

## Multiple lineare Regression

In diesem Skript sollen die theoretischen Grundlagen zur multiplen linearen Regression sowie deren Umsetzung in MS-Excel dargestellt werden.

Für die Unterstützung bei der Erstellung dieses Skripts gilt Herrn Prof. Dr. W. Kristof, Ph.D., mein besonderer Dank.

### 1. Vorüberlegungen

Ein lineares Gleichungssystem  $\underline{A}\underline{x} = \underline{b}$  heißt homogen, wenn  $\underline{b} = \underline{0}$ . Ist  $\underline{b} \neq \underline{0}$ , dann heißt ein solches Gleichungssystem inhomogen.

Ein homogenes Gleichungssystem hat nur dann eine Lösung  $\underline{x} \neq \underline{0}$ , wenn  $\underline{A}$  nicht vollen Rang hat. Dies kann man sich verdeutlichen, wenn man sich die Definition der linearen Unabhängigkeit anschaut<sup>1</sup>.

Ein inhomogenes Gleichungssystem ist genau dann lösbar, wenn  $\underline{b}$  im Raum liegt, der durch die Spaltenvektoren von  $\underline{A}$  aufgespannt wird. Dann ist  $\text{Rang } \underline{A} = \text{Rang } (\underline{A}|\underline{b})$ .

Ist  $\underline{A}$  quadratisch und hat vollen Rang, dann hat  $\underline{A}\underline{x} = \underline{b}$  die eindeutige Lösung

$$(1) \quad \underline{x} = \underline{A}^{-1}\underline{b}.$$

Im Fall der multiplen Regression werden inhomogene Gleichungssysteme betrachtet, bei denen  $\underline{A}$  vom Format  $(n, m)$  ist,  $n \geq m$  und  $\text{Rang } \underline{A} = m$ . In der Regel liegt  $\underline{b}$  nicht im Raum, der von den Spaltenvektoren von  $\underline{X}$  aufgespannt wird. Ein solches System ist dann inkonsistent, denn es existiert keine exakte Lösung.

---

<sup>1</sup>  $m$  Vektoren  $\underline{a}_1, \dots, \underline{a}_m$  sind linear unabhängig, wenn  $q_1\underline{a}_1 + q_2\underline{a}_2 + \dots + q_m\underline{a}_m = \underline{0}$  nur durch  $q_j = 0$  für alle  $q_j$  erfüllt werden kann.

Wir gehen jedoch zunächst davon aus, daß ein solches Gleichungssystem  $\underline{A}\underline{x} = \underline{b}$  (mit  $\underline{A}$  vom Format  $(n, m)$ ,  $n \geq m$  und  $\text{Rang } \underline{A} = m$ ) eine exakte Lösung hat.

Multiplikation der Gleichung  $\underline{A}\underline{x} = \underline{b}$  von links mit  $\underline{A}'$  liefert

$$(2) \quad \underline{A}'\underline{A}\underline{x} = \underline{A}'\underline{b}.$$

$\underline{A}'\underline{A}$  ist quadratisch und vom Format  $(m, m)$ . Aus der linearen Algebra ist bekannt, daß  $\text{Rang } \underline{A} = \text{Rang } \underline{A}' = \text{Rang } \underline{A}'\underline{A}$ . Folglich ist  $\text{Rang } \underline{A}'\underline{A} = m$  und damit  $\underline{A}'\underline{A}$  invertierbar. Gleichung (2) von links mit  $(\underline{A}'\underline{A})^{-1}$  multipliziert liefert

$$(\underline{A}'\underline{A})^{-1}\underline{A}'\underline{A}\underline{x} = (\underline{A}'\underline{A})^{-1}\underline{A}'\underline{b}. \quad \text{Demnach ist}$$

$$(3) \quad \underline{x} = (\underline{A}'\underline{A})^{-1}\underline{A}'\underline{b}.$$

Dies ist die (exakte) Lösung für das Gleichungssystem, sofern eine Lösung existiert<sup>2</sup>.

Im folgenden soll die „Lösung“ eines inkonsistenten und inhomogenen Gleichungssystems gesucht werden. In diesem Fall liegt  $\underline{b}$  nicht im Raum, der in den Spaltenvektoren von  $\underline{A}$  aufgespannt wird. Gesucht ist die „beste Lösung im Sinn der kleinsten Quadrate“ – dies wird durch das Symbol „ $\doteq$ “ ausgedrückt - für die  $x_j$  in folgendem Gleichungssystem:

$$(4) \quad \begin{array}{rcl} a_{11}x_1 + \dots + a_{1m}x_m & \doteq & b_1 \\ \vdots & & \vdots \\ a_{n1}x_1 + \dots + a_{nm}x_m & \doteq & b_n \end{array}$$

Man verwendet

$$(5) \quad \underline{A} = \begin{pmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{pmatrix}, \underline{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}, \text{ mit } \text{Rang } \underline{A} = m \text{ und } n \geq m.$$

---

<sup>2</sup> Dies ist gleichzeitig die allgemeine Lösung für alle konsistenten Gleichungssysteme. Falls  $\underline{A}$  quadratisch, ist  $\underline{x} = (\underline{A}'\underline{A})^{-1}\underline{A}'\underline{b} = \underline{A}^{-1}\underline{A}^{-1}\underline{A}'\underline{b} = \underline{A}^{-1}\underline{b}$ , wie bereits in Gleichung (1) dargestellt.

