

Lineare Regression - von der Black-Box zur White-Box

Anwendungen von Modulen werden unterrichtlich oft nur dann akzeptiert, wenn die zugrunde liegenden Algorithmen verstanden sind, und die zugehörigen elementaren Operationen (jederzeit) ohne Hilfsmittel realisiert werden können. Ein in diesem Sinne verständiger Umgang mit den verfügbaren Regressionen ist aber schulisch kaum zu vermitteln. Dennoch sollte auf Regressionen als „Black-Boxes“ nicht verzichtet werden. Sie gewähren Schülerinnen und Schülern Einsichten in funktionale Zusammenhänge zwischen Größenpaaren, die auf herkömmliche Weise nur mit erhöhtem Zeit- und Arbeitsaufwand zu bewirken sind. So führen Weg-Zeit-Experimente, Abkühlungsversuche, Optimierungsaufgaben, Wachstumsuntersuchungen, Zufallsversuche u.a. zu endlich vielen Messpunkten $P(x_1; y_1), P(x_2; y_2), \dots, P(x_n; y_n)$. Sie werden im Data-Menü des TI-89/92 Plus/Voyage 200 tabellarisch erfasst und durch Regressionen ausgewertet, ohne dass das Verfahren, wohl aber die Wahl der speziellen Regression kritisch zu hinterfragen ist.

Eine lineare Regression

Eine Messwertreihe mit 16 Messpunkten:

	F1	F2	F3	F4	F5	F6	F7
	Plot	Setup	Cell	Header	Calc	Util	Stat
DATA	1	x1	y1				
	c1	c2	c3	c4	c5	c6	c7
1	1.00	19	12				
2	2.00	17	12				
3	3.00	15	11				
4	4.00	18	14				
5	5.00	16	14				
6	6.00	21	16				
7	7.00	24	15				
8	8.00	13	12				
9	9.00	12	10				
10	10.00	22	17				
11	11.00	11	12				
12	12.00	20	14				
13	13.00	22	14				
14	14.00	9	9				
15	15.00	10	11				
16	16.00	11	9				
17							
r11c1=11.							
MAIN	DEG EXACT						

Abb. 1

(Die Betitelung der Listen c1, c2, ... mit den Namen x_i, y_i, \dots soll die Übersichtlichkeit verbessern.)

Graphische und arithmetische Auswertung der 16 Messpunkte $P(x_1; y_1), P(x_2; y_2), \dots, P(x_{16}; y_{16})$ mit den Möglichkeiten, die der Data/Matrix Editor bietet: Die Lineare Regression wird im Menue $\mathbf{\frac{1}{2}}$ Calc aktiviert. c2 und c3 sind Listennamen. Die Regressionsgleichung wird unter $y_1(x)$ gespeichert.

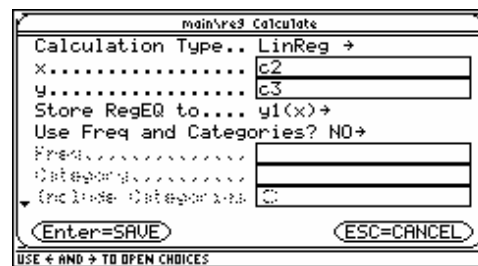


Abb. 2

Die Auswertung liefert:



Abb. 3

Die graphische Darstellung der Punktwolke mit der zugehörigen linearen Funktion.

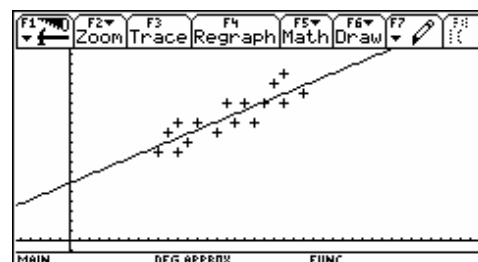


Abb. 4

Einordnung

Die Ausleuchtung der durchgeführten Regression ist Anlass, lineare Regressionsverfahren als innermathematische Modellierung aufzubereiten.

