

Statistik II für Betriebswirte

Vorlesung 9

Dr. Andreas Wünsche

TU Bergakademie Freiberg
Institut für Stochastik

9. Dezember 2019



Einfache lineare Regression durch den Koordinatenursprung

- ▶ Bei bestimmten Problemstellungen ist es sinnvoll zu fordern, dass die Regressionsgerade durch den Koordinatenursprung geht. Man spricht dann auch von einer **Regression ohne Absolutglied** oder einer **eigentlich-linearen Regression**.
- ▶ Man erhält nun als Modellansatz

$$Y_i = b_1 x_i + \varepsilon_i, \quad i = 1, \dots, n;$$

als Schätzung für den Parameter b_1

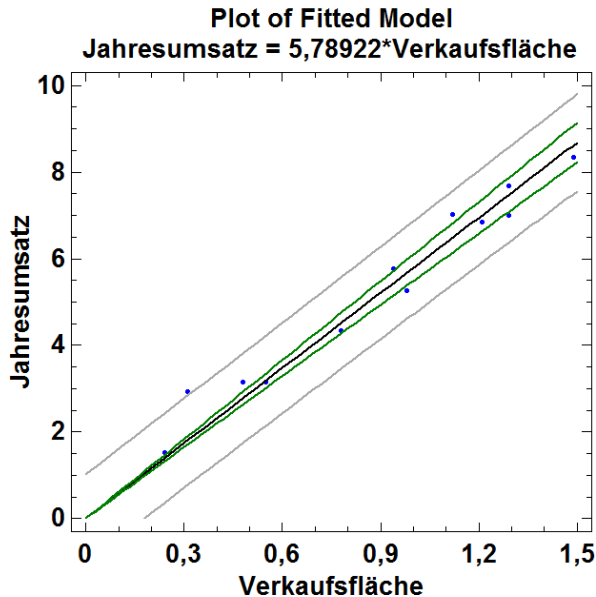
$$\hat{b}_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

und als Schätzung für die Varianz der zufälligen Fehler

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{mit} \quad \hat{y}_i = \hat{b}_1 x_i.$$



Regressionsgerade ohne Absolutglied im Beispiel 7.1



Transformationen auf Linearität

- ▶ Ist die gesuchte Abhängigkeitsfunktion eine nichtlineare Funktion (eine Gerade ist schlecht an die Daten anpassbar), kann man mitunter durch geeignete Variablentransformationen die Aufgabenstellung in eine der einfachen linearen Regression transformieren. Diese ist dann aber nicht äquivalent zur ursprünglichen Aufgabenstellung.
- ▶ Nichtlineare, in lineare transformierbare Funktionen sind z.B.

$$y = \alpha x^{\beta} \quad \Rightarrow \quad \ln y = \ln \alpha + \beta \ln x$$

$$y = \alpha e^{\beta x} \quad \Rightarrow \quad \ln y = \ln \alpha + \beta x$$

$$y = (\alpha + \beta x)^{-1} \quad \Rightarrow \quad y^{-1} = \alpha + \beta x$$

$$y = x(\alpha + \beta x)^{-1} \quad \Rightarrow \quad y^{-1} = \alpha x^{-1} + \beta$$

$$y = \alpha e^{\beta/x} \quad \Rightarrow \quad \ln y = \ln \alpha + \beta x^{-1}$$

$$y = (\alpha + \beta e^{-x})^{-1} \quad \Rightarrow \quad y^{-1} = \alpha + \beta e^{-x}$$



7.2. Multiple parameterlineare Regression

- ▶ Im Folgenden soll die Abhängigkeit eines Regressanden (einer Wirkungsgröße oder einer endogenen Variablen) Y von mehreren Regressoren (Einflussgrößen oder exogenen Variablen) X_1, \dots, X_k beschrieben werden, d.h. es soll gelten

$$Y \approx f(X_1, \dots, X_k)$$

mit einer geeigneten Funktion $f : \mathbb{R}^k \rightarrow \mathbb{R}$.

- ▶ Wir werden größtenteils wieder annehmen, dass die Regressoren deterministisch sind (z.B. mit exakt einstellbaren Werten), und dies durch kleine Buchstaben x_1, \dots, x_k in den Gleichungen kennzeichnen.
- ▶ Man erhält dann als Modellgleichung

$$Y(x_1, \dots, x_k) = f(x_1, \dots, x_k) + \varepsilon$$

mit einem zufälligen Fehler ε .