Folglich ist das Gleichungssystem (4) in Matrixschreibweise

$$(6) \quad \underline{Ax} \doteq \underline{b}.$$

Ein inkonsistentes Gleichungssystem hat keine exakte Lösung, also fordert man für

$$(7) \quad \begin{array}{l} a_{11}x_1 + \dots + a_{1m}x_m - b_1 = d_1 \\ \vdots \\ a_{n1}x_1 + \dots + a_{nm}x_m - b_n = d_n \end{array},$$

daß  $d_1^2 + \dots + d_n^2$  durch die Wahl der  $x_j$  minimal werde. Man steht somit vor dem Problem, die Funktion aller  $x_j$ ,

$$(8) \quad f(x_1, \dots, x_m) = (a_{11}x_1 + \dots + a_{1m}x_m - b_1)^2 + \dots + (a_{n1}x_1 + \dots + a_{nm}x_m - b_n)^2$$

zu minimieren. Minimierung von  $f(x_1, \dots, x_m)$  bedeutet

$$\frac{\partial f(x_1, \dots, x_m)}{\partial x_j} = 0, \quad j = 1, \dots, m.$$

$$\frac{\partial f}{\partial x_j} = 2(a_{11}x_1 + \dots + a_{1m}x_m - b_1)a_{1j} + \dots + 2(a_{n1}x_1 + \dots + a_{nm}x_m - b_n)a_{nj} = 0$$

$$(9) \quad \Leftrightarrow (a_{11}a_{1j} + \dots + a_{n1}a_{nj})x_1 + \dots + (a_{1m}a_{1j} + \dots + a_{nm}a_{nj})x_m = b_1a_{1j} + \dots + b_na_{nj} = 0.$$

Dies ergibt,  $j = 1, \dots, m$ , eingesetzt, folgende  $m$  Gleichungen:

$$(10) \quad \begin{array}{l} (a_{11}a_{11} + \dots + a_{n1}a_{n1})x_1 + \dots + (a_{1m}a_{11} + \dots + a_{nm}a_{n1})x_m = b_1a_{11} + \dots \\ + b_na_{n1} \\ \vdots \\ (a_{11}a_{1m} + \dots + a_{n1}a_{nm})x_1 + \dots + (a_{1m}a_{1m} + \dots + a_{nm}a_{nm})x_m = b_1a_{1m} + \dots \\ + b_na_{nm} \end{array}$$

Dies ist

$$(11) \quad \underbrace{\begin{pmatrix} a_{11}^2 + \dots + a_{n1}^2 & \dots & a_{1m}a_{11} + \dots + a_{nm}a_{n1} \\ \vdots & & \vdots \\ a_{11}a_{1m} + \dots + a_{n1}a_{nm} & \dots & a_{1m}^2 + a_{nm}^2 \end{pmatrix}}_{\underline{A}'\underline{A}} \underbrace{\begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}}_{\underline{x}} = \underbrace{\begin{pmatrix} a_{11} & \dots & a_{n1} \\ \vdots & & \vdots \\ a_{1m} & \dots & a_{nm} \end{pmatrix}}_{\underline{A}'} \underbrace{\begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}}_{\underline{b}}, \text{ also}$$

$$(12) \quad \underline{A}'\underline{A}\underline{x} = \underline{A}'\underline{b}, \text{ daher}$$

$$(13) \quad \underline{x} = (\underline{A}'\underline{A})^{-1} \underline{A}'\underline{b}, \quad \text{mit Format } (\underline{A}'\underline{A}) = (m, m), \text{ Rang } (\underline{A}'\underline{A}) = m \\ \text{und } (\underline{A}'\underline{A})^{-1} \text{ existent nach Voraussetzung.}$$

Gleichung (13) ist also identisch mit der Gleichung (3). Dies ist ein allgemeines Ergebnis. Es ist die (exakte) Lösung für  $\underline{x}$ , sofern das Gleichungssystem lösbar ist, andernfalls die „beste Lösung“ für  $\underline{x}$  im Sinn der kleinsten Quadrate.

## 2. Multiple Regression

Es ist zu untersuchen, wie der Zusammenhang einer Variablen  $Y$  mit den Variablen  $X_1, \dots, X_m$  aussieht.

Die Variable  $Y$  wird in diesem Fall als Kriterium oder *Regressand* bezeichnet, die Variablen  $X_i$  als Prädiktoren oder *Regressoren*.

Kann ein solcher Zusammenhang beschrieben werden, so können folgende Aufgaben durchgeführt werden:

- Erkennen eines statistischen Zusammenhangs
- Schätzung der Parameter eines statistischen Zusammenhangs
- Trendanalysen
- Prognose von  $Y$  aus  $X_1, \dots, X_m$

Gegeben sind die Variablen  $Y$  (Regressand) sowie die  $X_1, \dots, X_m$  (Regressoren).

Der Regressand hat metrisches Skalenniveau, die Prädiktoren haben in der Regel ebenfalls metrisches Skalenniveau, können aber auch dichotome Variablen sein. Werden zusätzlich statistische Testverfahren angewendet, müssen bestimmte Normalverteilungsannahmen erfüllt sein, die jedoch von Fall zu Fall verschieden sind. Diese Verteilungsannahmen sind jedoch keine Voraussetzung für die Durchführung der multiplen Regression.

In vielen Fällen wird zu den Variablen eine Konstante  $b_0$  hinzugenommen. Die Aufgabenstellung ist analog zu (4) die Annäherung von  $Y$  im Sinn der kleinsten Quadrate durch Linearkombination  $b_0 + b_1X_1 + \dots + b_mX_m$ , also für tatsächliche Meßwerte

$$(14) \quad \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \stackrel{\circ}{=} \begin{pmatrix} b_0 + b_1X_{11} + \dots + b_mX_{1m} \\ \vdots \\ b_0 + b_1X_{n1} + \dots + b_mX_{nm} \end{pmatrix},$$

wir betrachten also  $n$  Objekte (Personen) und  $m$  Variablen (Regressoren). Die Bezeichnung der Variablen entspricht der allgemein üblichen Bezeichnung im Falle der multiplen Regression - im Gegensatz zur Variablenbezeichnung in den Vorüberlegungen sind also hier nicht die  $x_j$  unbekannt bzw. gesucht, sondern die  $b_j$ !

