

5 Assoziationsmessung in Kontingenztafeln

5.1 Multivariate Merkmale

Die Analyse eindimensionaler Merkmale ist nur der allererste Schritt zur Beschreibung der Daten. Meist ist die Analyse von *Zusammenhängen* zwischen Merkmalen von grösserem inhaltlichen Interesse.

Bsp. 5.1.

Beispiele für typische Fragestellung:

- Beeinflusst das Geschlecht das Erwerbseinkommen?
- Gibt es einen Zusammenhang zwischen Schichtzugehörigkeit (als etwas veralteter, dennoch klassischer soziologischer Begriff) und sozialem Engagement?
- Spielt die Stärke der Kirchenbindung eine Rolle bei der Parteienpräferenz?
- Haben Studierende mit guten Mathematikvorkenntnissen bessere Statistiknoten?

Hierzu werden an jeder Einheit *mehrere* Merkmale erhoben und ihre Ausprägungen auch *gemeinsam* analysiert (z.B. wird das Geschlecht der i -ten Person mit ihrem Einkommen in Beziehung gesetzt).

Hat man z.B. die Merkmale X, Y, Z und analysiert sie gemeinsam, so nennt man das Paar (X, Y) bzw. das Tripel (X, Y, Z) (trivariates) Merkmal. Allgemein spricht man von mehrdimensionalen Merkmalen.

$$(X, Y) : \Omega \longrightarrow (W_X \times W_Y)$$

$$\omega \longmapsto (X(\omega), Y(\omega))$$

Achtung:

- Die später folgenden statistischen Verfahren messen die Stärke von Zusammenhängen, aber erlauben keine Aussagen über Kausalität!
- Ob eine kausale Interpretation des Zusammenhangs zulässig ist, hängt davon ab, wie die Daten erhoben wurden.
- Statistische Zusammenhangsmaße können nicht klären:
 - die Richtung des Zusammenhangs (was ist Ursache, was Wirkung?)
⇒ Längsschnitt-Studie, „cross-lag“ Design
 - ob eine dritte, evtl. unbeobachtete Variable den Zusammenhang verursacht ⇒ Experiment

5.2 Kontingenztafeln und bedingte Verteilungen

Wir betrachten in diesem Kapitel diskrete Merkmale, also typischerweise Merkmale mit wenigen verschiedenen Ausprägungen, die zudem nominal oder ordinal skaliert seien. Man spricht dann bei Zusammenhängen von „*Assoziation*“, im Gegensatz zur „Korrelation“ bei stetigen Merkmalen (vgl. Kap. 6).

5.2.1 Gemeinsame Verteilung, Randverteilung, Kontingenztafel

Betrachtet wird ein zweidimensionales Merkmal (X, Y) bestehend aus den diskreten Merkmalen X und Y und die zugehörige Urliste

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n).$$

Ferner sei wieder angenommen, dass X und Y nur endlich viele (wenige), verschiedene

Werte annehmen können:

$$a_1, \dots, a_i, \dots, a_k \quad \text{bzw.} \quad b_1, \dots, b_j, \dots, b_m$$

Anmerkung: In vielen Büchern (v.a. zur induktiven Statistik) wird statt a_1, \dots, a_k auch x_1, \dots, x_k und analog statt b_1, \dots, b_m auch y_1, \dots, y_m geschrieben. Bei uns sind aber die (x_i, y_i) Werte der Urliste, x_i also der Wert der i -ten Einheit. Daraus ergibt sich zwar die Doppeldeutigkeit der Laufindizes i und j , wir bleiben jedoch bei dieser Notation um Einheitlichkeit mit Fahrmeir et al. (2009) und Jann (2005) herzustellen.

Bsp. 5.2.

(fiktiv):

$$\begin{aligned} \text{Fahrzeugmodell } X &= \begin{cases} 1, & \text{Kompaktklasse} \\ 2, & \text{Mittelklasse} \\ 3, & \text{Oberklasse} \end{cases} \\ \text{aggressives Fahrverhalten } Y &= \begin{cases} 1, & \text{ja} \\ 2, & \text{nein} \end{cases} \end{aligned}$$

Typische Urliste des zweidimensionalen Merkmals (X, Y) :

$(3, 1), (2, 2), (2, 1), (3, 1), (3, 2), (3, 1), (1, 2), (1, 1), (1, 1), (1, 2), (3, 1), (3, 1)$

Einheit	X	Y
1	3	1
2	2	2
3	2	1
4	3	1
5	3	2
6	3	1
7	1	2
8	1	1
9	1	1
10	1	2
11	3	1
12	3	1

Achtung:

- Tupel sind – im Gegensatz zu Mengen – *geordnete* Anordnungen von Zahlen
- Die Tupel sind „gemeinsam indiziert“

Gemeinsame relative und absolute Häufigkeitsverteilung:

$$h_{ij} = h(a_i, b_j), \quad i = 1, \dots, k, \quad j = 1, \dots, m,$$

Anzahl von Beobachtungen mit $x = a_i$ und $y = b_j$.

$$f_{ij} = h_{ij}/n = f(a_i, b_j), \quad i = 1, \dots, k, \quad j = 1, \dots, m,$$

Anteil von Beobachtungen mit $x = a_i$ und $y = b_j$.

Man nennt (h_{ij}) , bzw. (f_{ij}) , $i = 1, \dots, k, j = 1, \dots, m$ die *gemeinsame Verteilung* von (X, Y) in absoluten bzw. relativen Häufigkeiten.

Kontingenztabelle / Kontingenztabelle / Kreuztabelle: Darstellung der Häufigkeiten in Form einer $(k \times m)$ -dimensionalen Häufigkeitstabelle

	b_1	\cdots	b_j	\cdots	b_m	
a_1	h_{11}	\cdots	h_{1j}	\cdots	h_{1m}	$h_{1\bullet}$
a_2	h_{21}	\cdots	h_{2j}	\cdots	h_{2m}	$h_{2\bullet}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
a_i	h_{i1}	\cdots	h_{ij}	\cdots	h_{im}	$h_{i\bullet}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
a_k	h_{k1}	\cdots	h_{kj}	\cdots	h_{km}	$h_{k\bullet}$
	$h_{\bullet 1}$	\cdots	$h_{\bullet j}$	\cdots	$h_{\bullet m}$	n

Der Punkt steht für Summation über den entsprechenden Index. Es ergeben sich die *Randverteilungen* $h_{i\bullet} = \sum_{j=1}^m h_{ij} = h_{i1} + \dots + h_{im} = h_X(a_i)$, $i = 1, \dots, k$, für X und $h_{\bullet j} = \sum_{i=1}^k h_{ij} = h_{1j} + \dots + h_{kj} = h_Y(b_j)$, $j = 1, \dots, m$, für Y . Es gilt also:

- $h_{i\bullet}$ ist die absolute Häufigkeit von a_i ,
- $h_{\bullet j}$ ist die absolute Häufigkeit von b_j .

$$h_{i\bullet} = \sum_{j=1}^m h_{ij},$$
$$h_{\bullet j} = \sum_{i=1}^k h_{ij}$$

Also ist $h_{i\bullet}$ die i -te Zeilensumme, $h_{\bullet j}$ die j -te Spaltensumme (daher der Name Randhäufigkeiten).

