frankreich</td><td>5</td><td>25</td></tr> <tr><td>England</td><td>2</td><td>10</td></tr> <tr><td>Deutschland</td><td>2</td><td>9</td></tr> </tbody> </table> <p>Zeichnen Sie ein Streudiagramm und die Regressionsgerade in eine Zeichnung. Welche Länder weichen am weitesten von der Geraden ab? Warum? Welchen Wert für die Produktionsmenge würde man für ein Land mit 10 Mio. ha Anbaufläche erwarten?</p>	Land	Fläche in Mio. ha	Produktion in Mio. t	Sowjetunion	57	87	USA	32	76	China	28	63	Indien	22	38	Kanada	13	28	Australien	12	9	Türkei	9	18	Frankreich	5	25	England	2	10	Deutschland	2	9	<p>6) Folgende Messwerte von Meereshöhe und Jahresmitteltemperatur sind gegeben:</p> <table border="1"> <thead> <tr> <th>Station</th> <th>Seehöhe</th> <th>Temp.</th> </tr> </thead> <tbody> <tr><td>Berlin</td><td>49</td><td>9,1</td></tr> <tr><td>Prag</td><td>374</td><td>7,9</td></tr> <tr><td>Wien</td><td>203</td><td>9,1</td></tr> <tr><td>Feuerkogel</td><td>1592</td><td>3,3</td></tr> <tr><td>Salzburg</td><td>437</td><td>8,6</td></tr> <tr><td>Budapest</td><td>130</td><td>10,9</td></tr> <tr><td>Zugspitze</td><td>2962</td><td>-5,0</td></tr> <tr><td>Innsbruck</td><td>579</td><td>8,4</td></tr> <tr><td>Säntis</td><td>2496</td><td>-2,3</td></tr> <tr><td>Sonnblick</td><td>3106</td><td>-6,4</td></tr> </tbody> </table> <p>Ermitteln Sie die Regressionsgerade und das Streudiagramm. Wie groß wird demnach die Temperatur in 800 m Meereshöhe sein? Welche Meereshöhe sollte man bevorzugen, wenn man ein Jahresmittel von 15° haben will?</p>	Station	Seehöhe	Temp.	Berlin	49	9,1	Prag	374	7,9	Wien	203	9,1	Feuerkogel	1592	3,3	Salzburg	437	8,6	Budapest	130	10,9	Zugspitze	2962	-5,0	Innsbruck	579	8,4	Säntis	2496	-2,3	Sonnblick	3106	-6,4
Ehemann	Ehefrau																																																																																					
24	32																																																																																					
18	20																																																																																					
37	35																																																																																					
26	30																																																																																					
19	22																																																																																					
21	26																																																																																					
20	24																																																																																					
31	38																																																																																					
Land	Fläche in Mio. ha	Produktion in Mio. t																																																																																				
Sowjetunion	57	87																																																																																				
USA	32	76																																																																																				
China	28	63																																																																																				
Indien	22	38																																																																																				
Kanada	13	28																																																																																				
Australien	12	9																																																																																				
Türkei	9	18																																																																																				
Frankreich	5	25																																																																																				
England	2	10																																																																																				
Deutschland	2	9																																																																																				
Station	Seehöhe	Temp.																																																																																				
Berlin	49	9,1																																																																																				
Prag	374	7,9																																																																																				
Wien	203	9,1																																																																																				
Feuerkogel	1592	3,3																																																																																				
Salzburg	437	8,6																																																																																				
Budapest	130	10,9																																																																																				
Zugspitze	2962	-5,0																																																																																				
Innsbruck	579	8,4																																																																																				
Säntis	2496	-2,3																																																																																				
Sonnblick	3106	-6,4																																																																																				

7) Eine Datenliste des **PKW-Bestandes** von Österreich ist gegeben:

Jahr	1950	1955	1960	1965	1970
PKW in 1000 Stk.	51	143	404	791	1197

Machen Sie eine Prognose für 1975-1990 und vergleichen Sie mit den realen Daten:

Jahr	1975	1980	1985	1990
PKW in 1000 Stk.	1721	2247	2531	2991

8) Durch Messung des Wasserabflaufs im **Kaffeefilter** einer Kaffeemaschine mit einer Filtertiefe von 10cm ergibt sich folgende Tabelle:

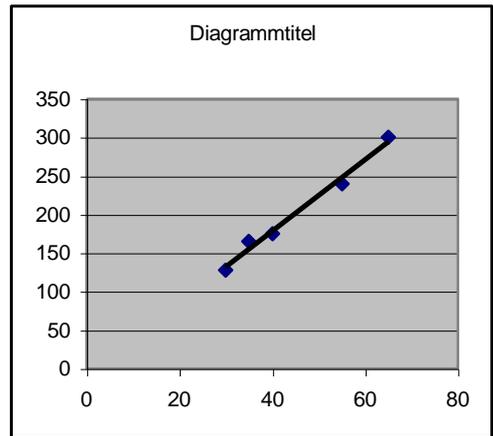
Zeit in sec	Wasserstand in cm
0	10
10	8
20	6
33	4
50	2

9) Der **Stromverbrauch** in Österreich in den Jahren 1975-1980 ist in der folgenden Tabelle gegeben. Legen Sie eine Regressionsgerade durch diese Daten und ermitteln Sie die Verbrauchswerte für 1981-1983

Jahr	Stromverbrauch in Md. kWh
1975	30,7
1976	33,1
1977	33,7
1978	35,3
1979	36,8
1980	38

Lösungen:

Stückmenge	Kosten	x*y	x ²	y ²	
30	128	3840	900	16384	
35	165	5775	1225	27225	
40	175	7000	1600	30625	
55	240	13200	3025	57600	
65	300	19500	4225	90000	
Summe	225	1008	49315	10975	221834
durch 5	45	201,6	9863	2195	44367
	= xq	= yq			
quadriert und verschoben			9072	2025	40643
Differenz			791	170	3724,2
Wurzel			= sxy	13,04	61,03
				= sx	= sy



Geradensteigung $k = sxy : sx^2 = 4,65$
 Geradenabstand $d = yq - k \cdot xq = -7,78$
 Korrelation $r = sxy : (sx \cdot sy) = 0,99$

1) Bestimmtheitsmaß $r^2 = 0,99$

Die Regressionsgerade ist $y = 4,65 \cdot x - 7,78$, für $x = 35$ ergibt das $y = 4,65 \cdot 35 - 7,78 = 155$

Die Differenz zu den realen Kosten ist: 10. Bei 88 Stück sind die Kosten ca. 400 GE

2) $y = 0,44x + 27,42$ $r = 0,76$ $r^2 = 0,58$ (58% = „Genügend“) Ölpreise: (8 → 30,94\$) (9 → 31,38\$)

Im 12. Monat übersteigt der Preis die 32\$-Marke

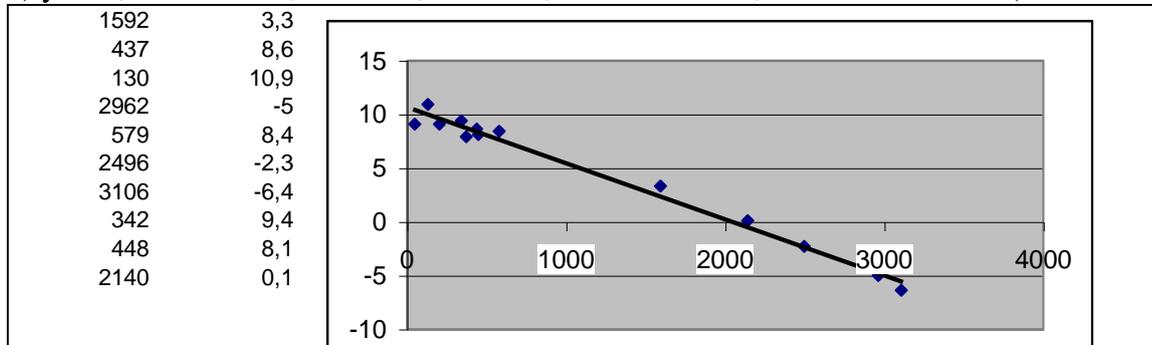
3) $y = 0,88x - 89,2$ $r = 0,41$ $r^2 = 0,17$ (wenig!) $x = 165 \rightarrow y = 56$ kg; $y = 100$ kg → 215 cm!!