Gleichung (14) aufgelöst nach  $b_0$  ergibt

$$(15) \quad \begin{pmatrix} b_0 \\ \vdots \\ b_0 \end{pmatrix} \stackrel{\circ}{=} \begin{pmatrix} y_1 - b_1X_{11} - \dots - b_mX_{1m} \\ \vdots \\ y_n - b_1X_{n1} - \dots - b_mX_{nm} \end{pmatrix}.$$

Aus der Statistik 1 ist bekannt, daß die Zahl, die die kleinste Abweichung im Sinn der kleinsten Quadrate von einer Menge von Werten darstellt, das arithmetische Mittel dieser Werte ist. Also ist

$$(16) \quad b_0 = \bar{y} - b_1 \bar{x}_1 - \dots - b_m \bar{x}_m$$

Aus (14), (15) und (16) folgt, daß

$$(17) \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix} \stackrel{=}{=} \begin{pmatrix} b_1(x_{11} - \bar{x}_1) + \dots + b_m(x_{1m} - \bar{x}_m) \\ \vdots \\ b_1(x_{n1} - \bar{x}_1) + \dots + b_m(x_{nm} - \bar{x}_m) \end{pmatrix}.$$

Wendet man die Lösung (13) auf das Gleichungssystem (17) an, so erhält man – unter Beachtung der unterschiedlichen Bezeichnungsweise – die Lösung

$$(18) \underline{b} = \left[ \begin{pmatrix} x_{11} - \bar{x}_1 & \dots & x_{n1} - \bar{x}_1 \\ \vdots & & \vdots \\ x_{1m} - \bar{x}_m & \dots & x_{nm} - \bar{x}_m \end{pmatrix} \begin{pmatrix} x_{11} - \bar{x}_1 & \dots & x_{1m} - \bar{x}_m \\ \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & \dots & x_{nm} - \bar{x}_m \end{pmatrix} \right]^{-1} \begin{pmatrix} x_{11} - \bar{x}_1 & \dots & x_{1m} - \bar{x}_m \\ \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & \dots & x_{nm} - \bar{x}_m \end{pmatrix} \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}$$

Man erkennt, daß es sich bei dem Ausdruck

$$\left[ \begin{pmatrix} x_{11} - \bar{x}_1 & \dots & x_{n1} - \bar{x}_1 \\ \vdots & & \vdots \\ x_{1m} - \bar{x}_m & \dots & x_{nm} - \bar{x}_m \end{pmatrix} \begin{pmatrix} x_{11} - \bar{x}_1 & \dots & x_{1m} - \bar{x}_m \\ \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & \dots & x_{nm} - \bar{x}_m \end{pmatrix} \right]$$

um das n-fache der Kovarianzmatrix der Regressoren handelt und bei dem Ausdruck

$$\begin{pmatrix} x_{11} - \bar{x}_1 & \dots & x_{1m} - \bar{x}_m \\ \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & \dots & x_{nm} - \bar{x}_m \end{pmatrix} \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}$$

um das n-fache des Kovarianzvektors der n Regressanden

mit dem Regressor.

Praktisch kann man die Konstante jedoch relativ einfach berücksichtigen, indem man eine Variable  $x_j$ , die nur aus Einsen besteht, in den Satz der Variablen aufnimmt. Die Gleichungen werden dann rechnerisch so behandelt, als wäre keine Konstante vorhanden (siehe Abschnitt 6.2).

Im Falle von standardisierten Variablen ist die Kovarianzmatrix gleich der Korrelationsmatrix. In diesem Fall gilt

$$(19) \underline{b} = \underline{R}^{-1} \underline{r} \quad \text{bzw.} \quad \underline{r} = \underline{R} \underline{b}.$$

### 3. Geometrischer Ansatz

Einem Ansatz von Kristof folgend, läßt sich die Aufgabe auch geometrisch lösen. Der Vorteil der geometrischen Lösung besteht darin, daß man allein mit Mitteln der linearen Algebra auskommt.

Im Falle eines inkonsistenten Systems ist die Lösung im Sinn der kleinsten Quadrate für  $\underline{b}$  in der Gleichung

$$(20) \quad \underline{X}\underline{b} \stackrel{\circ}{=} \underline{y},$$

gesucht. Die Abweichungen der Werte der Variablen  $Y$  von den Werten, die durch  $\underline{X}\underline{b}$  geschätzt werden, wird angegeben durch den Vektor  $\underline{\varepsilon}$

$$(21) \quad \underline{\varepsilon} = \underline{X}\underline{b} - \underline{y}.$$

Die nichtnegative skalare Funktion  $f(\underline{b}) = \underline{\varepsilon}'\underline{\varepsilon}$  soll so klein wie möglich werden; der Vektor  $\underline{b}$  sei die Lösung. Es sei

$$(22) \quad \underline{y} = \underline{g} + \underline{h},$$

wobei  $\underline{g}$  im durch die Spaltenvektoren von  $\underline{X}$  aufgespannten Raum liege und  $\underline{h}$  senkrecht dazu. Also ist

$$(23) \quad \underline{X}'\underline{h} = \underline{0} \quad \text{und} \quad \underline{g}'\underline{h} = \underline{0}.$$