Kontingenztafel der relativen Häufigkeitsverteilung:

	b_1	\cdots	b_j	\cdots	b_m	
a_1	f_{11}	\cdots	f_{1j}	\cdots	f_{1m}	$f_{1\bullet}$
a_2	f_{21}	\cdots	f_{2j}	\cdots	f_{2m}	$f_{2\bullet}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
a_i	f_{i1}	\cdots	f_{ij}	\cdots	f_{im}	$f_{i\bullet}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
a_k	f_{k1}	\cdots	f_{kj}	\cdots	f_{km}	$f_{k\bullet}$
	$f_{\bullet 1}$	\cdots	$f_{\bullet j}$	\cdots	$f_{\bullet m}$	1

mit der relativen Häufigkeiten $f_{ij} = \frac{h_{ij}}{n}$ und den *Randverteilungen*

$$f_{i\bullet} = \frac{h_{i\bullet}}{n} = f_{i1} + \dots + f_{im} = f_X(a_i), \quad i = 1, \dots, k, \quad \text{für } X$$

und

$$f_{\bullet j} = \frac{h_{\bullet j}}{n} = f_{1j} + \dots + f_{kj} = f_Y(b_j), \quad j = 1, \dots, m, \quad \text{für } Y.$$

Bsp. 5.3. *Aggressives Fahren und Fahrzeugmodell (fiktiv)*

Beachte: Aus der gemeinsamen Verteilung kann man die Randverteilungen berechnen (aber nicht umgekehrt, s.u., Kap. 5.2.2).

Bsp. 5.4. *Wahlabsicht und Bildungsabschluss (ALLBUS 2010: V327, V747)*

	CDU-CSU	SPD	FDP	GRÜNE	LINKE	NPD	ANDERE	NICHTW.	
OHNE ABSCHLUSS	4	5	1	2	1	0	0	6	19
VOLKS-,HAUPTSCHULE	203	203	33	77	75	10	12	120	733
MITTLERE REIFE	180	169	48	130	88	13	18	102	748
FACHHOCHSCHULREIFE	26	34	12	30	10	1	5	11	129
HOCHSCHULREIFE	124	120	34	173	41	3	15	31	541
ANDERER ABSCHLUSS	3	1	0	1	0	0	0	3	8
NOCH SCHUELER	4	2	3	3	0	0	0	0	12
	544	534	131	416	215	27	50	273	2190

Bem. 5.5.

Ist $k = m = 2$ so spricht man von einer Vierfeldertafel. Dabei vereinfachen sich die Tabellen wesentlich, mit der Angabe der Häufigkeit in einer Zelle sind bei gegebenen Randhäufigkeiten auch die Häufigkeiten in den anderen Zellen bestimmt.

	1	2	
1	h_{11}	h_{12}	$h_{1\bullet}$
2	h_{21}	h_{22}	$h_{2\bullet}$
	$h_{\bullet 1}$	$h_{\bullet 2}$	n

z.B. gegeben h_{11}

$\Rightarrow h_{12} = h_{1\bullet} - h_{11}$ etc.

Unabhängige und abhängige Variable:

Hat man eine Vermutung über die Richtung einer potentiellen Wirkung, so bezeichnet man die Variablen entsprechend als *unabhängige* (wirkende, erklärende) und *abhängige* (bewirkte, erklärte) Variable.

In der Statistik ist es üblich, die unabhängige Variable mit X zu bezeichnen und die abhängige Variable mit Y . (Geht man von einem deterministischen Zusammenhang aus, so ist dann, wie gewohnt, Y eine Funktion von X .)

z.B:

	unabhängige (X)		abhängige (Y)
möglicherweise:	Automodell	→	aggressives Fahrverhalten
<u>eindeutig:</u>	Geschlecht	→	Einkommen
allgemein:	unabhängige	→	abhängige Variable

Damit werden die Häufigkeitsverteilungen für feste Werte der unabhängigen Variablen in den Zeilen angegeben.

$\uparrow \leftarrow Y \rightarrow$

X

\downarrow

Vorsicht: In einigen sozialwissenschaftlichen Büchern wird auf eine andere Konvention zurückgegriffen. Dort wird die unabhängige Variable in den Spalten und die abhängige in den Zeilen abgetragen.

5.2.2 Ökologischer Fehlschluss

Es gibt sehr viele gemeinsame Verteilungen, die zu denselben Randhäufigkeiten passen.
Im Beispiel oben passen u.a.:

Man sieht also, wie wichtig es zur Feststellung potentieller Zusammenhänge ist, die *gemeinsame* Verteilung h_{ij} zu kennen, also tatsächlich die Paare (x_i, y_i) zu betrachten.

Der *unzulässige* Schluss

- von (Eigenschaften der) Randverteilungen auf Eigenschaften der gemeinsamen Verteilung,
- also von zwei univariaten Analysen auf eine bivariate Aussage,
- von der Kollektiv- auf die Individualebene,

heißt *ökologischer Fehlschluss*.

Kommen zwei Eigenschaften (verschiedene Merkmale) häufig vor, heißt dies nicht notwendig, dass sie gemeinsam häufig vorkommen.

5.2.3 Grafische Darstellung der gemeinsamen Verteilung

Verschiedene Darstellungsarten:

- Mosaikplots: Darstellung der gemeinsamen Häufigkeiten h_{ij} als *flächentreue* Kachelung
- 3D-Säulendiagramm der gemeinsamen Häufigkeiten h_{ij}
- „normale“ Säulendiagramme nach einer Variable aufgespalten, d.h. für jeden Wert a_i von X werden jeweils die Häufigkeiten h_{ij} bzw. f_{ij} aufgetragen.

