



## Multiple lineare Regression

1

### Worum geht es in diesem Modul?

- Das Konzept der multiplen linearen Regression
- Spezialfälle
- Eigenschaften des Residuums
- Darstellung in Matrixnotation
- Identifikation der Regressionskoeffizienten
- Der multiple Determinationskoeffizient
- Multiple lineare Quasi-Regression
- Statistische Modelle zur multiplen linearen Regression
- Das Allgemeine Lineare Modell



## Multiple lineare Regression: Definition

2

Seien  $Y$  und  $X_1, \dots, X_m$  numerische Zufallsvariablen auf demselben Wahrscheinlichkeitsraum mit endlichen Erwartungswerten, positiven, endlichen Varianzen, sowie invertierbarer Kovarianzmatrix  $S_{xx}$ . Dann heißt die Regression  $E(Y | X_1, \dots, X_m)$  *linear in*  $(X_1, \dots, X_m)$ , falls

$$E(Y | X_1, \dots, X_m) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m.$$



## Multiple lineare Regression: Spezialfälle I

3

Bei der obigen Definition setzen wir also noch nicht voraus, dass die Regressoren  $X_1, \dots, X_m$  unabhängig voneinander definiert sind. In Kapitel 9 haben wir bereits darauf hingewiesen dass z. B. auch die einfache quadratische Regression  $E(Y|X)$  linear in  $(X, X^2)$  ist, falls also gilt:

$$E(Y|X) = \beta_0 + \beta_1 X + \beta_2 X^2.$$



## Multiple lineare Regression: Spezialfälle II

4

Ein weiterer Spezialfall, den wir bereits im Kapitel über die bedingte lineare Regression kennen gelernt haben, ist

$$E(Y|X_1, X_2) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 \cdot X_2.$$

Auch hier ist die Regression  $E(Y|X_1, X_2)$  linear in  $(X_1, X_2, X_1 \cdot X_2)$  und es gilt:

$$E(Y|X_1, X_2) = E(Y|X_1, X_2, X_1 \cdot X_2).$$



## Multiple lineare Regression: Spezialfälle III

5

Schließlich sei explizit auch noch einmal auf den Fall eines qualitativen Regressors  $X$  mit  $n$  Werten  $x_1, \dots, x_n$  hingewiesen. Folgendes Zellenmittelwertmodell

$$E(Y|X) = \beta_0 + \beta_1 \cdot I_1 + \beta_2 \cdot I_2 + \dots + \beta_n \cdot I_n, \quad \text{mit } \beta_0 = 0,$$

ist ebenfalls ein Spezialfall der multiplen linearen Regression. Die Regression  $E(Y|X)$  ist stets linear in  $(I_1, \dots, I_n)$ , auch dann, wenn die Regression  $E(Y|X)$  nicht linear in  $(X)$  ist.



## Eigenschaften des Residuums

6

Das Residuum

$$\mathbf{e} := Y - E(Y|X_1, \dots, X_m)$$

besitzt die bekannten Eigenschaften

$$E(\mathbf{e}) = 0,$$

$$E(\mathbf{e} | X_1, \dots, X_m) = 0,$$

$$\text{Cov}[\mathbf{e}, f(X_1, \dots, X_m)] = 0,$$

$$\text{Cov}(\mathbf{e}, X_i) = 0 \quad \text{für } i = 1, \dots, m,$$

wobei  $f(X_1, \dots, X_m)$  eine beliebige numerische Funktion der Regressoren bezeichnet.