Nun wird

$$\begin{aligned} f(\underline{b}) = \underline{\varepsilon}'\underline{\varepsilon} &= (\underline{X}\underline{b} - \underline{g} - \underline{h})'(\underline{X}\underline{b} - \underline{g} - \underline{h}) \\ &= ((\underline{X}\underline{b} - \underline{g})' - \underline{h}')(\underline{X}\underline{b} - \underline{g} - \underline{h}) \\ &= (\underline{X}\underline{b} - \underline{g})'(\underline{X}\underline{b} - \underline{g}) - (\underline{X}\underline{b} - \underline{g})'\underline{h} - \underline{h}'(\underline{X}\underline{b} - \underline{g}) + \underline{h}'\underline{h} \\ &= (\underline{X}\underline{b} - \underline{g})'(\underline{X}\underline{b} - \underline{g}) - 0 - 0 + \underline{h}'\underline{h}, \text{ also} \end{aligned}$$

$$(24) \quad f(\underline{b}) = (\underline{X}\underline{b}-\underline{g})'(\underline{X}\underline{b}-\underline{g}) + \underline{h}'\underline{h}.$$

Folglich wird  $f(\underline{b})$  minimal, wenn  $(\underline{X}\underline{b}-\underline{g})$  Null gesetzt wird. Dies wird erreicht, wenn  $\underline{X}\underline{b} = \underline{g}$ . Der Ausdruck ist wiederum die exakte Lösung nach (3). Dies ist möglich, da  $\underline{g}$  im durch die Spaltenvektoren von  $\underline{X}$  aufgespannten Raum liegt. Hierfür ist

$$(25) \quad \underline{b} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{g}.$$

Da nach (23)  $\underline{X}'\underline{h} = \underline{0}$ , gilt auch

$$(26) \quad \underline{0} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{h}.$$

Addiert man die Gleichungen (25) und (26), so erhält man

$$\underline{b} + \underline{0} = (\underline{X}'\underline{X})^{-1}\underline{X}'(\underline{g} + \underline{h}), \text{ bzw.}$$

$$(27) \quad \underline{b} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{y}$$

Sofern eine exakte Lösung existiert, ist dieses  $\underline{b}$  die Lösung von  $\underline{X}\underline{b} = \underline{y}$ , andernfalls ist es die Lösung von  $\underline{X}\underline{b} \stackrel{\circ}{=} \underline{y}$ , also die Lösung im Sinn der kleinsten Quadrate (vgl. (13)).

Unter (22) wurde der Vektor  $\underline{y}$  zerlegt in die Vektoren  $\underline{g} + \underline{h}$ , wobei  $\underline{g}$  im durch die Spaltenvektoren von  $\underline{X}$  aufgespannten Raum lag und  $\underline{h}$  senkrecht dazu. Man kann nun  $\underline{g}$  verwenden, um ein Maß zu bilden, das angibt, welcher Anteil des Vektors  $\underline{y}$  durch die Lösung  $\underline{X}\underline{b} \stackrel{\circ}{=} \underline{y}$  erfaßt wird. Der Quotient

$$(28) \quad \frac{\underline{g}'\underline{g}}{\underline{y}'\underline{y}}$$

gibt hierbei den erfaßten Anteil von  $\underline{y}$  an,  $\frac{\underline{h}'\underline{h}}{\underline{y}'\underline{y}}$  entsprechend den nicht erfaßten Anteil.



#### 4. Modellgüte

Gesucht ist ein Maß, das angibt, wie gut das Modell die gegebenen Kriterien  $Y$  schätzt.

Hierfür wird die Produkt-Moment-Korrelation zwischen dem Kriterium  $Y$  und den  $b_0 + b_1X_1 + \dots + b_mX_m$  betrachtet, also

$$(29) \quad R = \frac{S_{Y, b_1X_1 + \dots + b_mX_m}}{S_Y S_{b_1X_1 + \dots + b_mX_m}} \text{ bzw.}$$

$$(30) \quad R^2 = \frac{S^2_{Y, b_1X_1 + \dots + b_mX_m}}{S^2_Y S^2_{b_1X_1 + \dots + b_mX_m}} .^3$$

Man spricht hier vom multiplen  $R$  bzw.  $R^2$ .  $R^2$  ist der Anteil der gemeinsamen Varianz zwischen der Kriteriumsvariablen und den Prädiktorvariablen.  $R^2$  ist eine Schätzung des Teils der Varianz der Kriteriumsvariablen, der durch die Prädiktorvariablen vorhergesagt werden kann. Das multiple  $R^2$  entspricht dem Ausdruck (28) (ohne Beweis).

Möchte man eine Variable mittels der multiplen Regression in der Praxis vorherzusagen, sollte mindestens ein  $R^2$  von 0,7 erreicht werden.

Die Überprüfung des ermittelten  $R^2$  auf statistische Signifikanz gegenüber der Hypothese der Unabhängigkeit von Kriterium und Prädiktoren erfolgt – ohne Herleitung - durch einen F-Test:

$$(31) \quad F = \frac{R^2(n - k - 1)}{(1 - R^2)m}, \text{ df}_{\text{Zähler}} = m, \text{ df}_{\text{Nenner}} = n - m - 1.$$

---

<sup>3</sup> Die Konstante  $b_0$  entfällt, da die Produkt-Moment-Korrelation invariant gegenüber einer linearen Transformation ist.

## 5. Schrittweise Regression

Um zu entscheiden, welche Prädiktoren einen nennenswerten Beitrag für die Gesamtgüte leisten, wird zunächst noch einmal der multiple Korrelationskoeffizient  $R$  ( $0 \leq R \leq 1$ ) betrachtet.

Wird die Kriteriumsvariable  $Y$  durch nur zwei Prädiktorvariablen (1,2) vorhergesagt, so ist

$$(32) \quad R_{Y,12} = \sqrt{b_1 r_{1Y} + b_2 r_{2Y}}, \quad \text{bzw. allgemein bei } m \text{ Prädiktorvariablen}$$

$$(33) \quad R_{Y,12\dots m} = \sqrt{\sum_j b_j r_{jY}}, \quad \text{dieser Ausdruck entspricht dem multiplen } R,$$

wobei die  $b_j$  die *standardisierten* Koeffizienten bezeichnen.

