



Methodik der multiplen linearen Regression

Sibel Aydemir

**Statistisches Amt, Direktorium
Landeshauptstadt Munchen**

Regressionsanalyse: Schritt fur Schritt

- Schritt 1 Modellbildung
Auswahl der erklärenden Variablen
- Schritt 2 Schatzung der Regressionsgerade
- Schritt 3 Wie gut fittet die Regressionsgerade?
- Schritt 4 Ist das Gesamtmodell brauchbar?
- Schritt 5 Ist der Einfluss der erklärenden Variablen
statistisch signifikant?
- Schritt 6 Welche Variablen sind zur Erklärung der
Zielvariable tatsachlich erforderlich?
Variablenselektion



Schritt 1: Modellbildung

Schritt 1: Modellbildung

SB-Nr.	Wahlbeteiligung gesamt	Alleinerziehende in %	HH-Groe	RK %	EV %	Sonst. %	Arbeitslose in % 18-65	Auslander in %	Eingebur- gerte in %	Aus- siedler in %	mit Migr.- hint. in %	ohne Migr.- hint. in %	SGB II Bedarfs- gemein- sch. in % der HH	SGB II Per- sonen in % der Einw.
111	66,9%	2,6%	1,73	44,1%	15,6%	40,2%	4,7%	17,2%	5,2%	10,3%	32,7%	67,3%	3,6%	3,3%
121	71,6%	2,8%	1,57	42,9%	14,8%	42,3%	4,7%	19,8%	5,4%	5,8%	31,0%	69,0%	6,1%	5,3%
122	69,0%	2,0%	1,49	49,0%	11,6%	39,4%	2,1%	22,1%	5,9%	4,7%	32,7%	67,3%	3,2%	2,6%
131	78,8%	3,6%	2,14	48,5%	17,6%	33,9%	3,5%	18,6%	4,0%	11,3%	33,9%	66,1%	9,9%	6,3%
132	52,7%	4,2%	2,28	50,1%	17,3%	32,6%	4,6%	20,2%	5,8%	11,3%	37,3%	62,7%	9,9%	6,5%
141	80,2%	2,0%	2,11	59,0%	14,8%	26,2%	2,1%	2,8%	3,7%	4,7%	11,2%	88,8%	2,3%	1,3%

- Spezifikation der abhangigen Variable y
- Auswahl der erklarenden Variablen x_1, \dots, x_n aufgrund theoretischer Voruberlegungen
- Einschrankungen in der Auswahl, da nicht immer alle potentiellen Einflussvariablen verfugbar bzw. messbar sind
- **Beachte:** Ergebnis der Regressionsanalyse hangt von der Auswahl der unabhangigen Variablen ab

Beispiel: Variable „mit Migrationshintergrund“ wird alternativ zu den 3 Variablen „Auslander“, „Eingeburgerte“ und „Aussiedler“ verwendet

	Modell 1:	Modell 2:
abhangige Variable	Wahlbeteiligung	Wahlbeteiligung
unabhangige Variablen	Alleinerziehende Haushaltsgroe Romisch-katholisch Evangelisch Arbeitslose SGBII Bedarfsgemeinschaften SGBII Personen Mit Migrationshintergrund	Alleinerziehende Haushaltsgroe Romisch-katholisch Evangelisch Arbeitslose SGBII Bedarfsgemeinschaften SGBII Personen Auslander Eingeburgerte Aussiedler
signifikante Variablen	SGBII Bedarfsgemeinschaften Mit Migrationshintergrund Romisch-katholisch	SGBII Bedarfsgemeinschaften Auslander Aussiedler Evangelisch

Schritt 2:

Schatzung der Regressionsgerade

Schatzung der Regressionsgerade

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n + e$$

- Schatzung der Regressionskoeffizienten $\beta_0, \beta_1, \dots, \beta_n$ ber KQ-Methode
- Interpretation des Regressionskoeffizienten β_i im multiplen Regressionsmodell: β_i gibt den Einfluss der Variablen x_i bei Konstanthaltung des Einflusses aller anderen erklarenden Variablen wieder

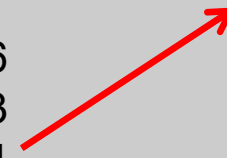
Beispiel: Interpretation der Regressionskoeffizienten im multiplen linearen Modell

	Modell 1
abhangige Variable	Wahlbeteiligung
unabhangige Variablen	Alleinerziehende Haushaltsgroe Romisch-katholisch Evangelisch Arbeitslose SGBII Bedarfsgemeinschaften SGBII Personen Mit Migrationshintergrund
signifikante Variablen	Romisch-katholisch SGBII-Bedarfsgemeinschaften Mit Migrationshintergrund

β -Koeffizient

Erhoung des Anteils der Personen mit Migrationshintergrund um 1% fuhrt zu einem Ruckgang der Wahlbeteiligung um 0,6% bei Konstanthaltung der Anteile der Variablen Rom.-kath. und SGBII-BG

- 0,26
- 0,83
- 0,64



Schritt 3:

Wie gut fittet die Regressionsgerade?