## Darstellung in Matrixnotation I

7

Fasst man die Regressoren zu dem Zeilenvektor  $\mathbf{x}' = (X_1 \dots X_m)$  und die Regressionskoeffizienten  $\beta_1, \dots, \beta_m$  zu einem  $m$ -dimensionalen Spaltenvektor  $\mathbf{b} = (\beta_1 \dots \beta_m)'$  zusammen, so lässt sich die multiple lineare Regression nun in Matrix- bzw. Vektornotation auch folgendermaßen schreiben:

$$E(\mathbf{y} | \mathbf{x}) = \beta_0 + \mathbf{x}' \mathbf{b} = \beta_0 + (X_1 \dots X_m) \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}.$$

Auch hier wird der Regressand als Vektor  $\mathbf{y} = (Y)$  aufgefasst, der eben nur aus einer einzigen Komponente, nämlich  $Y$ , besteht. Daher ist auch  $\beta_0$  eine reelle Zahl.



## Darstellung in Matrixnotation II

8

Definieren wir den Zeilenvektor  $\mathbf{z}' := (1 \ X_1 \dots X_m)$  und den Spaltenvektor  $\mathbf{g} := (\beta_0 \ \beta_1 \dots \beta_m)'$ , so können wir die letzte Gleichung auch wie folgt schreiben:

$$E(\mathbf{y} | \mathbf{x}) = E(\mathbf{y} | \mathbf{z}) = \mathbf{z}' \mathbf{g} = (1 \ X_1 \dots X_m) \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}.$$



## Identifikation der Regressionskoeffizienten I

9

Zur Identifikation von  $\beta_0$  und der Komponenten von  $\mathbf{b} = (\beta_1 \dots \beta_m)'$  greift man auf die Erwartungswerte des Regressanden und der Regressoren sowie die Kovarianzmatrizen  $S_{xx}$  und  $S_{xy}$  zurück. Für die Konstante  $\beta_0$  ergibt sich:

$$\beta_0 = E(y) - E(\mathbf{x})' \mathbf{b} = E(Y) - [E(X_1) \dots E(X_m)] \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}.$$



## Identifikation der Regressionskoeffizienten II

10

Multipliziert man beide Seiten der folgenden Gleichung

$$\begin{aligned} S_{xy} = \text{Cov}(\mathbf{x}, y) &= \text{Cov}(\mathbf{x}, \beta_0 + \mathbf{x}' \mathbf{b} + e) = \text{Cov}(\mathbf{x}, \beta_0 + \mathbf{b}' \mathbf{x} + e) \\ &= \text{Cov}(\mathbf{x}, \mathbf{x}) \mathbf{b} = S_{xx} \mathbf{b} \end{aligned}$$

mit  $S_{xx}^{-1}$  vor, so folgt

$$S_{xx}^{-1} S_{xx} \mathbf{b} = S_{xx}^{-1} S_{xy}.$$

Da  $S_{xx}^{-1} S_{xx} = \mathbf{I}$  die Einheitsmatrix ist, folgt daraus

$$\mathbf{b} = S_{xx}^{-1} S_{xy}.$$

Damit haben wir eine Formel zur Berechnung der Regressionskoeffizienten aus den Varianzen und Kovarianzen der Regressoren und des Regressanden.



## Identifikation der Regressionskoeffizienten bei zwei Regressoren I

11

Im Fall mit zwei Regressoren  $X_1$  und  $X_2$  erhält man die bereits aus Kapitel 9 bekannte Gleichung

$$\beta_0 = E(Y) - E(x)' b = E(Y) - [\beta_1 E(X_1) + \beta_2 E(X_2)] = E(Y) - \beta_1 E(X_1) - \beta_2 E(X_2).$$

Die Kovarianz und die Varianz-Kovarianzmatrix lauten

$$S_{xx} = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) \end{pmatrix} \text{ und } S_{xy} = \begin{pmatrix} \text{Cov}(X_1, Y) \\ \text{Cov}(X_2, Y) \end{pmatrix}.$$

Die Kramersche Regel aus dem letzten Kapitel liefert:

$$S_{xx}^{-1} = \frac{1}{\text{Var}(X_1) \text{Var}(X_2) - \text{Cov}(X_1, X_2)^2} \begin{pmatrix} \text{Var}(X_2) & -\text{Cov}(X_1, X_2) \\ -\text{Cov}(X_1, X_2) & \text{Var}(X_1) \end{pmatrix}.$$