Die Matrix  $\underline{R}$  ist die Matrix, die die Produkt-Moment-Korrelationen zwischen den Prädiktorvariablen enthält,  $\underline{r}$  sei der Vektor der Korrelationen der Prädiktoren mit dem Kriterium.

Nach (19) gilt für standardisierte Koeffizienten  $\underline{b}$ :

$$\underline{r} = \underline{R}\underline{b} \quad \text{bzw.}$$

$$(34) \quad \underline{b} = \underline{R}^{-1}\underline{r},$$

Man kann die Gleichung (34) verwenden, um standardisierte Koeffizienten aus nicht standardisierten zu gewinnen.

Die Prüfung der Regressionskoeffizienten auf Signifikanz erfolgt unter Verwendung standardisierter Koeffizienten  $\underline{b}$  mittels der Prüfgröße  $t$ , hier (ohne Herleitung)

$$(35) \quad t = \frac{b_j}{\sqrt{\frac{r_{jj}(1-R^2)}{n-m-1}}}, \quad df = n - m - 1,$$

wobei  $r^{jj}$  das Element  $(j, j)$  der invertierten Korrelationsmatrix  $\underline{R}^{-1}$  ist.

Nachdem man die standardisierten Regressionskoeffizienten auf Signifikanz geprüft hat, bildet man (für signifikante Korrelationen) die Produkte aus den standardisierten Regressionskoeffizienten und den Korrelationen der entsprechenden Prädiktoren mit dem Kriterium

$$(36) \quad b_j r_j.$$

Man kann nun zunächst diejenige Variable  $X_j$  mit signifikantem  $b_j$  aufnehmen, die den größten Beitrag  $b_j r_j$  leistet. Anschließend nimmt man weitere Variablen mit signifikantem  $b_j$  in das Regressionsmodell auf, bis der Zugewinn an Güte ( $R^2$ ) gegenüber dem vorherigen Modell einen zuvor festgelegten Wert (z.B. 0,05) nicht mehr übersteigt. Umgekehrt kann man zunächst von einem Modell ausgehen, daß alle  $m$  Prädiktorvariablen enthält und dann die Variablen mit nicht-signifikanten bzw. zu geringen Beiträgen ausschließen. Man spricht hier von einer „Rückwärtstechnik“, im ersten Fall von einer „Vorwärtstechnik“.

Es ist jedoch zu beachten, daß dieses Verfahren in der Praxis zu Problemen führen kann. Für Modelle, die nur einen Teil der  $m$  Prädiktorvariablen enthalten, ergeben sich in der Regel andere Gewichte  $b_j r_j$  als bei Einbeziehung aller  $m$  Variablen! Die Abweichung ist um so stärker, je stärker die Variablen  $X_j$  voneinander abhängig sind (Problem der Multikollinearität).

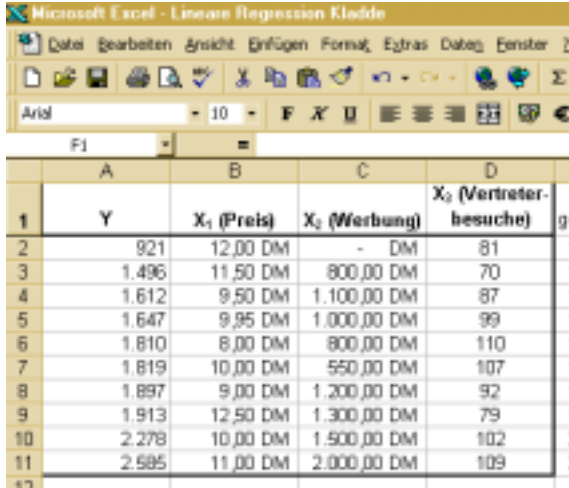
Die Verfahren der schrittweisen Regression sind daher keine exakten Verfahren, sondern müssen als explorative Verfahren verstanden werden.

## 6. Beispiel

Die Berechnung des Beispiels erfolgt auf Basis von MS-Excel. Die Durchführung erfolgt „manuell“ im Arbeitsblatt – die Erstellung von Makros bzw. VBA-Algorithmen wäre vorteilhaft, soll aber an dieser Stelle nicht explizit vorgestellt werden.

## 6.1 Ausgangsdaten

Gegeben sind folgende Ausgangsdaten:



	A	B	C	D
1	Y	X <sub>1</sub> (Preis)	X <sub>2</sub> (Werbung)	X <sub>3</sub> (Vertreterbesuche)
2	921	12,00 DM	- DM	81
3	1.496	11,50 DM	800,00 DM	70
4	1.612	9,50 DM	1.100,00 DM	87
5	1.647	9,95 DM	1.000,00 DM	99
6	1.810	8,00 DM	800,00 DM	110
7	1.819	10,00 DM	550,00 DM	107
8	1.897	9,00 DM	1.200,00 DM	92
9	1.913	12,50 DM	1.300,00 DM	79
10	2.278	10,00 DM	1.500,00 DM	102
11	2.585	11,00 DM	2.000,00 DM	109

Zu bestimmen ist eine Gleichung, die den Umsatz vorhersagt. (Die Tatsache, daß es sich bei  $n = 10$  Daten in der Praxis um eine recht geringe Menge handelt, sei außer Acht gelassen).

## 6.2 Bestimmung der Matrix $\underline{X}$

Anders als in Abschnitt 2 dargestellt, ist die erste Spalte der Matrix  $\underline{X}$  ein Vektor, der nur aus Einsen besteht. Hierdurch erreicht man, daß die Zahl  $b_0$  direkt im entsprechenden Lösungsvektor  $\underline{b}$  steht und nicht über Gleichung (16) bestimmt werden muß.