Mosaikplots

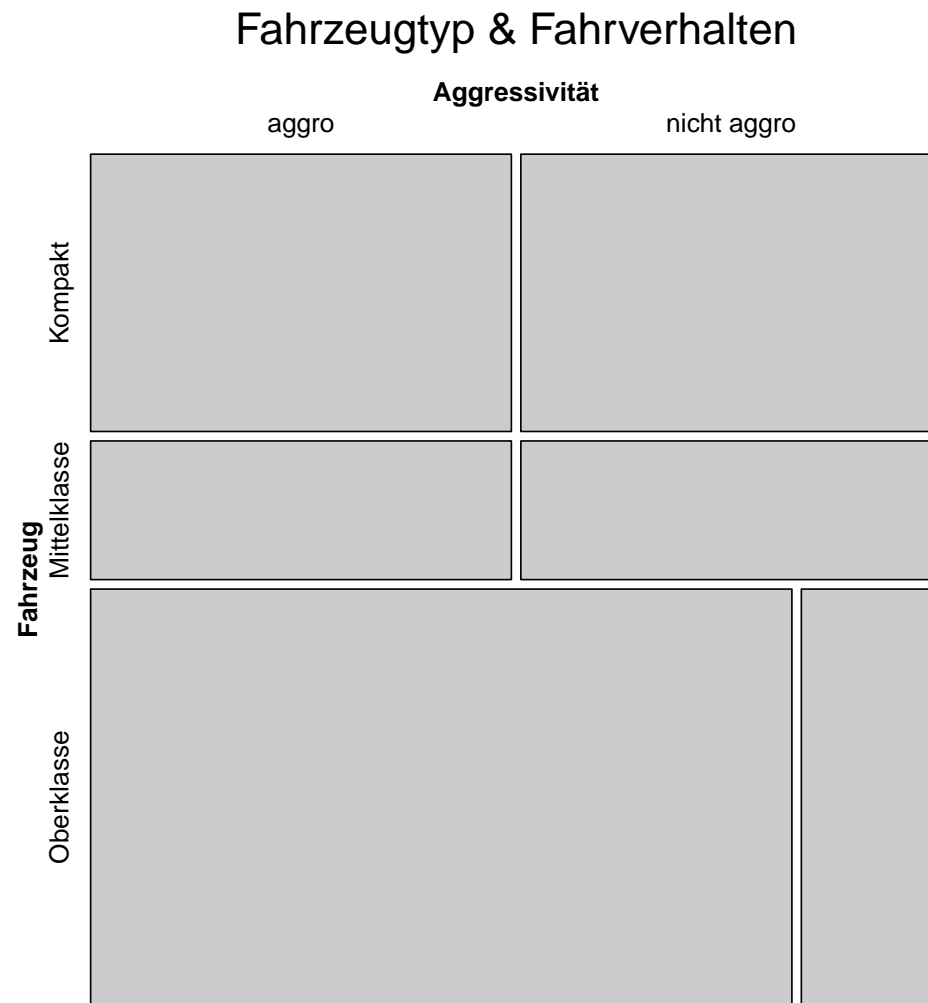
Grafische Darstellung der gemeinsamen Häufigkeiten zweier diskreter Merkmale.

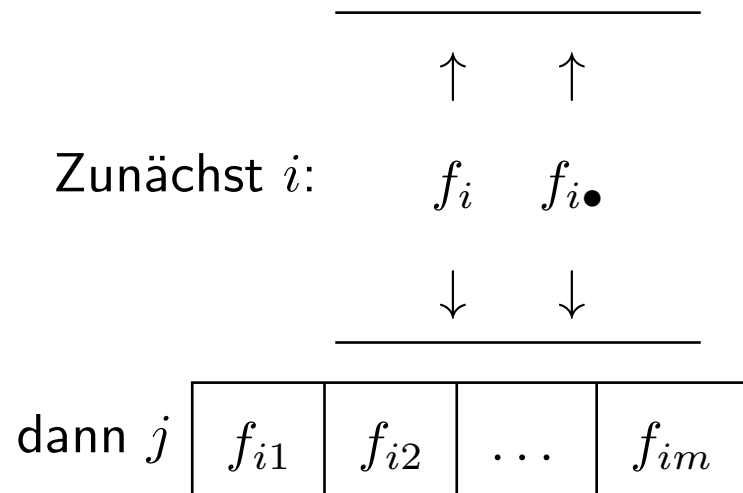
Idee: Teile Quadrat auf in Rechtecke, deren Flächeninhalte f_{ij} entsprechen.

Vorgehen:

1. Teile Einheitsquadrat auf in horizontale Streifen deren Höhe proportional zu $f_{i\bullet}$ ist.
2. Teile horizontale Streifen in Rechtecke deren Breite für jedes i proportional zu f_{ij} ist.

Bsp. 5.6. Aggressivität & Fahrverhalten





3D-Säulendiagramm & “Heatmaps” Grafische Darstellung der gemeinsamen Häufigkeiten zweier diskreter Merkmale, auch erweiterbar auf (quasi)-stetige.

Idee: Benutze Merkmalsausprägungen von X , Y als 2-D Koordinatensystem, gemeinsame Häufigkeiten für jede Kombination (a_i, b_j) werden graphisch über Höhe oder über Farbe dargestellt.

Bsp. 5.7.

Habilitationen nach Geschlecht und Fach (nach Fahrmeir et al., 2009).

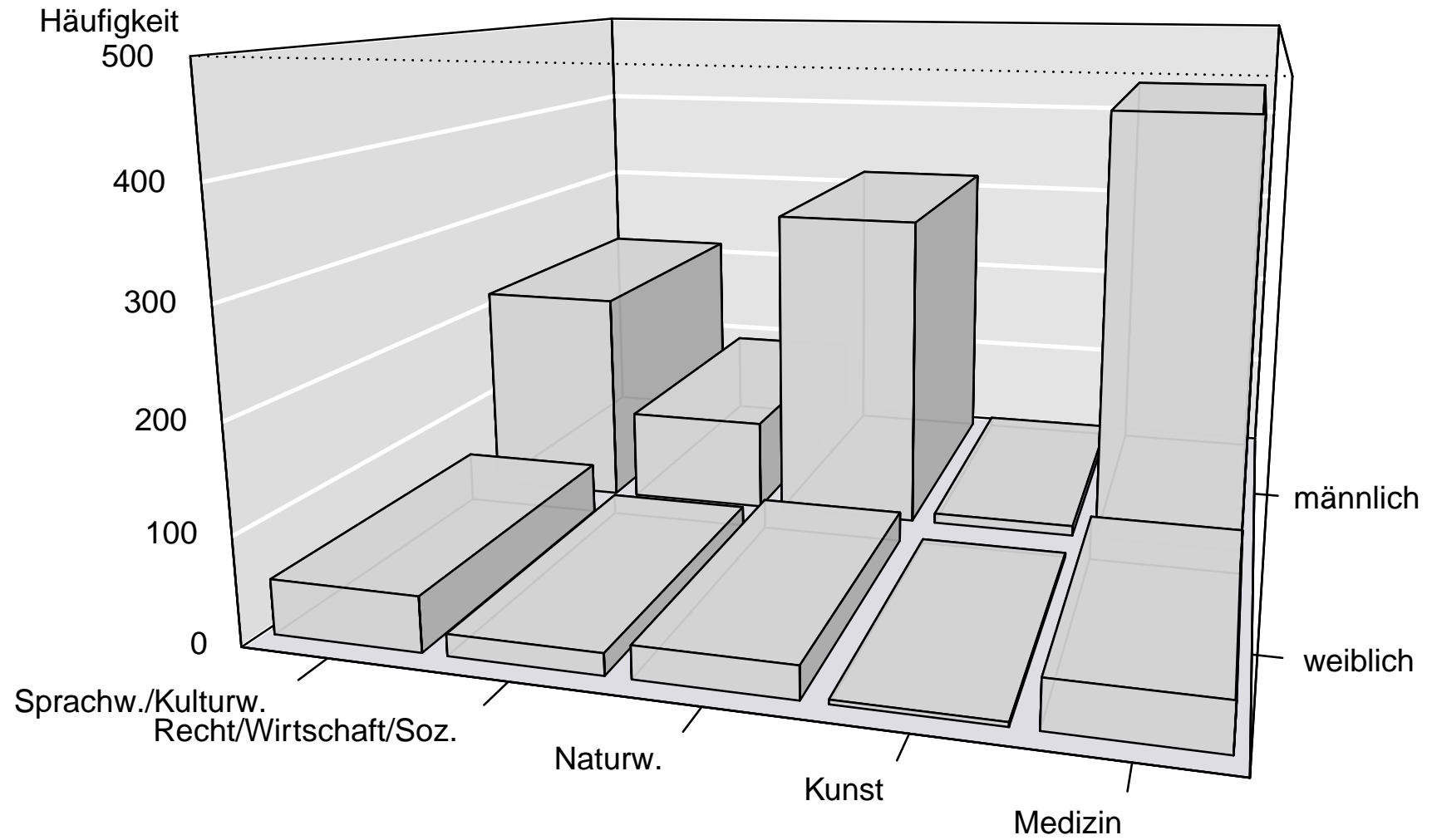
Grundgesamtheit: alle Habilitationen 1993