## Identifikation der Regressionskoeffizienten bei zwei Regressoren II

12

Multipliziert man  $S_{xx}^{-1}$  mit  $S_{xy}$  folgen:

$$\beta_1 = \frac{\text{Var}(X_2) \text{Cov}(X_1, Y) - \text{Cov}(X_2, Y) \text{Cov}(X_1, X_2)}{\text{Var}(X_1) \text{Var}(X_2) - \text{Cov}(X_1, X_2)^2},$$

$$\beta_2 = \frac{\text{Var}(X_1) \text{Cov}(X_2, Y) - \text{Cov}(X_1, Y) \text{Cov}(X_1, X_2)}{\text{Var}(X_1) \text{Var}(X_2) - \text{Cov}(X_1, X_2)^2}.$$



## Der multiple Determinationskoeffizient

13

Für die Varianz  $Var[E(y|x)]$  der Regression gilt:

$$\begin{aligned} Var[E(y|x)] &= Var(\beta_0 + x' b) = Var(x' b) = Var(b' x) = b' Var(x) b = \\ &= b' S_{xx} b, \end{aligned}$$

Dabei sind  $Var(x) = S_{xx}$  die  $(m \times m)$  Varianz-Kovarianzmatrix der Regressoren  $X_1, \dots, X_m$  und  $b$  der  $m$ -dimensionale Spaltenvektor der Regressionskoeffizienten  $\beta_1, \dots, \beta_m$ . Der multiple Determinationskoeffizient ergibt sich dann wie folgt:

$$R_{Y|X_1, \dots, X_m}^2 = [b' S_{xx} b] / Var(Y).$$



## Der multiple Determinationskoeffizient: Spezialfälle

14

Im Fall mit zwei Regressoren  $X_1$  und  $X_2$  erhält man die bereits aus Kapitel 9 bekannte Gleichung

$$R_{Y|X_1, X_2}^2 = [\beta_1^2 Var(X_1) + \beta_2^2 Var(X_2) + 2\beta_1 \beta_2 Cov(X_1, X_2)] / Var(Y).$$

Für den Spezialfall, dass alle Regressoren *paarweise unkorreliert* sind, vereinfacht sich der Ausdruck für den multiplen Determinationskoeffizienten zu

$$R_{Y|X_1, \dots, X_m}^2 = \left( \sum_{i=1}^m \beta_i^2 Var(X_i) \right) \frac{1}{Var(Y)}.$$



## Multiple lineare Quasi-Regression: Definition

15

**Definition 14.2.** Unter den gleichen Voraussetzungen wie in Definition 14.1 definieren wir die *multiple lineare Quasi-Regression*, die wir mit  $Q(Y|X_1, \dots, X_m)$  oder  $Q(y|x)$  bezeichnen, als diejenige Linearkombination  $\beta_0 + \beta_1 X_1 + \dots + \beta_m X_m = \beta_0 + \mathbf{b}'\mathbf{x}$  der Komponenten von  $\mathbf{x} = (X_1 \dots X_m)'$ , die folgendes erfüllt:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m + \mathbf{n}$$

mit

$$E(\mathbf{n}) = 0,$$

und

$$\text{Cov}(\mathbf{n}, X_1) = \dots = \text{Cov}(\mathbf{n}, X_m) = 0.$$



## Multiple lineare Quasi-Regression: Alternative Definition

16

**Definition 14.3.** Unter den gleichen Voraussetzungen wie zuvor, können wir  $Q(y|x)$  auch als diejenige Linearkombination  $\beta_0 + \mathbf{x}'\mathbf{b}$  definieren, welche die folgende Funktion von  $b_0$  und  $\mathbf{b}$ , das *Kleinst-Quadrat-Kriterium*, minimiert:

$$LS(b_0, \mathbf{b}) = E[[Y - (b_0 + \mathbf{x}'\mathbf{b})]^2].$$

Diejenige Zahl  $b_0$  und derjenige Vektor  $\mathbf{b}$ , für welche die Funktion  $LS(b_0, \mathbf{b})$  ein Minimum annimmt, seien mit  $\beta_0$  und  $\mathbf{b}$ , respektive, bezeichnet. Die multiple lineare Quasi-Regression ist dann definiert durch:

$$Q(y|x) = \beta_0 + \mathbf{x}'\mathbf{b}.$$





## Der Determinationskoeffizient der multiplen linearen Quasi-Regression

17

Für die Bestimmung der Regressionskoeffizienten der multiplen linearen Quasi-Regression gelten übrigens analog die gleichen Formeln wie für die entsprechenden Koeffizienten der echten multiplen linearen Regression. Ebenfalls gleich ist die Berechnungsformel für den Determinationskoeffizienten der linearen Quasi-Regression, d. h. es gilt

$$Q_{Y|X_1, \dots, X_m}^2 := \text{Var}[Q(\mathbf{y} | \mathbf{x})] / \text{Var}(Y) = [\mathbf{b} \hat{\mathbf{c}} \mathbf{S}_{xx} \mathbf{b}] / \text{Var}(Y).$$



## Statistische Modelle zur multiplen linearen Regression

18

Bisher haben wir nur ein Einzelexperiment betrachtet: Ziehen einer Beobachtungseinheit  $u$  aus der Population und Registrierung der Werte des Regressanden und der Regressoren. Statistische Modelle beziehen sich jedoch auf  $N$  Zufallsexperimente, in denen Informationen über die zu schätzenden Parameter gesammelt werden.



## Modelle mit stochastischen Regressoren

19

Modelle mit *stochastischen Regressoren* bestehen aus der  $N$ -maligen Wiederholung unseres bisher betrachteten Einzelexperiments. Dies führt dazu, dass man nicht mehr nur einen einzigen Regressanden  $Y$  und  $m$  Regressoren betrachten muss, sondern  $N$  Vektoren  $(Y_i, X_{i1}, \dots, X_{im})$ ,  $i = 1, \dots, N$ , die jeweils das Ergebnis des  $i$ -ten Zufallsexperiments repräsentieren. Über diese Vektoren kann man unterschiedliche Verteilungsannahmen machen, z. B. dass die  $(Y_i, X_{i1}, \dots, X_{im})$  unabhängig sind und jeder dieser Vektoren  $(m + 1)$ -*variater normalverteilt* ist.



## Modelle mit festen Regressoren

20

Mit einem anderen, weitaus häufiger verwendeten statistischen Modell, schätzt man innerhalb der Wertekombinationen  $x_1, \dots, x_m$  der Regressoren  $X_1, \dots, X_m$  die Erwartungswerte  $E(Y | X_1 = x_1, \dots, X_m = x_m)$  von  $Y$ , indem man innerhalb dieser Wertekombinationen den Regressanden  $Y$  mehrfach beobachtet. Die Werte  $x_1, \dots, x_m$  der Regressoren sind dabei also nicht mehr zufällig, sondern werden als feste Größen betrachtet, die das Design des Experiments charakterisieren. Man spricht daher auch von Modellen mit festen oder nicht-stochastischen Regressoren. Das wichtigste dieser Modelle mit festen Regressoren ist das *Allgemeine Lineare Modell*.



