

Formelsammlung Statistik

Andrey Behrens

August 2009

Das ist eine Formelsammlung für Statistik. Die Formelsammlung enthält alle Formeln aus dem Skript des Wintersemesters 2009/2010. Außerdem ein paar Sachen die mir sinnvoll erschienen und für die Klausur notwendig sein könnten, sowie Formblätter zum schnellen Ausfüllen während der Klausur.

Quellen sind (1) Statistikscript Prof. Dr. Müller, HS Wismar und (2) Taschenbuch der Wirtschaftsmathematik, Wolfgang Eichholz und Eberhard Vilkner.

Teil I.

Begriffe

Statistische Masse	Umfang der Einheiten einer statistischen Untersuchung
Statistische Einheit	Untersuchungsobjekt einer statistischen Untersuchung. Träger der interessanten Informationen.
Merkmal	Zu betrachtendes Attribut einer Einheit. Etwa Einkommen, Altern, ...
Merkmalstypen	<p>diskrete Merkmalstypen bestehen aus einer überschaubare, endliche Menge (etwa Geschlecht),</p> <p>stetige Merkmalstypen können in einem bestimmten Bereich jeden reellen Wert annehmen,</p> <p>quasi-stetige Merkmalstypen sind eigentlich diskret, enthalten aber sehr grosse Menge von möglichen Merkmalen</p>
Merkmalsausprägung	<p>Gruppierung Sortierung, gleiche Merkmalsausprägung</p> <p>Klassifizierung benachbarte Ausprägungen werden zu einer Klasse zusammengefasst. Übliche Schreibweise $[200; 400)$ mit der Bedeutung $200 \leq x < 400$.</p>
Skalenniveau	<p>nominal qualitativ (also keine Zahlen), etwa Geschlecht oder Studiengang. Darstellung als gruppierter Wert.</p> <p>ordinal Merkmalsausprägung mit objektiver Rangordnung, Abstände sind aber nicht bezifferbar (etwa Noten). Darstellung als gruppierter Wert.</p> <p>metrisch <i>Interval quantitativ</i>: reelle Zahlen, natürliche Rangfolge, eindeutige Abstände, etwa Sparsumme, <i>Verhältnis quantitativ</i>: reelle Zahlen, natürliche Rangfolge, eindeutige Abstände, absoluter Bezugspunkt (etwa Nullpunkt). Beispiel: Alter. Darstellung als klassierter Wert.</p>

Teil II.

Univariate Datenanalyse

Auto ¹	x_i	h_i	H_i	f_i	F_i	$x_i h_i$	$x_i \cdot f_i$	$(x_i)^2 \cdot f_i$
BMW	342	1	1	0,167	0,167	342	57,00	19494
Mercedes	549	1	2	0,167	0,333	549	91,50	50234
VW/Audi	1501	1	5	0,167	0,833	1501	250,17	3,76E5
Sonstige	1713	1	6	0,167	1,000	1713	285,50	4,89E5

Tabelle 0.1.: Beispiel gruppierter, nominaler Werte

x_i	h_i	H_i	f_i	F_i	$x_i h_i$	$x_i \cdot f_i$	$(x_i)^2 \cdot f_i$
280	1	1	0,1	0,1	280
340	2	3	0,2	0,3	680
740	1	9	0,1	0,9	740
1180	1	10	0,1	1,0	1180

Tabelle 0.2.: Beispiel gruppierter, ordinaler Werte

x_i	h_i	H_i	f_i	F_i	Δx_i	f_i^*	h_i^*
[200;400)	21	21	0,21	0,21	200	0,00105	0,1050
[700;1000)	19	96	0,19	0,96	300	0,00063	0,0633
[1000;1500)	2	98	0,02	0,98	500	0,00004	0,0040
[1500;2000)	2	100	0,02	1,00	500	0,00004	0,0040

Tabelle 0.3.: Beispiel klassierter, metrischer Werte

Math	TR	Formel	Erläuterung
h_i	hi		abs. Häufigkeit
H_i	shi	$h_1 + \dots + h_i = \sum_{j=1}^i h_j$	abs. Summenhäufigkeit cumsum(hi)
f_i	fi	$\frac{h_i}{N}$ mit $\sum_{i=1}^k f_i$	relative Häufigkeit gruppiert Stabdiagramm siehe Abbildung 0.3 auf Seite 13 klassiert Histogramm, siehe Abbildung 0.4 auf Seite 14
F_i	sfi	$f_1 + \dots + f_i = \sum_{j=1}^i f_j$	abs. Summenhäufigkeit
N	n	$\sum_{i=1}^k h_i$	Stat Masse
h_i^*	his	$\frac{h_i}{\Delta x_i}$	abs Häufigkeitsdichte
f_i^*	fis	$\frac{f_i}{\Delta x_i}$	rel Häufigkeitsdichte
$F(x)$	f(x)		Verteilungsfunktion, Funktion der relativen Summenhäufigkeit Als Beispiel für gruppierte Daten: $F(500)=0,30 \rightarrow$ Es wird nicht gerechnet, sondern aus dem Diagramm abgelesen, da es sich um gruppierte Werte handelt! Als grafische Lösung (Treppendiagramm, keine Zwischenwerte!) siehe Abbildung 0.1 auf Seite 13
	Gruppe	$F(x) = \begin{cases} 0 & x < x_1 \\ F_i & x_i \leq x < x_{i+1} \\ 1 & x \geq x_k \end{cases}$	Bei klassierten Daten: Klasse aus Diagramm ablesen (H_i), untere und obere Grenzen der Klasse herauslesen, in Formel einsetzen: $F(500) = 0.21 + \frac{0.56}{300}(500 - 400) = 0,397 = 39,7\%$ als grafische Lösung siehe Funktionsdiagramm 0.2 auf Seite 13
	Klasse	$F(x) = \begin{cases} 0 & x < x_1^u \\ F(x_i^u) + \frac{f_i}{\Delta x_i}(x - x_i^u) & x_i^u \leq x < x_i^o \\ 1 & x \geq x_k^o \end{cases}$	