4) $y = 0,86x + 7,28$ $r = 0,88$ $r^2 = 0,79$ Zusammenhang ist eher verkehrt quadratisch!

$x = 24,5 \pm 6,2$ $y = 28,3 \pm 5,9$ (Mittelwert ± Standardabweichung)

5) $y = 1,57x + 7,6$ $r = 0,93$ $r^2 = 0,87$ Australien unten (Wüste?) USA oben (no na!) 23,3 Mio.t

6) $y = -0,00525x + 10,63$ $r = -0,99$ $r^2 = 0,985$ 800m → 6,4° 15° → -832 m (unter dem Meeresspiegel)



7) $y = 58,8x - 3010,8$ $r = 0,98$ $r^2 = 0,95$

Jahr	1975	1980	1985	1990
PKW in 1000 Stk. real	1721	2247	2531	2991
Prognose mit Gerade	1399	1693	1987	2281

8) $y = -0,16x + 9,62$ $r = -0,99$ $r^2 = 0,98$ Wasserstand 0 nach 60 sec

9) $y = 1,4x + 31$ Verbrauchswerte 1981: 39,5, für 1982: 40,9 für 1983: 42,3 (real war es weniger!)

Anhang:

Berechnung der Korrelation und Regression mit Taschenrechner TI-30 II:



Wir wollen folgende Daten eingeben:

x	2	3	4
y	1	2	6

Zum Start wollen wir eine etwaige vorige Berechnung löschen: **2nd** **EXIT STAT**
Dann starten wir die Statistik-Eingabe mit **2nd** **STAT** und geben **2-VAR** **=**
für die Statistik mit 2 Variablen ein.

Nun kommen die Daten. Dazu starten wir mit **DATA**
Der Reihe nach geben wir nun die Daten ein und wechseln die Zeile mit der Pfeil-nach-unten-Taste **↓**

$X_1 = 2$ **↓** $Y_1 = 1$ **↓** $X_2 = 3$ **↓** $Y_2 = 2$ **↓** $X_3 = 4$ **↓** $Y_3 = 6$

Dann können wir schon die Ergebnisse sichten :

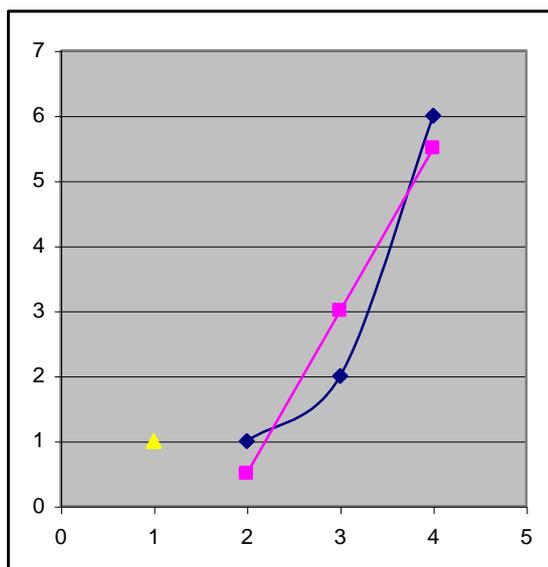
Mit **STATVAR** und der Pfeiltaste-nach-rechts **⇒** kommt man auf folgende Ergebnisse:

$n = 3$ $\bar{x} = 3$ $S_x = 1$ $\sigma_x = 0,81$ $\bar{y} = 3$ $S_y = 2,64$ $\sigma_y = 2,16$

und ganz hinten kommt unsere **Regression** : $a = 2,5$ $b = -4,5$ $r = 0,94$
und mit der **x²**-Taste bekommt man auch $r^2 = 0,89$ (89%-Güte)

Das ergibt die Ausgleichsgerade = Regressionlinie = **Trendgerade** $y = 2,5x - 4,5$

x	2	3	4
y	1	2	6
y-Trend: 2,5x-4,5	0,5	3	5,5



Zum Schluss wollen wir noch alles löschen: **2nd** **EXIT STAT**