Parameterlineare Ansätze

- ▶ Häufig werden bei solchen Aufgabenstellungen **parameterlineare Ansätze** verwendet, d.h. man setzt eine Beziehung

$$Y(x_1, \dots, x_k) = a_1 f_1(x_1, \dots, x_k) + \dots + a_r f_r(x_1, \dots, x_k) + \varepsilon$$

mit speziell gewählten, bekannten Funktionen f_1, \dots, f_r und zu bestimmenden Koeffizienten (Parametern) a_1, \dots, a_r (die linear in die Gleichung eingehen) voraus.

- ▶ Die einfache lineare Regression mit der Modellgleichung

$$Y(x) = b_0 + b_1 x + \varepsilon$$

ist ein Spezialfall davon, dort gelten $k = 1$, $r = 2$, $f_1(x) = 1$, $f_2(x) = x$, $a_1 = b_0$, $a_2 = b_1$.



Beispiele für parameterlineare Ansätze

- ▶ Im eigentlich nichtmultiplen Fall $k = 1$ (nur eine Einflussgröße) gilt bei der **quadratischen Regression**

$$Y(x) = b_0 + b_1x + b_2x^2 + \varepsilon$$

und allgemeiner bei der **polynomialen Regression vom Grade m**

$$Y(x) = b_0 + b_1x + \dots + b_mx^m + \varepsilon.$$

- ▶ Der **k -faktorielle Ansatz ohne Wechselwirkungen**

$$Y(x_1, \dots, x_k) = b_0 + b_1x_1 + \dots + b_kx_k + \varepsilon$$

wird zur Bestimmung der **Ausgleichsebene** (für die **ebene Regression**) genutzt.

- ▶ **Bemerkung:** Eine Gleichung der Form $y = b_0 + b_1x_1 + \dots + b_kx_k$ definiert eine (Hyper-)Ebene im $(k + 1)$ -dimensionalen Raum von Punkten mit Koordinaten (x_1, \dots, x_k, y) .



Weitere Beispiele für parameterlineare Ansätze

- ▶ Als Beispiel eines k -faktoriellen Ansatzes mit Wechselwirkungen werde hier noch der Fall einer multiplen quadratischen Regression vorgestellt:

$$\begin{aligned} Y(x_1, \dots, x_k) = & b_0 + b_1 x_1 + \dots + b_k x_k \\ & + b_{12} x_1 x_2 + \dots + b_{k-1,k} x_{k-1} x_k \\ & + b_{11} x_1^2 + \dots + b_{kk} x_k^2 \\ & + \varepsilon. \end{aligned}$$

- ▶ Auch höhere Polynomgrade oder andere Funktionen der Variablen x_1, \dots, x_k sind möglich und werden auch verwendet.

Regressionsansatz in Vektorschreibweise

- ▶ Es ist vorteilhaft, die Vektorschreibweise zu nutzen. Es seien

$$\underline{x} = (x_1, \dots, x_k)^T = \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix}, \quad \underline{a} = (a_1, \dots, a_r)^T = \begin{pmatrix} a_1 \\ \vdots \\ a_r \end{pmatrix},$$

$$\underline{f}(\underline{x}) = (f_1(\underline{x}), \dots, f_r(\underline{x}))^T = \begin{pmatrix} f_1(\underline{x}) \\ \vdots \\ f_r(\underline{x}) \end{pmatrix}.$$

- ▶ Der parameterlineare Ansatz kann dann geschrieben werden als

$$Y(\underline{x}) = \underline{a}^T \underline{f}(\underline{x}) + \varepsilon, \quad (1)$$

wobei die Definition der Multiplikation von Vektoren (als Spezialfall der Matrixmultiplikation oder als Skalarprodukt) genutzt wird.



Die Methode der kleinsten Quadrate

- ▶ Sind die (zufallsbeeinflussten) „Wirkungen“ y_i für $i = 1, \dots, n$ an den „Einflussstellen“ $\underline{x}_i = (x_{1i}, \dots, x_{ki})^T$ durch Messungen bestimmt worden, kann man mit Hilfe der Methode der kleinsten Quadrate eine geeignete Schätzung $\hat{\underline{a}}$ des **Vektors \underline{a} der Regressionskoeffizienten** im parameterlinearen Ansatz (1) finden.
- ▶ Die Schätzung $\hat{\underline{a}}$ ist ein Vektor von Regressionskoeffizienten \underline{a} , für den $\sum_{i=1}^n (y_i - \underline{a}^T \underline{f}(\underline{x}_i))^2$ minimal wird.