Mae fur die Modellgute

- **Bestimmtheitsma R^2**
- **Akaikes Informationskriterium AIC**
- **Schwarz'sche Bayes Kriterium (SBC oder BIC)**
- **Mallows C_p**

Das Bestimmtheitsma R^2

$$\text{Bestimmtheitsma } R^2 = \frac{\text{erklarte Varianz}}{\text{Gesamtvarianz}}$$

- R^2 = Anteil der durch das Regressionsmodell erklarten Varianz an der Gesamtvarianz
- R^2 nimmt Werte zwischen 0 und 1 an
- Je naher R^2 an 1 liegt, desto besser „passt“ die Regressionsgerade

Beispiel: $R^2 = 0,7$

D.h. 70% der Variation der abhangigen Variable y sind auf die erklarende Variable x zurckzufhren.

Die nicht erklarte Varianz von 30% resultiert u.a. durch nicht bercksichtigte Variablen.

Zur Interpretation des Bestimmtheitsmaes R^2

- Hohes R^2 sagt nichts ber die Erklarungskraft der einzelnen Koeffizienten aus
- Hhe des R^2 hangt stark von den Daten ab
- R^2 wachst mit zunehmender Anzahl von erklarenden Variablen
- Abhilfe: **Korrigiertes Bestimmtheitsma** (R^2 adjusted)

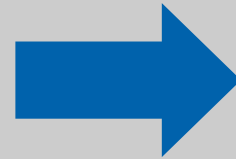
Beispiel

Abhangigkeit des Bestimmtheitsmaes von der Anzahl der Variablen

abhangige Variable	Modell 1: Wahlbeteiligung	R^2	R^2 korrigiert
unabhangige Variablen	Alleinerziehende Haushaltsgroe Romisch-katholisch Evangelisch Arbeitslose SGBII Bedarfsgemeinschaften SGBII Personen Mit Migrationshintergrund	0,880	0,871
signifikante Variablen	SGBII Bedarfsgemeinschaften Mit Migrationshintergrund Romisch-katholisch	0,875	0,872

Hohes Bestimmtheitsma R^2 nur Zufall?

**Deskriptive
Statistik**



**Induktive
Statistik**

Schritt 4:

**Ist das Gesamtmodell
brauchbar?**

Ist das Gesamtmodell brauchbar?

Spezifizierte Regressionsgleichung ist unbrauchbar



Gesamttest auf Signifikanz

Prufverfahren: F-Test

- $H_0 : \beta_i = 0$, fur alle $i=1, \dots, n$
- D.h. **keine** der berucksichtigten unabhangigen Variablen x_1, \dots, x_n besitzt einen Einfluss auf die abhangige Variable y
- Wird H_0 abgelehnt, so hat mindestens eine der erklarenden Variablen x_1, \dots, x_n einen Einfluss auf y
- Faustregel: H_0 wird abgelehnt, falls F -Wert > 10

Schritt 5:

**Ist der Einfluss der
erklarenden Variablen
statistisch signifikant?**

Welche Variablen sind statistisch signifikant?

Prufung der Regressionskoeffizienten β_1, \dots, β_n auf statistische Signifikanz



Prufverfahren: t-Test

- $H_0 : \beta_i = 0, i=1, \dots, n$
- D.h. die Variable x_i besitzt keinen Einfluss auf die abhangige Variable y
- Faustregel: H_0 wird abgelehnt, falls t -Wert $> |2|$
- Signifikanzniveau (Irrtumswahrscheinlichkeit) < 0.05

Nicht-signifikante Variablen

- **Vorsicht bei der Interpretation:** Ist eine erklärende Variable nicht signifikant, heißt das nicht unbedingt, dass sie keinen Einfluss auf die abhängige Variable y besitzt
- Korrelieren zwei (oder mehrere) unabhängige Variablen, so kann es sein, dass in der multiplen Regression eine Variable sich nicht durchsetzen kann, da sie keine zusätzliche Information zur Regression beiträgt
- Variablen, die in der multiplen Regression nicht-signifikant sind, können in der einfachen linearen Regression durchaus einen signifikanten Einfluss zeigen