Man kopiere die Inhalte der Zellen (b2:d11) in den Bereich (b13:d22) und füge Einsen in die Zellen (a13:a22) ein:

	A	B	C	D	E
13					X
14	1,00	12,00	0,00	81,00	
15	1,00	11,50	800,00	70,00	
16	1,00	9,50	1100,00	87,00	
17	1,00	9,95	1000,00	99,00	
18	1,00	8,00	800,00	110,00	
19	1,00	10,00	550,00	107,00	
20	1,00	9,00	1200,00	92,00	
21	1,00	12,50	1300,00	79,00	
22	1,00	10,00	1500,00	102,00	
23	1,00	11,00	2000,00	109,00	

Dieser Bereich ist die Matrix  $\underline{X}$ ; der Bereich (a2:a11) entspricht dem Vektor  $\underline{y}$ .

### 6.3 Transponieren der Matrix $\underline{X}$ , Bildung der Matrix $\underline{X}'\underline{X}$

Man Markiere die Zellen (a25:j28). Nun ist die Formel „=MTRANS(a13:d22)“ einzugeben. Damit die Formel für den gesamten Bereich der Matrix eingegeben wird, muß die Eingabe mit der Tastenkombination Strg-Shift-Enter erfolgen:

	A	B	C	D	E	F	G	H	I	J
24										
25	=MTRANS(A13	1	1	1	1	1	1	1	1	1
26	12	11,5	9,5	9,95	8	10	9	12,5	10	11
27	0	800	1100	1000	800	550	1200	1300	1500	2000
28	81	70	87	99	110	107	92	79	102	109

Um die Matrix  $\underline{X}'\underline{X}$  zu bilden, muß der Bereich (a30:d33) markiert werden und hier das Produkt der beiden Matrizen  $\underline{X}'$  und  $\underline{X}$  hineingeschrieben werden. Der Befehl lautet „=mmult(a25:j28;a13:d22)“. Die Eingabe – wie alle Eingaben von Matrix-Formeln – mit Strg-Shift-Enter bestätigen!

	A	B	C	D	E
24					
25	1	1	1	1	
26	12	11,5	9,5	9,95	
27	0	800	1100	1000	
28	81	70	87	99	
29					
30	=MMULT(A25:J28,A13:D22)	103,45	10250	936	XX
31	103,45	1087,7525	105550	9573,05	
32	10250	105550	13172500	981650	
33	936	9573,05	981650	89370	

### 6.4 Bestimmung der Inversen von $\underline{X}'\underline{X}$

Markierung des Bereichs (a:35;d38). Die Formel zur Berechnung der Inversen lautet „=MINV(a30;d33)“.

	A	B	C	D	E
29					
30	10	103,45	10250	936	XX
31	103,45	1087,7525	105550	9573,05	
32	10250	105550	13172500	981650	
33	936	9573,05	981650	89370	
34					
35	=MINV(A30:D33)	-1,567664453	0,000669406	-0,15934446	(XX)^-1
36	-1,567664453	0,096491525	-3,65092E-05	0,006483671	
37	0,000669406	-3,65092E-05	4,33096E-07	-7,75259E-06	
38	-0,15934446	0,006483671	-7,75259E-06	0,001020696	

### 6.5 Bestimmung des Lösungsvektors $\underline{b}$

Die Berechnung des Lösungsvektors  $\underline{b}$  muß in MS-Excel in zwei Schritten erfolgen, da sich nur jeweils zwei Matrizen pro Operation multiplizieren lassen. Zunächst bilde man die Matrix  $(\underline{X}'\underline{X})^{-1}\underline{X}'$ , indem man in den Bereich (a40:j43) die folgende Formel einträgt („=mmult(a35:d38;a25:j28)“):

A	B	C	D	E	F	G	H	I	J	K
1	1	1	1	1	1	1	1	1	1	1
12	11,5	9,5	9,95	8	10	9	12,5	10	11	
0	800	1100	1000	800	550	1200	1300	1900	2000	
81	70	87	99	110	107	92	79	102	109	
10	103,45	10250	936							
103,45	1087,7525	109550	9573,05							
10290	106990	13172500	981660							
936	9673,05	981660	89370							
30,5981362	-1,567654453	0,000699406	-0,15934446							
-1,967654453	0,006491525	-3,65092E-05	0,006483671							
0,000699406	-3,65092E-05	4,33096E-07	-7,75299E-06							
-0,15934446	0,006483671	-7,75299E-06	0,001070698							
=MMULT(A35:J38,A25:J28)	1,901521918	2,525796924	-0,157721732	1,014534122	-1,807993	2,578842491	-0,770529484	-0,384434646	-2,73787129	
0,115421159	-0,033362343	-0,127065754	0,002189601	0,111725898	0,070933481	-0,146544086	0,103237616	0,003631386	0,127454001	
-0,000406664	4,33463E-05	0,0001145	-3,82702E-05	0,000138876	-0,00029701	0,000137301	0,000153612	0,000153196	0,000289665	
0,005186088	-0,01603549	-0,013126749	0,003414532	0,004039964	0,01529296	-0,011790366	-0,000291836	0,000024614	0,013126774	

Anschließend wird dieses Matrizenprodukt mit dem Vektor  $y$  multipliziert („=mmult(a40:j43;a2:a11)“):

A	B	C	D	E
-1,162619464	1,901521918	2,525796924	-0,157721732	1,014534122
0,115421159	-0,033362343	-0,127065754	0,002189601	-0,111725898
-0,000406664	4,33463E-05	0,0001145	-3,82702E-05	0,000138876
0,005186088	-0,01603549	-0,013126749	0,003414532	0,004039964
400,8900395				
32,47162872				
0,644467666				
12,84393663				

### 6.6 Bestimmung der Güte

Man kann nun die Schätzwerte der Variablen  $Y$  bestimmen, indem man in die Zellen (e2:e11) die entsprechende Formel („=\$a\$45+b2\*\$a\$46+c2\*\$a\$47+d2\*\$a\$48“) einträgt. Wird die folgende Formel in die erste Zelle eingetragen, so kann man sie einfach in die folgenden Zellen kopieren. Die Eingabe von „\$“ bewirkt, daß es sich für die entsprechenden Zellen um absolute Bezüge handelt.