- ▶ Die **geschätzte Regressionsfunktion** ist dann

$$\hat{y}(\underline{x}) = \hat{a}_1 f_1(\underline{x}) + \dots + \hat{a}_r f_r(\underline{x}) = \hat{\underline{a}}^T \underline{f}(\underline{x}) = \underline{f}(\underline{x})^T \hat{\underline{a}}.$$

- ▶ Im Weiteren genutzte Bezeichnungen sind $\underline{y} = (y_1, \dots, y_n)^T$ und

$$F = (\underline{f}(\underline{x}_1), \dots, \underline{f}(\underline{x}_n))^T = \begin{pmatrix} f_1(\underline{x}_1) & \dots & f_r(\underline{x}_1) \\ \vdots & \ddots & \vdots \\ f_1(\underline{x}_n) & \dots & f_r(\underline{x}_n) \end{pmatrix}.$$

Das Normalgleichungssystem

- ▶ Die Schätzung $\hat{\underline{a}}$ des Vektors \underline{a} der Regressionskoeffizienten kann dann mit Hilfe des sogenannten **Normalgleichungssystems** gefunden werden:

$$F^T F \hat{\underline{a}} = F^T \underline{y}. \quad (2)$$

Dies ist ein lineares Gleichungssystem zur Bestimmung der Komponenten von $\hat{\underline{a}}$.

- ▶ Ist die Matrix $F^T F$ regulär, dann ist (2) eindeutig auflösbar und es gilt

$$\hat{\underline{a}} = \left(F^T F\right)^{-1} F^T \underline{y}. \quad (3)$$



Eigenschaften der Schätzung

- ▶ Unter der Annahme, dass die beobachteten Werte y_i Realisierungen der Zufallsgrößen

$$Y_i = a_1 f_1(\underline{x}_i) + \dots + a_r f_r(\underline{x}_i) + \varepsilon_i$$

sind, wobei die zufälligen Fehler ε_j unabhängige normalverteilte Zufallsgrößen mit Erwartungswert 0 und konstanter Varianz σ^2 sind, ist die Schätzung $\hat{\underline{a}}$ aus (3) erwartungstreu und konsistent.



Beispiel 7.2 Jahresumsatz

- ▶ Fortsetzung Beispiel 6.4 (und 7.1)
Daten aus BLEYMÜLLER ET AL, Statistik für
Wirtschaftswissenschaftler, 2004, Kap. 20.

- ▶ i Filiale
 x_{1i} Verkaufsfläche in Tsd. qm
 x_{2i} Passantenfrequenz in Tsd. Passanten pro Tag
 y_i Jahresumsatz in Mio. €

i	1	2	3	4	5	6
x_{1i}	0.31	0.98	1.21	1.29	1.12	1.49
x_{2i}	10.24	7.51	10.81	9.89	13.72	13.92
y_i	2.93	5.27	6.85	7.01	7.02	8.35
i	7	8	9	10	11	12
x_{1i}	0.78	0.94	1.29	0.48	0.24	0.55
x_{2i}	8.54	12.36	12.27	11.01	8.25	9.31
y_i	4.33	5.77	7.68	3.16	1.52	3.15



Fortsetzung Beispiel 7.2 Jahresumsatz

- Wir wählen als Ansatz mit $\underline{x} = (x_1, x_2)^T$

$$Y(\underline{x}) = a_1 + a_2 x_1 + a_3 x_2 + \varepsilon = (a_1, a_2, a_3) \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} + \varepsilon.$$

- Dann erhält man

$$F^T = \begin{pmatrix} 1 & 1 & 1 & \dots \\ 0.31 & 0.98 & 1.21 & \dots \\ 10.24 & 7.51 & 10.81 & \dots \end{pmatrix},$$

$$F^T F = \begin{pmatrix} 12.00 & 10.68 & 127.83 \\ 10.68 & 11.41 & 118.97 \\ 127.83 & 118.97 & 1410.14 \end{pmatrix},$$

$$\hat{\underline{a}}^T = (-0.83, 4.74, 0.175)$$

und so $\hat{y}(x_1, x_2) = -0.83 + 4.74 x_1 + 0.175 x_2$

als geschätzte Regressionsfunktion.