Geschlecht: X

Fächergruppe: Y

		Sprachw. Kulturw.	Rechtsw. Wirts., Soz.	Naturw.	Kunst	Medizin	
		1	2	3	4	5	
weiblich	1	51	20	30	4	44	149
männlich	2	216	92	316	10	433	1067
		267	112	346	14	477	1216

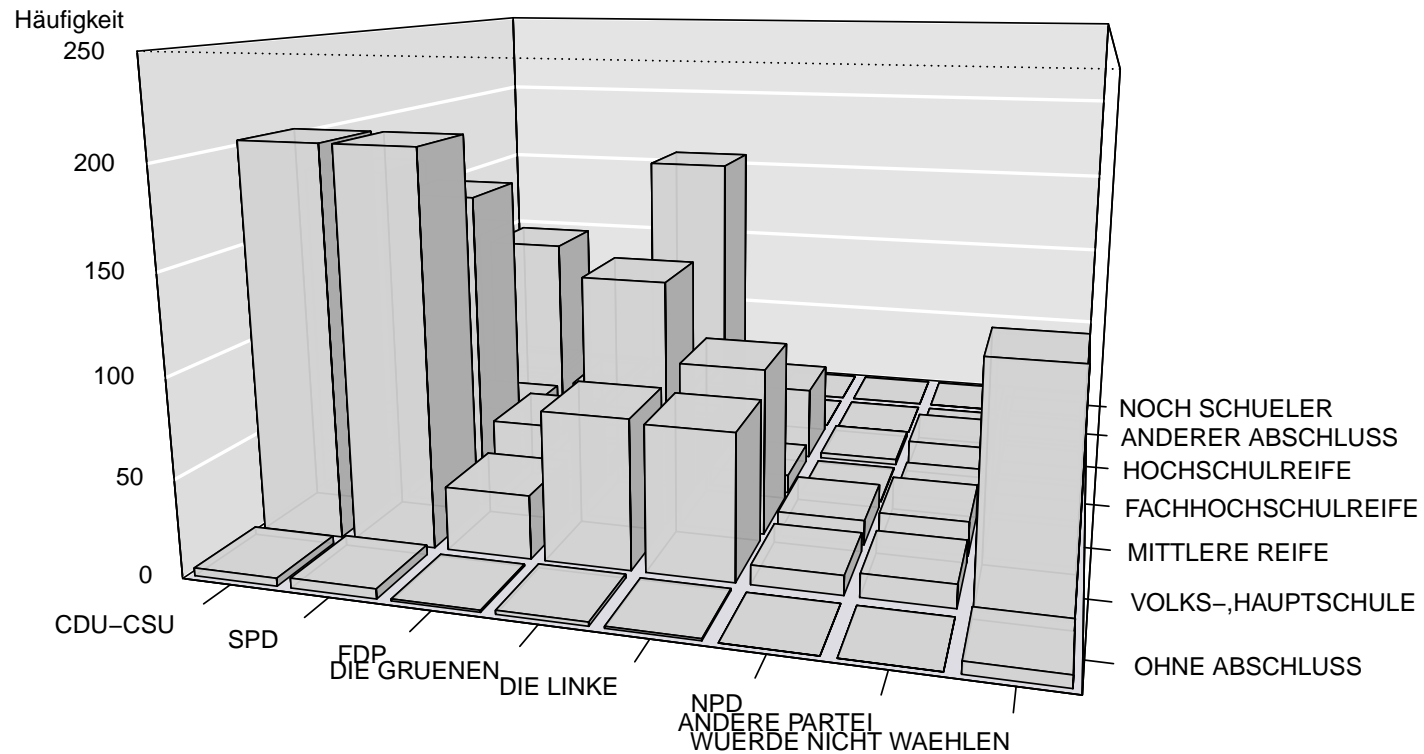
3D-Säulendiagramm:



Pro: relativ intuitiv

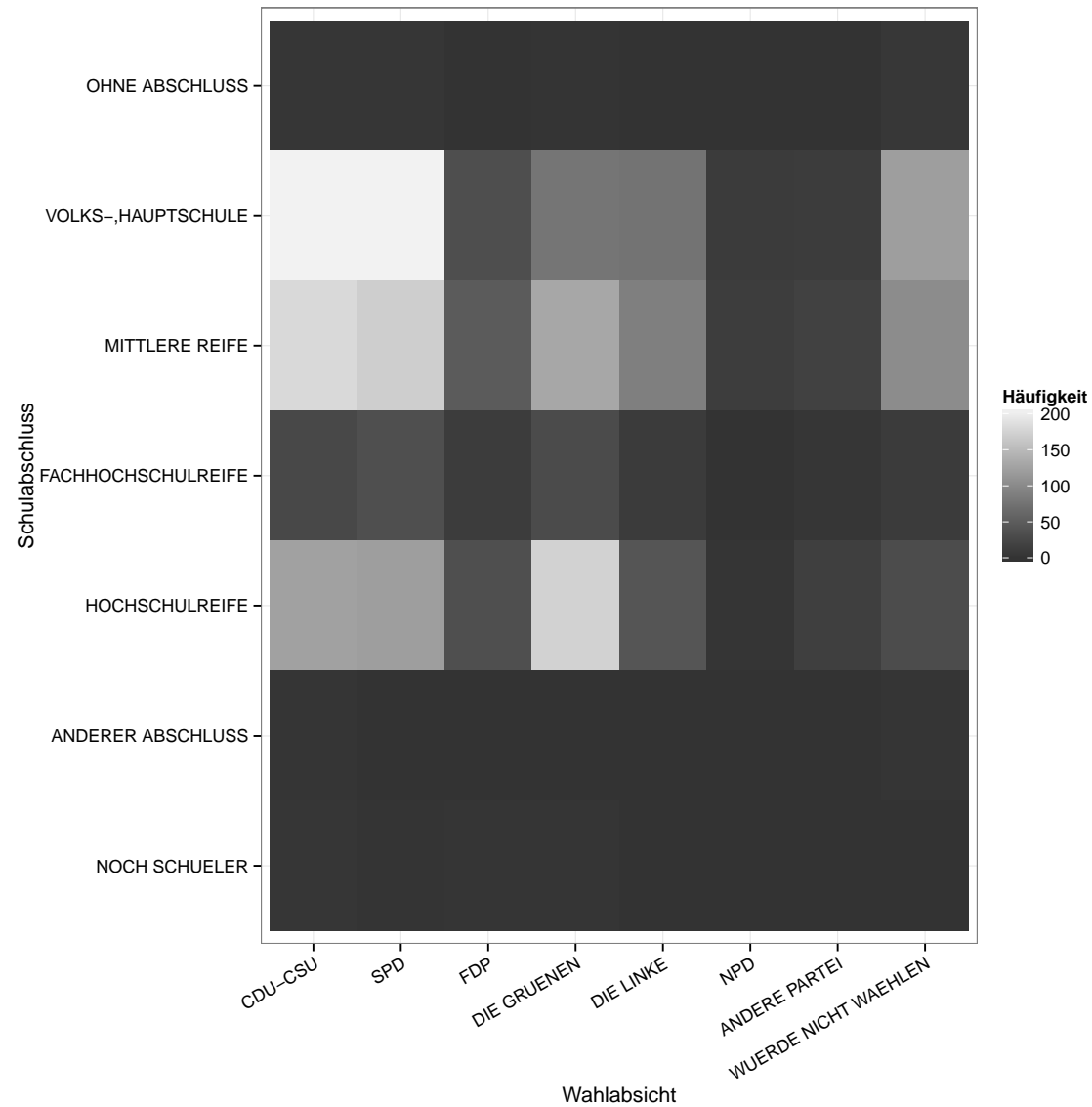
Kontra: 3D- Darstellung kann Muster verdecken (niedrige Balken hinter hohen Balken)

Bsp. 5.8. *Wahlabsicht und Bildungsabschluss*

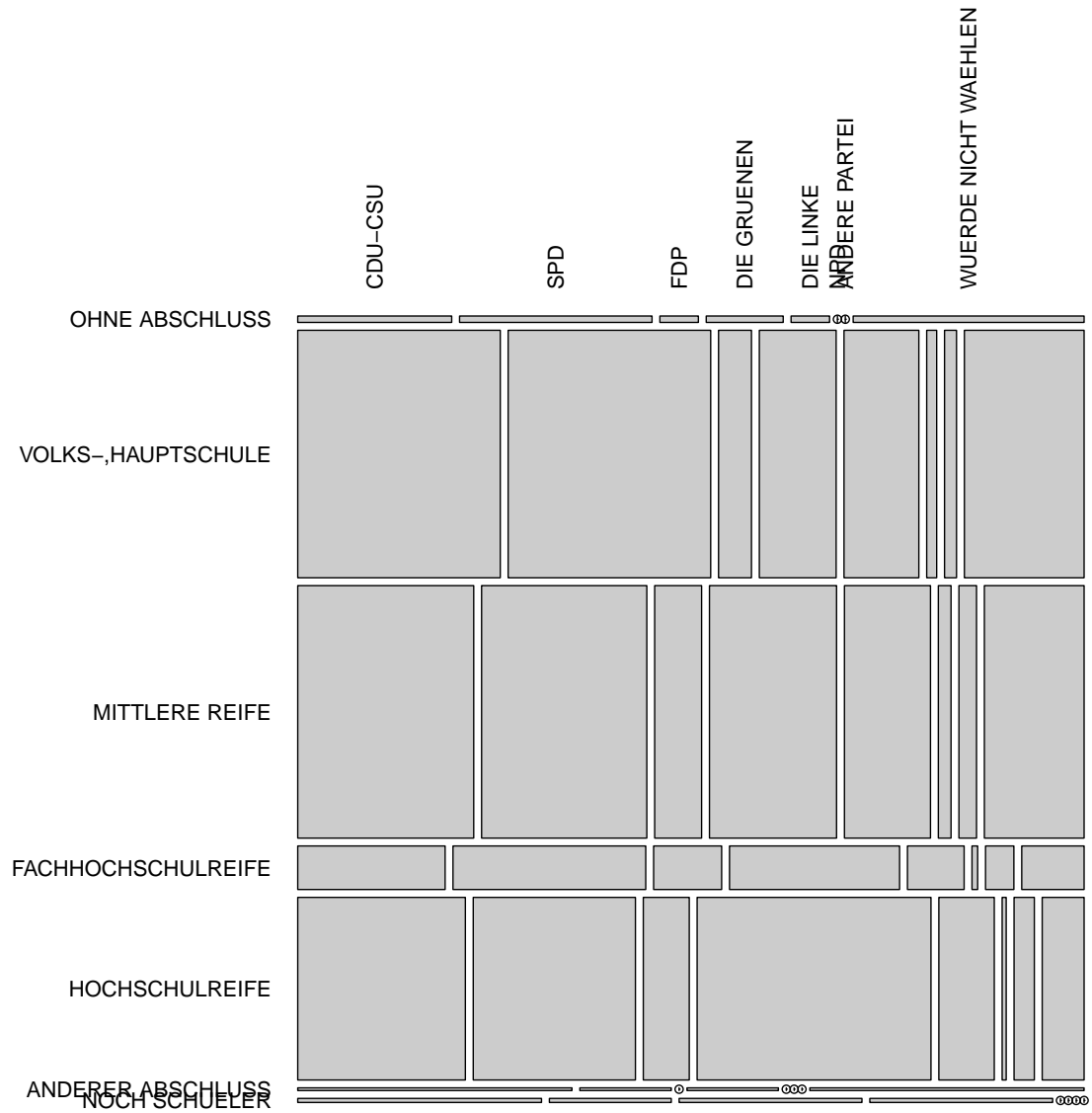


3D-Säulendiagramm hier ungünstig.

