

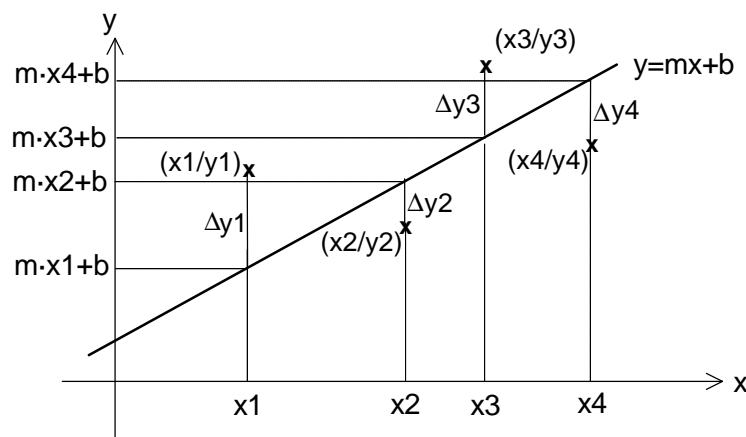
Wie funktioniert eigentlich eine lineare Regression?

Wenn wir zum Beispiel in der Physik den Zusammenhang von Strom und Spannung an einem Widerstand untersuchen wollen, messen wir bei verschiedenen Spannungen die Stromstärke. Wir erhalten eine Reihe von Messwertpaaren, die wir in ein Koordinatensystem eintragen. In diesem Fall können wir aus der Lage der Punkte vermuten, dass es einen linearen Zusammenhang zwischen Strom und Spannung gibt. Die Punkte werden aber wegen der Messfehler nie exakt auf einer Geraden liegen. Bei der linearen Regression versucht man eine Gerade zu finden, die "möglichst nahe" an allen Punkten liegt. Es handelt sich also um ein typisches Optimierungsproblem.

Zuerst müssen wir definieren, was wir unter "Nähe" verstehen wollen. Der Abstand ist die senkrechte Strecke von einem Punkt zur Geraden. Die Berechnung dieser Strecke bedeutet aber einen erheblichen Rechenaufwand. Wir beschränken uns auf den Abstand in y-Richtung, in der Zeichnung mit Δy_1 .. Δy_4 bezeichnet. Wir müssen die Gerade so bestimmen, dass die Summe aller Abstände minimal wird.

Da die Punkte mal über, mal unter der Geraden liegen, würde man bei der Berechnung durch eine einfache Differenz $\Delta y_1 = m \cdot x_1 + b - y_1$ mal positive und mal negative Vorzeichen bekommen. Bei der Summe könnten sich dann die Abstände aufheben.

Eine Möglichkeit besteht darin, Beträge zu verwenden. Dann ist der Abstand in y-Richtung $\Delta y_1 = |m \cdot x_1 + b - y_1|$. Mit Beträgen ist das Rechnen aber recht unbequem. Eine andere Möglichkeit ist, die Differenzen zu quadrieren, dann kann kein negatives Vorzeichen auftreten. Die Idee dahinter ist: wenn das Quadrat der Differenz minimal ist, dann ist auch der Betrag minimal, oder: je kleiner x^2 , desto kleiner $|x|$.



Mit diesen Voraussetzungen können wir uns daran machen, eine Regressionsgerade zu bestimmen. Dabei wollen wir folgende Werte verwenden:

x	1	2	3	4
y	5,2	7,4	12,3	15,9

Für die Summe der Abstandsquadrate gilt

$$(m \cdot x_1 + b - y_1)^2 + (m \cdot x_2 + b - y_2)^2 + (m \cdot x_3 + b - y_3)^2 + (m \cdot x_4 + b - y_4)^2$$

und mit den Werten aus der Tabelle erhalten wir:

$$(m \cdot 1 + b - 5,2)^2 + (m \cdot 2 + b - 7,4)^2 + (m \cdot 3 + b - 12,3)^2 + (m \cdot 4 + b - 15,9)^2 \\ = 4b^2 + 20b \cdot m - 81,6b + 30m^2 - 241m + 485,9$$

Aufgabe 1: Überprüfe das Ergebnis mit dem TI-92.

Aufgabe 2: Bestimme irgendwie (z.B. graphisch) eine möglichst gute Regressionsgerade und berechne dafür die Summe der Abstandsquadrate.

Dieser Term hängt von zwei Variablen ab. Wir können uns das als Funktionenschar vorstellen: $f_m(b) = 4b^2 + 20b \cdot m - 81,6b + 30m^2 - 241m + 485,9$

Aufgabe 3: Um welchen Kurventyp handelt es sich?

Aufgabe 4: Bestimme den Tiefpunkt der Kurve K_m .

Die y-Koordinate des Tiefpunktes gibt die Summe der Abstandsquadrate für das entsprechende m an.

Aufgabe 5: Für welches m hat K_m den niedrigsten Tiefpunkt?

Dieser Punkt ist der kleinste Funktionswert, der bei $f_m(b)$ überhaupt möglich ist. Das ist also das Minimum der Summe der Abstandsquadrate.

Aufgabe 6: Bestimme nun die Regressionsgerade.

Zum Vergleich wollen wir nun eine Regressionsgerade mit dem TI-92 berechnen.

(1) Wähle **O** - 6:Data/Matrix Editor - 3:New

(2) Gib in die Spalte c1 die x-Werte ein und in c2 die y-Werte ein.

(3) Tippe F5:Calc und wähle als "Calculation Type" 5:LinReg

(4) Gib im x-Fenster c1 und im y-Fenster c2 ein.

(5) Wähle bei "Store RegEQ to" $y_1(x)$, dann zweimal mit **↵** bestätigen.

Es erscheint ein Fenster mit den Werten für m und b (hier heißen sie a und b).

(6) Um das Ergebnis im Graphen zu sehen, gehe wieder in die Tabelle.

(7) Wähle F2:Plot Setup und dann F1:Define

(8) Wähle als "Plot Type" Scatter, als "Mark" Cross und für x und y wieder c1 und c2.

(9) Wähle im Window-Fenster einen vernünftigen Ausschnitt für das Koordinatensystem.

(10) Schalte nun in das Graph-Fenster um.