Beispiel

Multiple lineare Regression vs. einfache lineare Regression

	Modell 1:	Multiple lin. Regression	Einfache lin. Regression
abhangige Variable	Wahlbeteiligung		
unabhangige Variablen	Alleinerziehende	--	signifikant
	Haushaltsgroe	--	--
	Romisch-katholisch	signifikant	signifikant
	Evangelisch	--	--
	Arbeitslose	--	signifikant
	SGBII Bedarfsgemeinschaften	signifikant	signifikant
	SGBII Personen	--	signifikant
	Mit Migrationshintergrund	signifikant	signifikant

„Wichtigkeit“ einer erklärenden Variable

- **Vorsicht:** Das Signifikanzniveau ist nicht ausreichend, um Aussagen über die „Wichtigkeit“ einer erklärenden Variable machen zu können
- Die „Wichtigkeit“ einer Variable lässt sich an der (standardisierten) Koeffizientenschätzung erkennen
- Denn: Eine im Vergleich „weniger signifikante“ Variable, kann evtl. mehr zur Erklärung/Vorhersage von y beitragen (→ Beispiel)

Beispiel: „Wichtigkeit“ der erklärenden Variablen

abhängige Variable	Modell 2:	p-Wert	β	β stand.
unabhängige Variablen	Wahlbeteiligung Alleinerziehende Haushaltsgröße Römisch-katholisch Evangelisch Arbeitslose SGBII Bedarfsgemeinschaften SGBII Personen Ausländer Eingebürgerte Aussiedler			
signifikante Variablen	SGBII Bedarfsgemeinschaften Ausländer Aussiedler Evangelisch	.029 .007 .000 .001	- 0.62 - 0.28 - 0.75 0.37	- 0.22 - 0.19 - 0.47 0.12

P(Evangelisch) < P(SGBII-BG)

Aber: Variable SGBII-BG hat einen nahezu doppelt so großen Einfluss auf die Wahlbeteiligung als Variable Evangelisch

Angenommen, die Regressionsschatzung zeigt:

Einige Variablen besitzen einen statistisch signifikanten Einfluss auf y , andere Variablen besitzen keinen Einfluss.

Wie geht es weiter?

1.Moglichkeit: Schatzung eines Endmodells nur mit den signifikanten erklärenden Variablen

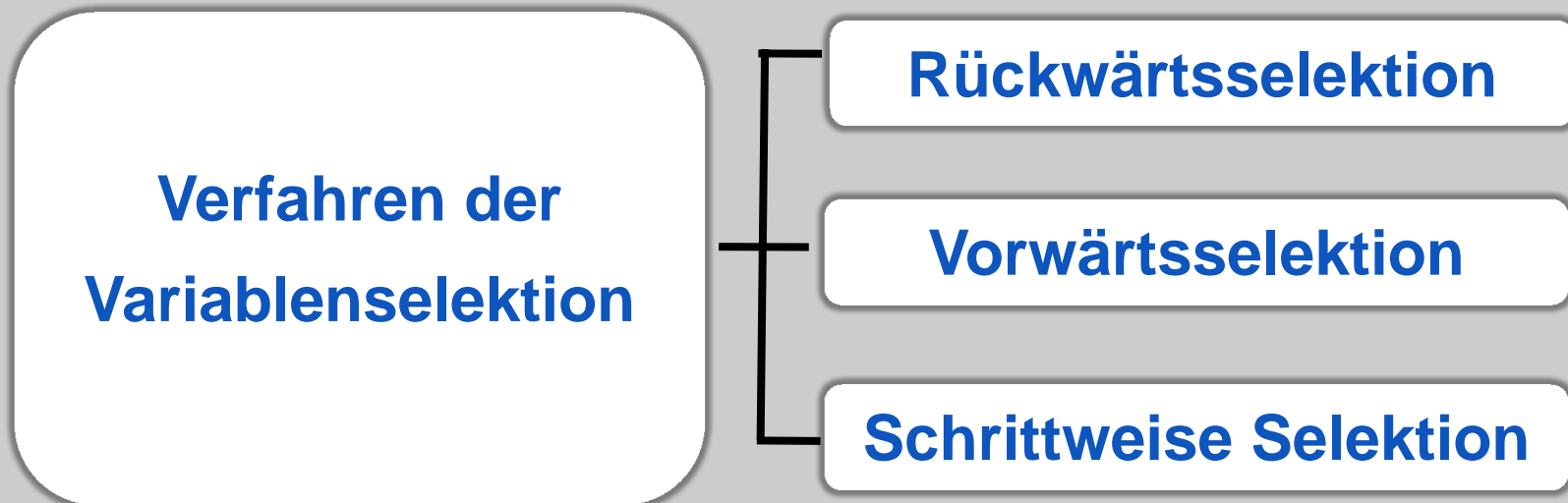
2.Moglichkeit: Schatzung eines Endmodells mittels Variablenselektion