Beispiel 7.2 in Statgraphics

Multiple Regression - Jahresumsatz

Dependent variable: Jahresumsatz (Mio Euro)

Independent variables:

Verkaufsfläche (1000 qm)

Passantenfrequenz (1000/Tag)

		<i>Standard</i>	<i>T</i>	
<i>Parameter</i>	<i>Estimate</i>	<i>Error</i>	<i>Statistic</i>	<i>P-Value</i>
CONSTANT	-0,831946	0,412663	-2,01604	0,0746
Verkaufsfläche	4,74295	0,226498	20,9403	0,0000
Passantenfrequenz	0,174988	0,0448676	3,90009	0,0036

Analysis of Variance

<i>Source</i>	<i>Sum of Squares</i>	<i>Df</i>	<i>Mean Square</i>	<i>F-Ratio</i>	<i>P-Value</i>
Model	52,8765	2	26,4383	384,18	0,0000
Residual	0,619359	9	0,0688177		
Total (Corr.)	53,4959	11			

R-squared = 98,8422 percent

Standard Error of Est. = 0,262331



Streuungszerlegung und Bestimmtheitsmaß

- ▶ Wie im Fall der einfachen linearen Regression gilt für den parameterlinearen Ansatz die **Quadratsummenzerlegung (Streuungszerlegung)** $SST = SSE + SSR$ (bei Schätzung der Regressionskoeffizienten mit der Methode der kleinsten Quadrate).

- ▶ Dabei sind wieder

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2, \text{ die Totalvariabilität (Totalvarianz);}$$

$$SSE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, \text{ die „erklärte“ Variabilität (erklärte Varianz);}$$

$$SSR = \sum_{i=1}^n (y_i - \hat{y}_i)^2, \text{ die Restvariabilität (Restvarianz).}$$

- ▶ Das **Bestimmtheitsmaß** ist

$$B = \frac{SSE}{SST} = 1 - \frac{SSR}{SST} = r_{Y|(f_1(\underline{X}), \dots, f_r(\underline{X}))}^2.$$



Streuungszerlegung und Bestimmtheitsmaß im Beispiel 7.2 in Statgraphics

Analysis of Variance

<i>Source</i>	<i>Sum of Squares</i>	<i>Df</i>	<i>Mean Square</i>	<i>F-Ratio</i>	<i>P-Value</i>
Model	52,8765	2	26,4383	384,18	0,0000
Residual	0,619359	9	0,0688177		
Total (Corr.)	53,4959	11			

R-squared = 98,8422 percent

Standard Error of Est. = 0,262331

- ▶ $SSE = 52.8765$
- ▶ $SSR = 0.619359$
- ▶ $SST = 53.4959$
- ▶

$$B = \frac{SSE}{SST} = \frac{52.8765}{53.4959} = 0.9884215$$



Schätzung der Fehlervarianz

- ▶ Eine konstante Varianz der zufälligen Fehler ε_i (und damit der Zufallsgrößen $Y(\underline{x}_i)$) kann analog zum Fall der einfachen linearen Regression durch

$$\hat{\sigma}^2 = s_{Rest}^2 = \frac{SSR}{n - r}$$

geschätzt werden. Der Nenner $n - r$ ist durch die Schätzung von r Parametern bedingt.

- ▶ Im Beispiel 7.2:

$$\begin{aligned}\hat{\sigma}^2 &= s_{Rest}^2 = \frac{SSR}{n - r} = \frac{0.619359}{12 - 3} = 0.0688177 \\ \hat{\sigma} &= \sqrt{0.0688177} = 0.262331\end{aligned}$$

- ▶ Statgraphics im Beispiel 7.2:

Standard Error of Est. = 0,262331



Konfidenzschätzungen I

- ▶ Für die folgenden Aussagen zu Konfidenzschätzungen und Tests setzen wir wieder voraus, dass die zufälligen Fehler ε_i unabhängige normalverteilte Zufallsgrößen mit Erwartungswert 0 und konstanter Varianz σ^2 sind.
- ▶ Mit m_i wird das i -te Diagonalelement der Matrix $(F^T F)^{-1}$ bezeichnet.
- ▶ Die Schätzung für die Varianz von \hat{a}_i ist damit: $s_{a_i}^2 = s_{Rest}^2 m_i$.
- ▶ Konfidenzintervall zum Niveau $1 - \alpha$ für die Komponente a_i von \underline{a} :
 $I = [\hat{a}_i - s_{a_i} t_{n-r; 1-\alpha/2} ; \hat{a}_i + s_{a_i} t_{n-r; 1-\alpha/2}]$.
- ▶ Konfidenzintervall zum Niveau 0.95 für a_3 (vgl. Folie 15):
 $\hat{a}_3 = 0.175 ; \quad s_{a_3} = 0.0448676 ;$
 $t_{9; 0.975} = 2.26 \quad \Rightarrow \quad I = [0.074; 0.276]$.