Unterrichtliche Umsetzung

Die Suche nach einer sinnvollen Ausgleichsgeraden führt unterrichtlich nach Schülerintentionen meist zur Berechnung der zugehörigen Gleichung aus zwei „mittleren Punkten“ bzw. den Start- und Endwerten. Auch die Idee der Minimierung der Abstände wird gelegentlich im Zusammenhang mit der graphischen Darstellung des Problems genannt. Die auf Legendre und Gauß zurückgehende Methode, aus gegebenen Messwerten Werte zu berechnen, für welche die

Summe der Fehlerquadrate $\sum_{i=1}^n (f(x_i) - y_i)^2$ minimal

ist, muss als Modellierungsansatz vorgegeben und von anderen Möglichkeiten abgegrenzt werden. Im speziellen Fall ist zur Beschreibung der Ausgleichsgeraden eine Gleichung $f(x) = mx + b$ zu

finden, die der Setzung genügt, dass

$$\sum_{i=1}^{16} (m \cdot x_i + b - y_i)^2 \text{ minimal ist.}$$

Alle Bemühungen, den Summenterm zur Bestimmung von m und b brauchbar umzuwandeln, führen zu komplexen Ausdrücken, die eine elementare Lösung nicht zulassen.

Mit $x_q = \frac{1}{16} \sum_{i=1}^{16} x_i$ und $y_q = \frac{1}{16} \sum_{i=1}^{16} y_i$ wird stattdessen

der sogenannte Regressionsschwerpunkt $P(x_q; y_q)$ als ein Punkt der Ausgleichsgeraden festgelegt.

Zur Bestimmung der Geradensteigung wird eine Koordinatentransformation so durchgeführt, dass der Regressionsschwerpunkt $P(x_q; y_q)$ Ursprung des neuen Koordinatensystems ist. Die Richtung der Koordinatenachsen und der Maßstab bleiben unverändert.

Es gelten die Transformationsgleichungen:

$$x_i - x_q = x_k \text{ und } y_i - y_q = y_k.$$

Mit $c4 = \text{cumSum}(c2)$ und $c5 = \text{cumSum}(c3)$ lassen sich $x_q = 260/16$ und $y_q = 202/16$ einfach bestimmen und damit die transformierten Koordinaten x_k und y_k einfach in $c6$ und $c7$ berechnen.

	F1	F2	F3	F4	F5	F6	F7
	Plot	Setup	Cell	Header	Calc	Util	Stat
DATA	i	x _i	y _i	Sx _i	Sy _i	x _k	y _k
	c1	c2	c3	c4	c5	c6	c7
1	1.00	19	12	19.0	12.0	2.75	-6.3
2	2.00	17	12	36.0	24.0	1.750	-6.3
3	3.00	15	11	51.0	35.0	-1.3	-1.6
4	4.00	18	14	69.0	49.0	1.75	1.38
5	5.00	16	14	85.0	63.0	-2.25	1.38
6	6.00	21	16	106.79	80.475	3.38	
7	7.00	24	15	130.94	95.075	2.38	
8	8.00	13	12	143.106	107.3	-6.3	
9	9.00	12	10	155.116	117.3	-2.6	
10	10.0	22	17	177.133	134.38		
11	11.0	11	12	188.145	146.3	-6.3	
12	12.0	20	14	208.159	160.38		
13	13.0	22	14	230.173	174.38		
14	14.0	9	9	239.182	183.36		
15	15.0	10	11	249.193	194.38		
16	16.0	11	9	260.202	203.36		
17							

Abb. 5

Bei einer Koordinatentransformation werden alle Punkte der Ebene in zwei verschiedenen Koordinatensystemen dargestellt, $c6$ und $c7$ sind die Listen der Koordinaten im transformierten System. Die zum neuen Listenpaar ($c6, c7$) zugehörige Ausgleichsgerade verläuft durch den Ursprung und genügt der Gleichung $g(x) = mx$. Mit dem TI-92/Voyage 200 lassen sich die zum Listenpaar ($c6, c7$) zugehörigen Punkte nun im selben Koordinatensystem darstellen. Es handelt sich um eine Abbildung der Ebene (Punkte).