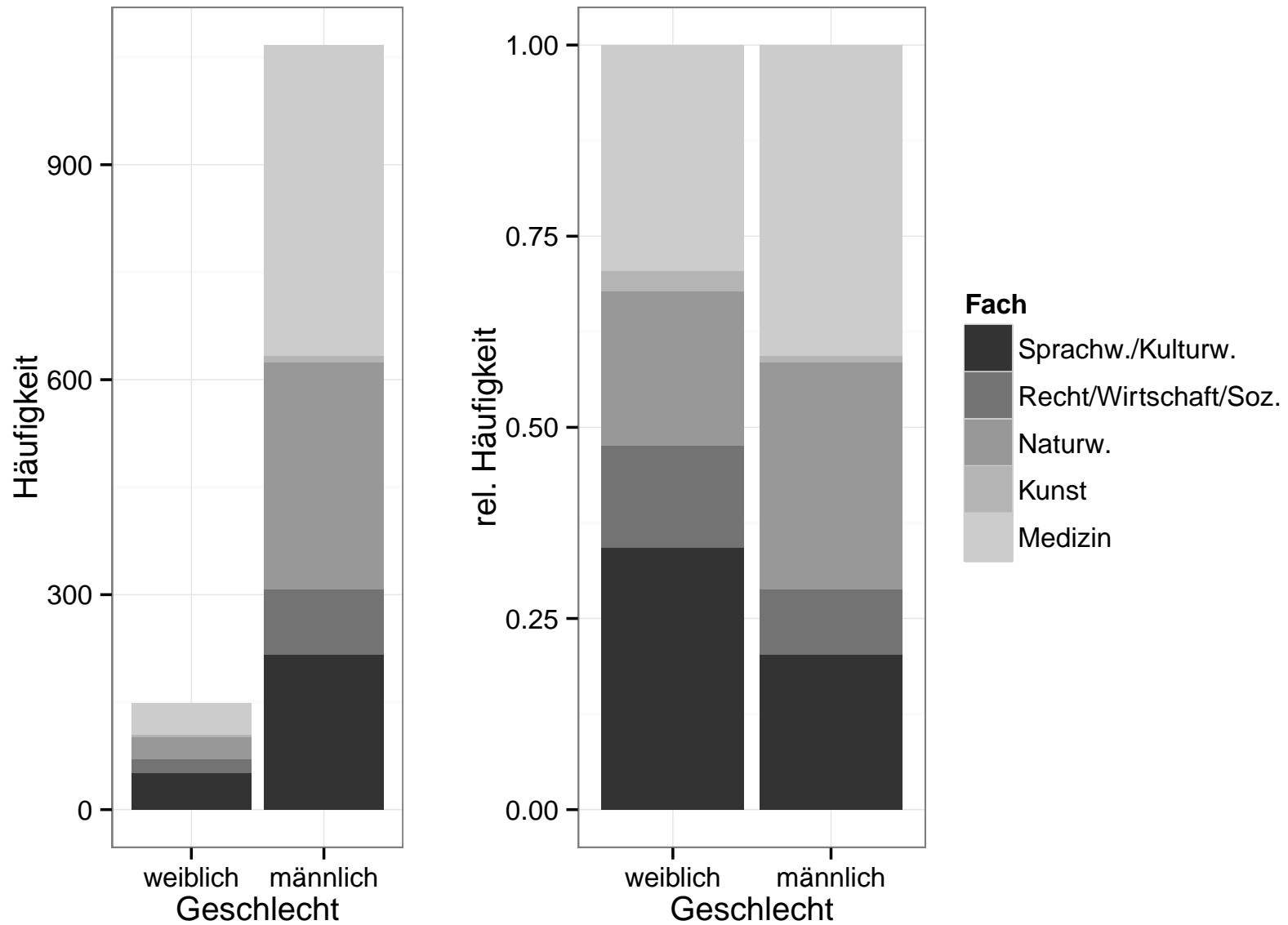
Oft bessere Alternative: Heatmaps, wie Kreuztabelle (Häufigkeit wird als Farbe kodiert)

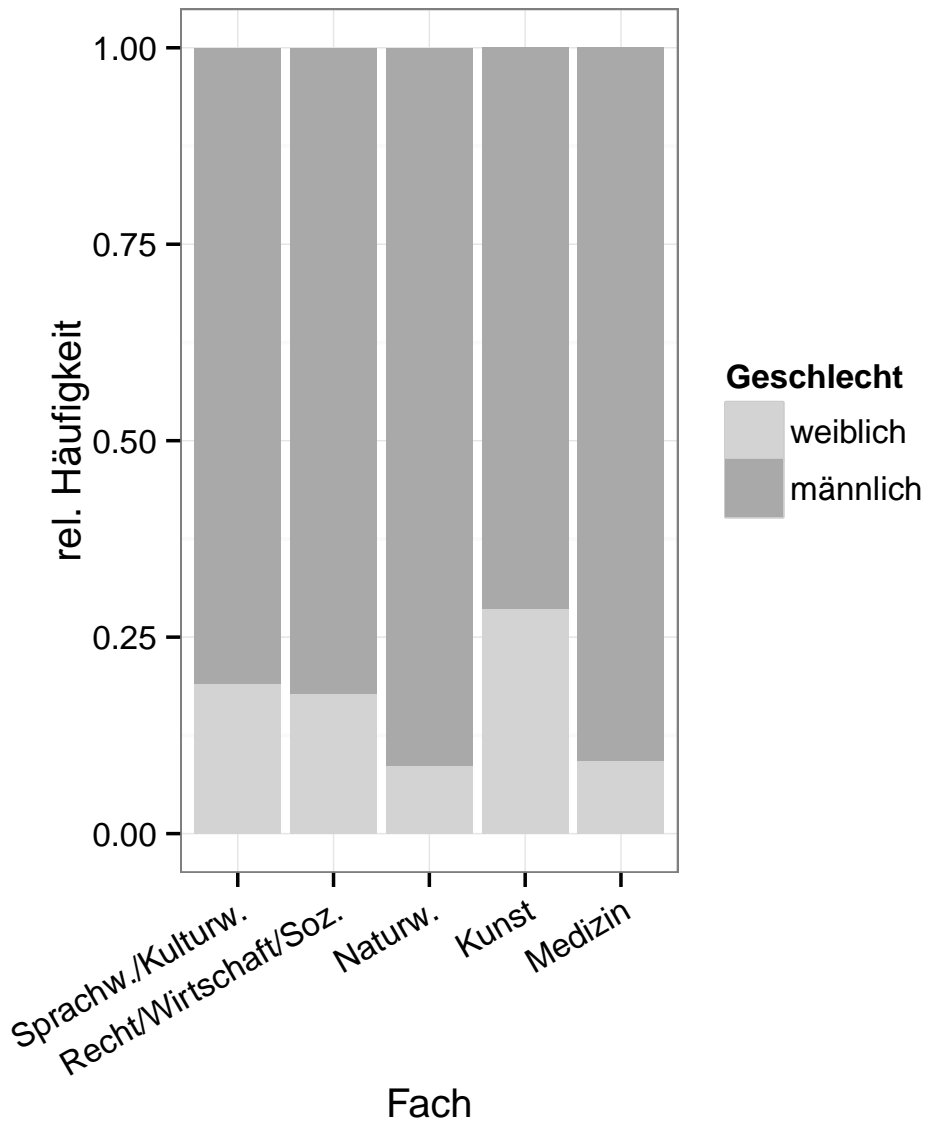
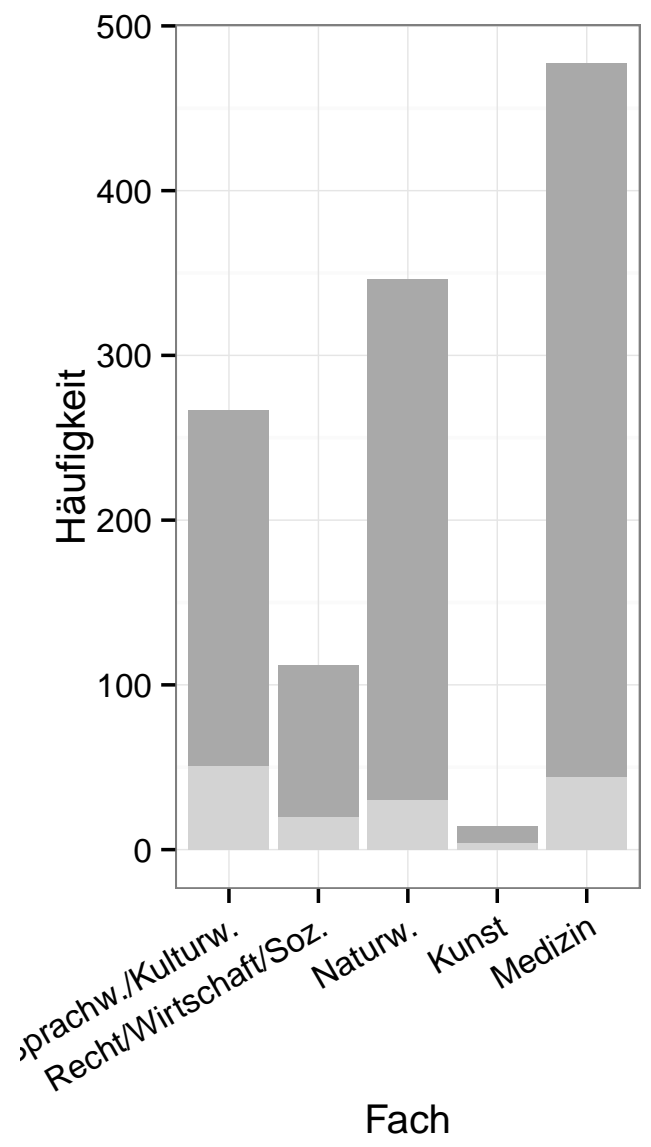