Konfidenzschätzungen II

- ▶ Ein Konfidenzintervall für die Fehlervarianz σ^2 ist

$$\left[\frac{(n-r)\hat{\sigma}^2}{\chi_{n-r;1-\alpha/2}^2} ; \frac{(n-r)\hat{\sigma}^2}{\chi_{n-r;\alpha/2}^2} \right] = \left[\frac{\text{SSR}}{\chi_{n-r;1-\alpha/2}^2} ; \frac{\text{SSR}}{\chi_{n-r;\alpha/2}^2} \right].$$

- ▶ Konfidenzintervall zum Niveau $1 - \alpha$ für die Regressionsfunktion $\underline{f}(\underline{x})^T \underline{a}$:

$$I = \left[\underline{f}(\underline{x})^T \hat{\underline{a}} - t_{n-r;1-\alpha/2} \sqrt{s_{\text{Rest}}^2 \underline{f}(\underline{x})^T (F^T F)^{-1} \underline{f}(\underline{x})} ; \right. \\ \left. \underline{f}(\underline{x})^T \hat{\underline{a}} + t_{n-r;1-\alpha/2} \sqrt{s_{\text{Rest}}^2 \underline{f}(\underline{x})^T (F^T F)^{-1} \underline{f}(\underline{x})} \right].$$

- ▶ Auch Prognoseintervalle können konstruiert werden.

t -Test für einzelne Parameter

▶ **Hypothesen:** $H_0 : a_i = a_i^{(0)}$, $H_A : a_i \neq a_i^{(0)}$.

▶ **Testgröße:** $T = \frac{\hat{a}_i - a_i^{(0)}}{s_{a_i}}$.

Diese Testgröße ist unter H_0 t -verteilt mit $n - r$ Freiheitsgraden.

▶ **Kritischer Bereich zum Niveau α :**

$$K = \{t \in \mathbb{R} : |t| > t_{n-r; 1-\alpha/2}\}.$$

▶ Analog können einseitige Tests durchgeführt werden.

▶ Test mit $H_0 : a_3 = 0$, $H_A : a_3 \neq 0$, $\alpha = 0.05$ im Beispiel 7.2.

$$t = \frac{0.175}{0.0448676} = 3.9 > 2.26 = t_{9; 0.975}$$

⇒ H_0 wird abgelehnt, d.h. der Koeffizient a_3 (der die Abhängigkeit des Umsatzes von der Passantenfrequenz beschreibt) ist signifikant verschieden von 0.



F-Test für das Modell (Varianzanalyse)

- ▶ Wir setzen voraus, dass $f_1(x) = 1$ gilt, d.h. a_1 die Konstante im Modell ist.
- ▶ **Hypothesen:** $H_0 : a_2 = \dots = a_r = 0$, $H_A : a_i \neq 0$ für ein $i > 1$.
- ▶ **Testgröße:** $T = \frac{\text{MSE}}{\text{MSR}}$ mit $\text{MSE} = \frac{\text{SSE}}{r-1}$, $\text{MSR} = \frac{\text{SSR}}{n-r}$.

Diese Testgröße ist unter H_0 F -verteilt mit $(r-1; n-r)$ Freiheitsgraden.

- ▶ **Kritischer Bereich zum Niveau α :**
 $K = \{t \in \mathbb{R} : t > F_{r-1; n-r; 1-\alpha}\}$.



t - und F -Tests in Statgraphics im Beispiel 7.2

t -Test für a_3 hervorgehoben:

		<i>Standard</i>	<i>T</i>	
<i>Parameter</i>	<i>Estimate</i>	<i>Error</i>	<i>Statistic</i>	<i>P-Value</i>
CONSTANT	-0,831946	0,412663	-2,01604	0,0746
Verkaufsfläche	4,74295	0,226498	20,9403	0,0000
Passantenfrequenz	0,174988	0,0448676	3,90009	0,0036

F -Test:

Analysis of Variance

<i>Source</i>	<i>Sum of Squares</i>	<i>Df</i>	<i>Mean Square</i>	<i>F-Ratio</i>	<i>P-Value</i>
Model	52,8765	2	26,4383	384,18	0,0000
Residual	0,619359	9	0,0688177		
Total (Corr.)	53,4959	11			