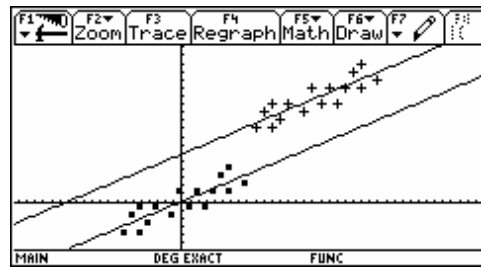


Abb. 6

Somit ist die Summe der Fehlerquadrate

$$\sum_{k=1}^{16} (g(x_k) - y_k)^2 = \sum_{k=1}^{16} (mx_k - y_k)^2 \text{ zu minimieren.}$$

Diese Aufgabe ist einfacher, denn

$$\sum_{k=1}^{16} (mx_k - y_k)^2 = m^2 \cdot \sum_{k=1}^{16} x_k^2 - 2m \cdot \sum_{k=1}^{16} x_k \cdot y_k + \sum_{k=1}^{16} y_k^2$$

ist ein quadratischer Term in m und sein Graph entspricht einer nach oben geöffneten Parabel. Die erste Koordinate des Scheitelpunktes bestimmt also die Steigung m so, dass diese Summe minimal ist.

Damit ergibt sich der Regressionskoeffizient der

$$\text{Linearen Regression als } m_{\min} = \frac{\sum_{k=1}^{16} x_k \cdot y_k}{\sum_{k=1}^{16} x_k^2}.$$

	F1	F2	F3	F4	F5	F6	F7
	Plot	Setup	Cell	Header	Calc	Util	Stat
DATA	x _k	y _k	x _k y _k	x _k ²	Sx _k	Sy _k	Sx _k ²
	c6	c7	c8	c9	c10	c11	
1	2.750	-6.25	-1.72	7.563	-1.72	7.563	
2	.7500	-.625	-.469	.5625	-2.19	8.125	
3	-1.25	-1.63	2.031	1.563	-1.56	9.688	
4	1.750	1.375	2.406	3.063	2.250	12.75	
5	-.250	1.375	-.344	.0625	1.906	12.81	
6	4.750	3.375	16.03	22.56	17.94	35.38	
7	7.750	2.375	18.41	60.06	36.34	95.44	
8	-3.25	-.625	2.031	10.56	38.38	106.0	
9	-4.25	-2.63	11.16	18.06	49.53	124.1	
10	5.750	4.375	25.16	33.06	74.69	157.1	
11	-5.25	-.625	3.281	27.56	77.97	184.7	
12	3.750	1.375	5.156	14.06	83.13	198.8	
13	5.750	1.375	7.906	33.06	91.03	231.8	
14	-7.25	-3.63	26.28	52.56	117.3	284.4	
15	-6.25	-1.63	10.16	39.06	127.5	323.4	
16	-5.25	-3.63	19.03	27.56	146.5	351.0	
17							

Abb. 7

Im Data/Matrix Editor wird in Spalte $c8$ die Folge der $x_k y_k$ und in Spalte $c9$ die der x_k^2 eingetragen. Die entsprechenden kumulativen Summen stehen in Spalten $c10$ und $c11$. Damit berechnet sich der Regressionskoeffizient als $m = 146,5/351 = 0,417379$. Zur Bestimmung des fehlenden Achsenabschnittes wird entsprechend der Transformation die Gleichung $f(x) = m(x - x_q) + y_q$ herangezogen. Dies ergibt $f(x) = 0,417379(x - 16,25) + 12,625 = 0,417379x + 5,84259$, in Übereinstimmung mit der vom TI-92/Voyage 200 ermittelten Regressionsgleichung (vgl. Abb. 3).

Die Güte der Angleichung

Wie aber sind $corr$ und R^2 zu interpretieren? Je besser unterrichtlich ermittelte Regressionsgraphen die Messpunkte treffen, desto mehr nähern sich $corr$

bzw. R^2 dem Wert 1. Sie geben offensichtlich Auskunft über die Güte der Angleichung (vgl. Abb. 3). Hier ist zu thematisieren, dass die Wurzel aus dem Produkt der Steigungen der Ausgleichsgeraden zu den jeweiligen Listenpaaren (x;y) bzw. (y;x) als Maß für die Abhängigkeit der Größen x und y vereinbart wird und Korrelationskoeffizient heißt. Es gilt:

$$R^2 = \frac{\sum x_k y_k \cdot \sum x_k y_k}{\sum x_k^2 \cdot \sum y_k^2}$$

Im Falle strenger Linearität genügen die Funktionen zu den Listenpaaren (x;y) und (y;x) der Relation „ist Umkehrfunktion von“, so dass die Wurzel den Wert 1 liefert.