Vorteil: Multikollinearitat wird berucksichtigt

Schritt 6:

**Welche Variablen sind zur
Erklärung der Zielvariable
tatsächlich erforderlich?**

Variablenselektion



Ruckwartsselektion

- **Start:** vollstandiges Modell mit allen unabhangigen Variablen
- Sukzessive werden diejenigen Variablen entfernt, die zum geringsten Ruckgang des Bestimmtheitsmaes R^2 fuhren wurden.
- **Stopp:** Verfahren bricht ab, falls sich beim Entfernen einer (bzw. der nachsten) Variable das Bestimmtheitsma R^2 signifikant verkleinert.

Beispiel: Ruckwartsselektion (RS)

abhangige Variable	Modell 1	p-Wert	RS	p (Endmodell)
	Wahlbeteiligung			
unabhangige Variablen	Alleinerziehende	0,296	3	
	Haushaltsgroe	0,384	2	
	Romisch-katholisch	0,002		0,001
	Evangelisch	0,635	1	
	Arbeitslose	0,171	5	
	SGBII Bedarfsgemeinschaften	0,081		0,000
	SGBII Personen	0,362	4	
	Mit Migrationshintergrund	0,000		0,000
signifikante Variablen	Romisch-katholisch			Romisch-katholisch
	Mit Migrationshintergrund			SGBII Bedarfsgemeinschaften Mit Migrationshintergrund

Vorwartsselektion

- **Start:** Modell ohne unabhangige Variablen, also $y = \beta_0$
- Bestimme diejenige erklarende Variable, die mit y am starksten korreliert ist und berechne das Bestimmtheitsma R^2 .
- Ist R^2 signifikant, wird diese Variable in das Modell aufgenommen.
- In weiteren Schritten werden sukzessive die Variablen ins Modell aufgenommen, die zum groten Anstieg von R^2 fuhren.
- **Stopp:** Verfahren bricht ab, falls sich bei Hinzunahme einer weiteren Variable das Bestimmtheitsma R^2 nicht signifikant vergroert.

Beispiel: Vorwartsselektion (VS)

	Modell 2	R^2	VS
abhangige Variable	Wahlbeteiligung		
unabhangige Variablen	Alleinerziehende		
	Haushaltsgroe		
	Romisch-katholisch		
	Evangelisch	0,877	4
	Arbeitslose		
	SGBII Bedarfsgemeinschaften	0,765	1
	SGBII Personen		
	Auslander	0,866	3
	Eingeburgerte		
	Aussiedler	0,850	2

**Signifikante
Variablen bei VS**

Lineares Regressionsmodell ohne Variablenselektion

Signifikante Variable: **Aussiedler**

Schrittweise Selektion

- **Kombination aus Vorwarts- und Ruckwartsselektion**
- **Es wird eine Vorwartsselektion durchgefuhrt, bei der nach jedem Schritt untersucht wird, ob bei Entfernen einer der bisher aufgenommenen Variablen das Bestimmtheitsma signifikant abnehmen wurde (=Ruckwartsselektion).**



Variablenselektion

**Vorwartsselektion, Ruckwartsselektion
und schrittweise Selektion fuhren (meist)
zum selben Ergebnis.**



Überblick

Regressionsmodelle

Regressionsmodelle

Ziel

**Untersuchung des Zusammenhanges
zwischen einer abhangigen Variable
und mehreren unabhangigen Variablen**

Die bekanntesten Regressionsmodelle

- **Lineare Regression (einfach / multipel)**
Abhangige Variable y = metrisch
- **Logistische Regression**
Abhangige Variable y = binar (0/1-Kodierung)
- **Cox-Regression**
Abhangige Variable y = Zeitdauer
(z.B. Wohndauer, Ehedauer)



**Vielen Dank fur Ihre
Aufmerksamkeit**

Schritt 3: Wie gut fittet die Regressionsgerade?

Bestimmtheitsma R^2

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\text{erklarte Varianz}}{\text{Gesamtvarianz}}$$