## Das Allgemeine Lineare Modell

21

Das Allgemeine Lineare Modell (ALM) ist durch die folgenden Annahmen definiert:

$$\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{e}$$
$$\mathbf{e} \sim \mathcal{N}(\mathbf{0}, s^2 \mathbf{I}),$$

Dabei bezeichnet  $\mathbf{y} = (Y_1 \dots Y_i \dots Y_N)'$  nun den Spaltenvektor der für eine Stichprobe des Gesamtumfangs  $N$  zu erhebenden „abhängigen“ Variablen. Die so genannte *Designmatrix*  $\mathbf{X}$  besteht aus  $N \times (m + 1)$  festen Zahlen. Dabei besteht jede Zeile von  $\mathbf{X}$  aus den Vektoren  $\mathbf{x}_i' := (1 \ x_{i1} \dots x_{im})$ , eben den Wertekombinationen der Regressoren, innerhalb derer die Beobachtung  $Y_i$  erhoben wird und der vorangestellten Konstanten 1, die dazu führt, dass die Regressionskonstante  $\beta_0$  die erste Komponente von  $\boldsymbol{\beta} = (\beta_0 \ \beta_1 \dots \beta_m)'$  ist. Der Vektor  $\mathbf{e} = (e_1 \dots e_i \dots e_N)'$  schließlich enthält die Residuen  $Y_i - (\mathbf{x}_i' \boldsymbol{\beta})$ .



## Das Allgemeine Lineare Modell II

22

Die zweite Annahme besagt, dass  $\mathbf{e}$  mit Erwartungswertvektor  $E(\mathbf{e}) = \mathbf{0}$  und der  $N \times N$ -Kovarianzmatrix  $\mathbf{S}_{ee} = s^2 \mathbf{I}$  multivariat normalverteilt ist. Die Residuen  $e_i$  sind also unkorreliert und haben gleiche Varianzen. Letzteres ist die so genannte Homoskedastizitätsannahme.

Folgerungen aus den Annahmen des ALM sind zunächst:

$$E(\mathbf{y}) = \mathbf{X} \boldsymbol{\beta},$$

und

$$\mathbf{S}_{yy} = s^2 \mathbf{I}.$$



Ist  $\mathbf{X}'\mathbf{X}$  invertierbar, so gilt

$$\hat{b} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' y.$$

zur Schätzung des Vektors der Regressionskoeffizienten. Diese Formel erhält man durch die Minimierung der Kleinst-Quadrat-Funktion

$$LS(\mathbf{b}) = (\mathbf{y} - \mathbf{X}\mathbf{b})' (\mathbf{y} - \mathbf{X}\mathbf{b}).$$

Weiter ist noch

$$S_{\hat{b}\hat{b}} = s^2 (\mathbf{X}'\mathbf{X})^{-1}$$

von Bedeutung, die Kovarianzmatrix dieser Schätzer. Die Wurzeln aus den diagonalen Komponenten von  $S_{\hat{b}\hat{b}}$  sind die Standardschätzfehler der Regressionskoeffizienten.



Schließlich sei noch die Formel zur Schätzung des Determinationskoeffizienten genannt:

$$\hat{R}^2 = \frac{\mathbf{y}'\mathbf{X}\hat{\mathbf{b}} - N \cdot \bar{Y}^2}{\mathbf{y}'\mathbf{y} - N \cdot \bar{Y}^2} = \frac{\text{QuadratsummederRegression}}{\text{QuadratsummeGesamt}},$$

wobei  $\bar{Y} = (1/N) \cdot \sum_{i=1}^N Y_i$ .



## Formulierung von Hypothesen I

25

Wie wir in diesem und den vorangegangenen Kapiteln gesehen haben, lassen sich mit der multiplen linearen Regression durchaus auch komplexe und nichtlineare Abhängigkeiten beschreiben. Dabei gibt es zwei grundsätzliche Strategien.