Die Funktionen der Regressionsgeraden zu den Listenpaaren (c2;c3) und (c3;c2) der Ausgangsmesswerte (Abb. 1) genügen nicht der Relation „ist Umkehrfunktion von“, d.h. das Produkt ihrer Steigungskoeffizienten ist ungleich 1. Darstellung der Ausgleichsgeraden zu den Listenpaaren (c2;c3) und (c3;c2):

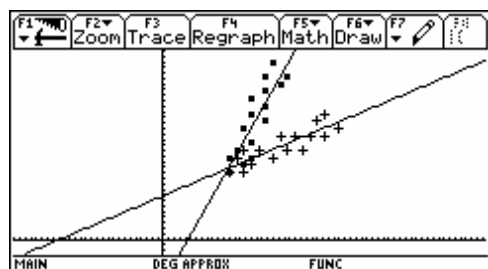


Abb. 8

In der Tabelle wird noch in Spalte c12 die Folge der y_k^2 und in Spalte c13 ihre kumulative Summe berechnet.

	F1	F2	F3	F4	F5	F6	F7
	Plot	Setup	Cell	Header	Calc	Util	Stat
DATA	xkyk	xkxk	Sxkyk	Sxkxk	ykuk	Sykuk	
	c8	c9	c10	c11	c12	c13	
1	-1.72	7.563	-1.72	7.563	.3906	.3906	
2	-.469	.5625	-2.19	8.125	.3906	.7813	
3	2.031	1.563	-1.156	9.688	2.641	3.422	
4	2.406	3.063	2.250	12.75	1.891	5.313	
5	-.344	.0625	1.906	12.81	1.891	7.203	
6	16.03	22.56	17.94	35.38	11.39	18.59	
7	18.41	60.06	36.34	95.44	5.641	24.23	
8	2.031	10.56	38.38	106.0	.3906	24.63	
9	11.16	18.06	49.53	124.1	16.891	31.52	
10	25.16	33.06	74.69	157.1	19.14	50.66	
11	3.281	27.56	77.97	184.7	.3906	51.05	
12	5.156	14.06	83.13	198.8	1.891	52.94	
13	7.906	33.06	91.03	231.8	1.891	54.83	
14	26.28	52.56	117.3	284.4	13.14	67.97	
15	10.16	39.06	127.5	323.4	2.641	70.61	
16	19.03	27.56	146.5	351.0	13.14	83.75	
17							
	Rr11c8=3.28125						
	MAIN	DEG APPROX	PAR				

Abb. 9

Im Hauptbildschirm erhält man dann die in Abb. 3 berechneten Korrelationskoeffizienten.

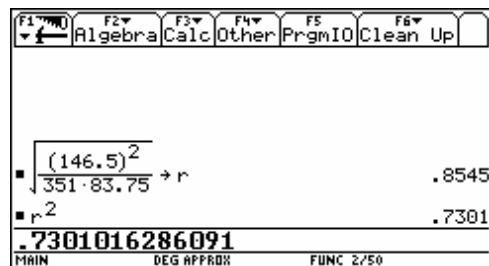


Abb. 10

Anregungen und Anwendungsbeispiele unter www.studienseminar-leer.de
Fachleiter – Mathematik.

Der Autor:
Reinhold Hoffmann
Glansdorfer Str.47
26817 Rhaderfehn
email: K.R.Hoffmann@t-online.de

Anmerkung der Redaktion:
Nach Angaben des Autors wurde diese Unterrichtssequenz in der 9. Jahrgangsstufe, also in der Sek I behandelt.
Für die Sek II kann man vom gleichen Datenmaterial ausgehen, aber die CAS-Fähigkeiten des Voyage 200 ausnutzen. Man speichert die Koordinaten der Messpunkte in 2 Listen xi (für die Abszisse) und yi (für die Ordinate). Dann deklariert man einen Term $m \cdot x + b$ als Funktion $g(x)$. Die Summe der quadratischen Abweichung $(g(x_i) - y_i)^2$ ist ein biquadratischer Ausdruck in m und b. Ihn weist man einer Funktion, z.B. d(m, b) zu. Setzt man nun die Ableitungen von d nach m bzw. nach b gleich Null, so erhält man ein System linearer Gleichungen in m und b. Die Lösung dieses Gleichungssystems liefert genau die vom CAS berechnete Regressionsgerade.