Man kann das R bzw.  $R^2$  (hier in den Zellen f2 bzw. g2) bestimmen, indem man die Produkt-Moment-Korrelation der Reihe der Schätzwerte (e2:e11) und der Ausgangswerte (a2:a11) berechnen lässt. Der entsprechende Excel-Befehl lautet („=korrel(e2:e11;a2:a11)“).  $R^2$  erhält man durch Quadrieren der Zelle f2 („=f2^2“):

	D	E	F	G	H
1	X <sub>3</sub> (Vertreterbesuche)	geschätzt	R	R <sup>2</sup>	
2	81	1.029,13 DM	0,96637222	0,933875268	
3	70	1.387,18 DM			
4	87	1.733,93 DM			
5	99	1.838,22 DM			
6	110	1.787,29 DM			
7	107	1.652,58 DM			
8	92	1.846,36 DM			
9	79	1.857,48 DM			
10	102	2.200,61 DM			
11	109	2.645,22 DM			

### 6.7 Überprüfung des $R^2$ auf Signifikanz

Die Überprüfung des  $R^2$  auf Signifikanz erfolgt nach (31). In Excel wird der F-Wert mit der entsprechenden Formel in Zelle h2 berechnet:

	G	H	I
1	R <sup>2</sup>	F	
2	0,933875268	28,2458693	

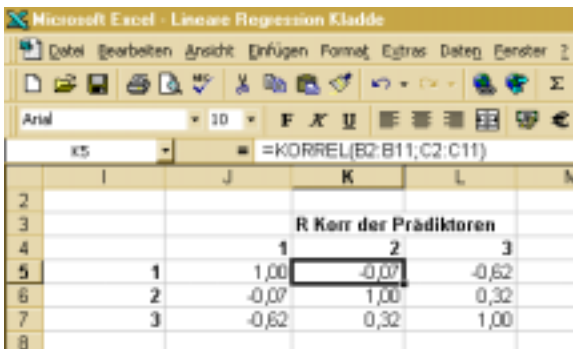
Die Überprüfung von F auf Signifikanz erfolgt in Excel durch die Funktion „Fwert(F; df<sub>Zähler</sub>; df<sub>Nenner</sub>)“ (Zelle i2). Die Prüfgröße ist mit 0,00062 hoch signifikant.

	G	H	I
1	R <sup>2</sup>	F	Signifikanz
2	0,933875268	28,2458693	0,000616573
3			



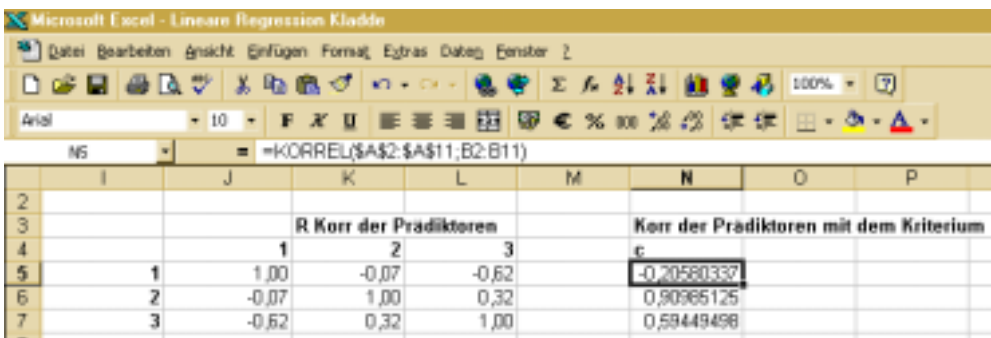
## 6.8 Bestimmung der standardisierten Koeffizienten

Um eine schrittweise Korrelation durchzuführen, muß man die standardisierten Koeffizienten  $b_i$  bestimmen. Hierzu benötigt man zunächst die Korrelationsmatrix der Prädiktoren. Hierfür werden die Korrelationen zwischen den Variablen  $X_1 \dots X_3$  paarweise mit dem Befehl „korrel“ (s.o.) berechnet. Die Korrelationsmatrix ist hier in den Bereich (j5:l7) eingetragen:



	1	2	3
1	1,00	-0,07	-0,62
2	-0,07	1,00	0,32
3	-0,62	0,32	1,00

Anschließend berechnet man die Korrelationen zwischen den Prädiktoren und dem Kriterium, die hier in den Bereich (n5:n7) eingetragen werden:



	1	2	3	c
1	1,00	-0,07	-0,62	-0,20580337
2	-0,07	1,00	0,32	0,80965125
3	-0,62	0,32	1,00	0,58448498

Nun ist die Inverse der Korrelationsmatrix der Prädiktoren zu bestimmen:

	I	J	K	L	M
2					
3			R Korr der Prädiktoren		
4			1	2	3
5	1	1,00	-0,07	-0,62	
6	2	-0,07	1,00	0,32	
7	3	-0,62	0,32	1,00	
8					
9		=MINV(J5:L7)	-0,24982902	1,140031093	R <sup>-1</sup> (1)
10			-0,249829017	1,15474246	-0,53113237
11			1,140031093	-0,53113237	1,88485595
12					

Die standardisierten Regressionskoeffizienten (hier im Bereich k13:k15) erhält man durch Multiplikation von  $\underline{R}^{-1}$  mit  $\underline{r}_c$ .