Erste Strategie: Vergleich zwischen

$$Q(Y | X_1, X_2, \dots, X_{m-p}) = \gamma_0 + \gamma_1 X_1 + \dots + \gamma_{m-p} X_{m-p}$$

und

$$E(Y | X_1, \dots, X_m) = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m.$$

Die zu testende Nullhypothese kann dann in drei Versionen formuliert werden:

$$H_0: \beta_{m-p+1} = \beta_{m-p+2} = \dots = \beta_m = 0 \quad \text{Version 1}$$

$$H_0: Q(Y | X_1, \dots, X_{m-p}) = E(Y | X_1, X_2, \dots, X_m). \quad \text{Version 2}$$

$$H_0: R_{Y|X_1, \dots, X_m}^2 - Q_{Y|X_1, \dots, X_{m-p}}^2 = 0 \quad \text{Version 3}$$



## Formulierung von Hypothesen II

26

Die zweite Strategie ist noch allgemeiner: Man formuliert eine  $H_0$  in der Form

$$H_0: \mathbf{A} \mathbf{b} - \mathbf{d} = \mathbf{0},$$

der *Allgemeinen Linearen Hypothese* und testet diese mit einem Programm wie z. B. Systat oder SPSS (über Syntax) direkt, indem man die Matrix  $\mathbf{A}$  und den Vektor  $\mathbf{d}$  gemäß seiner Hypothese spezifiziert. Die Matrix  $\mathbf{A}$  muss  $p \leq m$  linear unabhängige Zeilen enthalten.



## Signifikanztests im ALM I

27

Will man die Nullhypothese  $H_0: \beta_{m-p+1} = \beta_{m-p+2} = \dots = \beta_m = 0$ , testen, schätzt man zunächst den Determinationskoeffizienten  $\hat{R}_E^2$  für die Regression und dann  $\hat{R}_Q^2$  für die multiple lineare Quasi-Regression. Dabei beachte man, dass dies jeweils mit unterschiedlichen Designmatrizen und unterschiedlichen Regressionskoeffizienten geschieht. Mit diesen Schätzungen  $\hat{R}_E^2$  bzw.  $\hat{R}_Q^2$  der beiden Determinationskoeffizienten geht man dann in die Formel

$$F = \frac{(\hat{R}_E^2 - \hat{R}_Q^2) / p}{(1 - \hat{R}_E^2) / (N - m - 1)},$$

die unter den Annahmen des ALM und der Gültigkeit der Nullhypothese eine  $F$ -verteilte Teststatistik liefert, mit den Zählerfreiheitsgraden  $df_1 = p$  und den Nennerfreiheitsgraden  $df_2 = N - m - 1$ . Dabei sind  $m$  die Anzahl der Regressoren in der multiplen linearen Regression,  $p$  die Anzahl der Parameter, die laut Nullhypothese gleich null sein sollen und  $N$  der Stichprobenumfang.



## Signifikanztests im ALM II

28

Bei der *zweiten Strategie* berechnet man für die jeweilige *Allgemeine Lineare Hypothese* (ALH)

$$H_0: \mathbf{A} \mathbf{b} - \mathbf{d} = \mathbf{0},$$

die Prüfgröße

$$F = \frac{\hat{Q}_h / p}{\hat{Q}_e / (N - m - 1)},$$

Dabei sind  $p$  die Anzahl der (linear unabhängigen) Zeilen der Matrix  $\mathbf{A}$  der ALH (und damit die Anzahl der simultan geprüften Einzelhypothesen),

$$\hat{Q}_h = (\mathbf{A}\hat{\mathbf{b}} - \mathbf{d})' [\mathbf{A}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{A}']^{-1} (\mathbf{A}\hat{\mathbf{b}} - \mathbf{d}) \quad \text{Hypothesenquadratsumme}$$

$$\text{und } \hat{Q}_e = \mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{X}\hat{\mathbf{b}}. \quad \text{Fehlerquadratsumme}$$

Auch die letztgenannte Prüfgröße  $F$  ist unter den Annahmen des ALM und der Gültigkeit der Nullhypothese eine  $F$ -verteilte Teststatistik, mit den Zählerfreiheitsgraden  $df_1 = p$  und den Nennerfreiheitsgraden  $df_2 = N - m - 1$ .