Zum Vergleich der Mosaikplot:



Gestapelte Balkendiagramme & gestapelte skalierte Balkendiagramme





5.2.4 Bedingte Häufigkeitsverteilungen

Bsp. 5.9.

Habilitationen nach Geschlecht und Fach (nach Fahrmeir et al., 2009).

Grundgesamtheit: alle Habilitationen 1993

Geschlecht: X

Fächergruppe: Y

		Sprachw. Kulturw.	Rechtsw. Wirts., Soz.	Naturw.	Kunst	Medizin	
		1	2	3	4	5	
weiblich	1	51	20	30	4	44	149
männlich	2	216	92	316	10	433	1067
		267	112	346	14	477	1216

Zur Interpretation:

Definition 5.10.

Gegeben sei ein bivariates Merkmal (X, Y) mit gemeinsamer Häufigkeitsverteilung (h_{ij}) , $i = 1, \dots, k$; $j = 1, \dots, m$.

Seien $h_{i\bullet} > 0$ und $h_{\bullet j} > 0$ für alle i, j . Für jedes $i = 1, \dots, k$ heißt

$$f_{Y|X}(b_1|a_i) := \frac{h_{i1}}{h_{i\bullet}} = \frac{h(a_i, b_1)}{h(a_i)}, \quad \dots, \quad f_{Y|X}(b_m|a_i) := \frac{h_{im}}{h_{i\bullet}} = \frac{h(a_i, b_m)}{h(a_i)}$$

bedingte (relative) Häufigkeitsverteilung von Y unter der Bedingung $X = a_i$.

Analog heißt für jedes $j = 1, \dots, m$

$$f_{X|Y}(a_1|b_j) := \frac{h_{1j}}{h_{\bullet j}} = \frac{h(a_1, b_j)}{h(b_j)}, \quad \dots, \quad f_{X|Y}(a_k|b_j) := \frac{h_{kj}}{h_{\bullet j}} = \frac{h(a_k, b_j)}{h(b_j)}$$

bedingte (relative) Häufigkeitsverteilung von X unter der Bedingung $Y = b_j$.

Bedingte Verteilungen werden immer als relative Häufigkeiten ausgedrückt. Für die Berechnung gilt

$$f_{X|Y}(a_i|b_j) = \frac{h_{ij}}{h_{\bullet j}} = \frac{\frac{h_{ij}}{n}}{\frac{h_{\bullet j}}{n}} = \frac{f_{ij}}{f_{\bullet j}}$$

und analog

$$f_{Y|X}(b_j|a_i) = \frac{h_{ij}}{h_{i\bullet}} = \frac{f_{ij}}{f_{i\bullet}}.$$

Die Verwechslung von gemeinsamer und bedingter Verteilung bzw. verschiedener bedingter Verteilungen ist eine häufige Fehlerquelle.

Konvention: Bei Vermutung über Richtung des Zusammenhangs betrachtet man vorwiegend die bedingte Verteilung der abhängigen („bewirkten“) Variablen gegeben die

festen Werte der unabhängigen („wirkenden“) Variable. In diese Richtung geht ja auch die „Prognose“! Man kennt den Wert der unabhängigen Variablen und will Aussagen über die abhängige machen.

Bsp. 5.11.

Bedingte Verteilung der Fächergruppen gegeben das Geschlecht ($f_{Y|X}(b_j|a_i)$ für verschiedene i).

		Sprachw. Kulturw.	Rechtsw. Wirts., Soz.	Naturw.	Kunst	Medizin	
		1	2	3	4	5	
weiblich	1	51	20	30	4	44	149
männlich	2	216	92	316	10	433	1067
		267	112	346	14	477	1216

Bedingte Verteilung der Fachgruppe gegeben das Geschlecht $(f_{Y|X}(b_j|a_i))_j$ für verschiedene j .

$Y : b_j$		Sprachw.	Rechtsw.	Naturw.	Kunst	Medizin
		Kulturw.	Wirts., Soz.			
$X : a_i$		1	2	3	4	5
weiblich	1					
männlich	2					

Kontrolle: Bei den Variablen aus der Bedingung ergibt sich immer eine Randsumme der relativen Häufigkeiten von 1.

$$f_{Y|X}(\text{Rechtsw.}|\text{weiblich}) =$$

$$f_{Y|X}(\text{Kunst}|\text{männlich}) =$$

Bedingte Verteilung des Geschlechts gegeben die Fachgruppe $((f_{X|Y}(a_i|b_j))_i)$ für verschiedene j).

		b_j				
		Sprachw.	Rechtsw.	Naturw.	Kunst	Medizin
a_i		Kulturw.	Wirts., Soz.			
		1	2	3	4	5
weiblich	1					
männlich	2					

$$f_{X|Y}(\text{weiblich}|\text{Rechtsw.}) =$$

$$f_{X|Y}(\text{männlich}|\text{Kunst}) =$$

Nochmals zur Interpretation:

1. $f_{X|Y}(\text{weiblich}|\text{Medizin}) = .$

2. $f_{Y|X}(\text{Medizin}|\text{weiblich}) = .$

3. $f_{15} = f(\text{Medizin und weiblich}) = .$

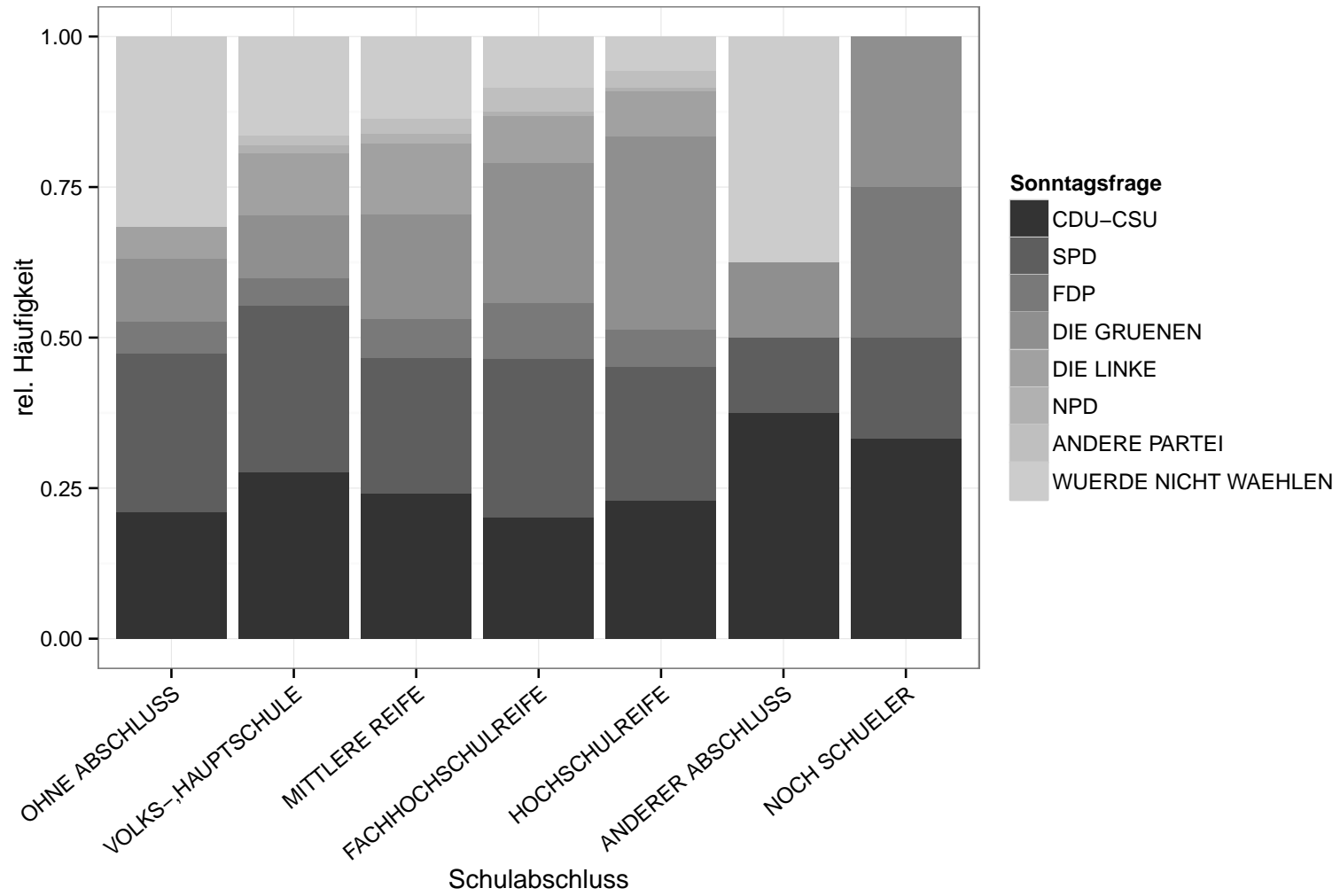
Es liegt jeweils eine andere Grundgesamtheit zu Grunde:

Bsp. 5.12. *Wahlabsicht und Bildung*

	CDU-CSU	SPD	FDP	GRÜNE	LINKE	NPD	ANDERE	NICHTW.	
OHNE ABSCHLUSS	4	5	1	2	1	0	0	6	19
VOLKS-,HAUPTSCHULE	203	203	33	77	75	10	12	120	733
MITTLERE REIFE	180	169	48	130	88	13	18	102	748
FACHHOCHSCHULREIFE	26	34	12	30	10	1	5	11	129
HOCHSCHULREIFE	124	120	34	173	41	3	15	31	541
ANDERER ABSCHLUSS	3	1	0	1	0	0	0	3	8
NOCH SCHUELER	4	2	3	3	0	0	0	0	12
	544	534	131	416	215	27	50	273	2190

Bedingt auf Bildung:

	CDU-CSU	SPD	FDP	DIE GRUENEN	DIE LINKE	NPD	ANDERE	NICHTW.
OHNE ABSCHLUSS	0.21	0.26	0.05	0.11	0.05	0.00	0.00	0.32
VOLKS-,HAUPTSCHULE	0.28	0.28	0.05	0.11	0.10	0.01	0.02	0.16
MITTLERE REIFE	0.24	0.23	0.06	0.17	0.12	0.02	0.02	0.14
FACHHOCHSCHULREIFE	0.20	0.26	0.09	0.23	0.08	0.01	0.04	0.09
HOCHSCHULREIFE	0.23	0.22	0.06	0.32	0.08	0.01	0.03	0.06
ANDERER ABSCHLUSS	0.38	0.12	0.00	0.12	0.00	0.00	0.00	0.38
NOCH SCHUELER	0.33	0.17	0.25	0.25	0.00	0.00	0.00	0.00
Gesamt	0.25	0.24	0.06	0.19	0.10	0.01	0.02	0.12



Bedingt auf Wahlabsicht:

	CDU-CSU	SPD	FDP	DIE GRUENEN	DIE LINKE	NPD	ANDERE	NICHTW.	Gesamt
OHNE ABSCHLUSS	0.01	0.01	0.01	0.00	0.00	0.00	0.00	0.02	0.01
VOLKS-,HAUPTSCHULE	0.37	0.38	0.25	0.19	0.35	0.37	0.24	0.44	0.33
MITTLERE REIFE	0.33	0.32	0.37	0.31	0.41	0.48	0.36	0.37	0.34
FACHHOCHSCHULREIFE	0.05	0.06	0.09	0.07	0.05	0.04	0.10	0.04	0.06
HOCHSCHULREIFE	0.23	0.22	0.26	0.42	0.19	0.11	0.30	0.11	0.25
ANDERER ABSCHLUSS	0.01	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.00
NOCH SCHUELER	0.01	0.00	0.02	0.01	0.00	0.00	0.00	0.00	0.01