	I	J	K	L	M	N	O	P
2								
3			R Korr der Prädiktoren			Korr der Prädiktoren mit dem Kriterium		
4			1	2	3	c		
5	1	1,00	-0,07	-0,62		-0,20580337		
6	2	-0,07	1,00	0,32		0,90985125		
7	3	-0,62	0,32	1,00		0,58449498		
8								
9								R <sup>-1</sup> (1)
10								
11								
12								
13		b1 standardisiert	=MMULT(J9:L11;N5:N7)					
14		b2	0,786304					
15		b3	0,40086371					
16								

Die standardisierten Koeffizienten  $b_i$  werden mit den Korrelationen der Prädiktoren mit dem Kriterium  $r_{jc}$  multipliziert, hier im Bereich (n13:n15). Die Summe dieser Produkte ergibt  $R^2$ , die Wurzel aus dieser Summe R:

	I	J	K	L	M	N	O
2							
3			R Korr der Prädiktoren			Korr der Prädiktoren	
4			1	2	3	c	
5	1	1,00	-0,07	-0,62		-0,20680337	
6	2	-0,07	1,00	0,32		0,90885125	
7	3	-0,62	0,32	1,00		0,59449498	
8							
9		1,694608282	-0,24982902	1,14003109	R <sup>2</sup> (-t)		
10		-0,249829017	1,15474246	-0,53113237			
11		1,140031093	-0,53113237	1,88485595			
12							
13		b1 standardisiert	0,10167942		b1r1c	-0,02092597	
14		b2	0,786304		b2r2c	0,71541968	
15		b3	0,40266371		b3r3c	0,23838155	
16							
17		b1r1c			R <sup>2</sup>	0,93387527	
18		b2r2c			R	0,96637222	
19		b3r3c					

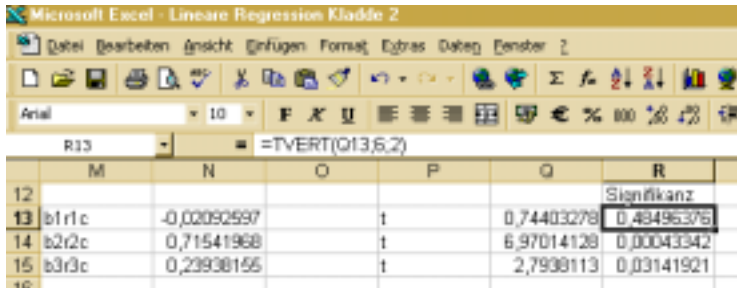
### 6.9 Prüfung der standardisierten Regressionskoeffizienten auf Signifikanz

Die Prüfung der standardisierten Regressionskoeffizienten auf Signifikanz erfolgt nach (35), wir unterstellen, daß die Variablen normalverteilt seien. In Excel erfolgt die Berechnung der t-Werte entsprechend in den Zellen im Bereich (q13:q15); die Werte ( $r^{ii}$ ) wurden zur besseren Berechnung in den Bereich (p9:p11) übertragen:

	I	J	K	L	M	N	O	P	Q
2									
3			R Korr der Prädiktoren			Korr der Prädiktoren mit dem Kriterium			
4			1	2	3	c			
5	1	1,00	-0,07	-0,62		-0,20680337			
6	2	-0,07	1,00	0,32		0,90885125			
7	3	-0,62	0,32	1,00		0,59449498			
8									
9		1,694608282	-0,24982902	1,14003109	R <sup>2</sup> (-t)		1,69460828		
10		-0,249829017	1,15474246	-0,53113237			1,15474246		
11		1,140031093	-0,53113237	1,88485595			1,88485595		
12									
13		b1 standardisiert	0,10167942		b1r1c	-0,02092597	t		0,74403228
14		b2	0,786304		b2r2c	0,71541968	t		6,97014128
15		b3	0,40266371		b3r3c	0,23838155	t		2,7938113
16									
17		b1r1c			R <sup>2</sup>	0,93387527			
18		b2r2c			R	0,96637222			
19		b3r3c							

Die Funktion „tvert(t-Wert;Freiheitsgrade;Seiten)“ prüft die t-Werte auf Signifikanz. Hier wurde die Funktion in die Zellen des Bereichs (r13:r15) eingetragen. Die Regressionskoeff-

fizienten  $b_2$  und  $b_3$  sind also auf einem Niveau von  $\alpha = 0,05$  (zweiseitig) signifikant; die Variablen „Werbeausgaben“ und „Vertreterbesuche“ leisten positive Beiträge von 0,72 bzw. 0,24, der Preis ist nicht signifikant und leistet einen negativen Beitrag zur Gesamtmodellgüte.



	M	N	O	P	Q	R
12						Signifikanz
13	b1r1c	-0,02092597		t	0,74403279	0,46495375
14	b2r2c	0,71541968		t	6,97014128	0,00043342
15	b3r3c	0,23938155		t	2,79981113	0,03141921

Man würde nun ein neues Modell erstellen, das nur die Variablen „Werbeausgaben“ und „Vertreterbesuche“ enthält.

### 6.10 Aufgabe

Erstellen Sie ein neues Modell, das nur die Prädiktoren „Werbeausgaben“ und „Vertreterbesuche“ enthält. Errechnen Sie die Modellgüte und die standardisierten Regressionskoeffizienten. Prüfen Sie die Regressionskoeffizienten auf Signifikanz und bestimmen Sie die Beiträge der beiden Variablen für das Modell!