Tabelle 0.4.: Überblick Häufigkeiten

Name	Math	TR	nominal	ordinal	metrisch	Vor- und Nachteile
Modal	x_D	xd	ja	ja	ja	Ist die Merkmalsausprägung, die am häufigsten vorkommt. Es kann mehrere Modalwerte geben.
Median	x_z	xz	?	ja	ja	Mitte aller Merkmalsträger, bzw. welcher Merkmalswert wird von der Hälfte aller Merkmalsträger nicht überschritten. Vorteil: Robust gegen Ausreißer.
Quantil	x_p	xp	?	?	?	ein Teil aller Merkmalsträger (etwa 0,25x oder 0,75x) bzw. welcher Merkmalswert wird von einem Teil aller Merkmalsträger nicht überschritten. Dabe ist das $x_p = x_{0.5} = x_z$
Arith. Mittelw.	\bar{x}	xs	nein	nein	ja	Der Durchschnitt oder Mittelwert aller Merkmale
Geom Mittelw.	x_G	xg	?	?	ja	Mittelwert für Produkte, etwa bei Verhältnissen oder Wachstumswerten. Nur für Zahlen >0 sinnvoll.

Tabelle 0.5.: Überblick Lageparameter

Math	TR	Formel	Erläuterung
x_D	xd	<p>Gruppen da x_i wo f_i am größten ist</p> <p>Klassen Mitte der modalen Klasse $x_D = \frac{x_i^u + x_i^o}{2} = x_i'$</p>	Ist die Merkmalsausprägung, die am häufigsten vorkommt. Es kann mehrere Modalwerte geben.
x_z	xz	<p>Gruppe $x_z = 0.5N$</p> <p>Klasse $x_z = x_i^u + \frac{0.5 - F(x_i^u)}{f_i} * \Delta x_i$</p>	Median bzw. Zentralwert ist der Wert, der in der Mitte der Variantsreihe liegt. Ist N gerade, wird der Mittelwert der zwei mittelsten Werte ermittelt. Beispiel: Zuerst Klasse bestimmen und dann $400 + \frac{0.5 - 0.21}{0.56} * 300 = 555.36$
x_p	xp	<p>Gruppe $x_p = p \cdot N$</p> <p>Klasse $x_p = x_i^u + \frac{p - F(x_i^u)}{f_i} * \Delta x_i$</p>	Wird eine Variationsreihe in gleich große Teile zerlegt, entstehen Quantile. Typisch sind 0,25, 05, 0,75. Der Quantilabstand ist $Q = x_{0.75} - x_{0.25}$. Das 0,5-Quantil ist gleich dem Median. Quantil ist gewissermaßen das Gegenüber der Verteilungsfunktion!
\bar{x}	xa	<p>Gruppe $\bar{x} = \sum_{i=1}^k x_i f_i$</p> <p>Klasse $\bar{x} = \sum_{i=1}^k x_i' f_i$</p>	Arithmetischer Mittelwert bzw. Durchschnitt. Durchschnitte werden mit dieser Formel addiert: $\bar{x} = \frac{\sum_{m=1}^k N_m * \bar{x}_m}{\sum_{m=1}^k N_m}$
x_G	xg	$x_G = \sqrt[N]{\prod_{i=1}^k x_i}$	Der Geometrische Mittelwert wird bei der Mittelung von Wachstumsraten oder multiplikativ verknüpften Daten angewendet.

Tabelle 0.6.: Lageparameter

Math	TR	Formel	Erläuterung
R	r	Gruppe $R = x_{max} - x_{min}$	Die Spannweite ist die Differenz zwischen größtem und kleinstem Merkmalswert.
		Klasse $R = x_k^o - x_1^u$	
Q	q	$Q = x_{0.75} - x_{0.25}$	Der Quantilsabstand ist der Abstand zwischen oberem und unterem Quantil.
s_x^2	s2x	Gruppe $s_x^2 = \left\{ \sum_{i=1}^k [(x_i)^2 \cdot f_i] \right\} - \bar{x}$	Die Varianz ist die mittlere quadratische Abweichung aller Merkmalsausprägungen vom arith. Mittelwert.
		Klasse $s_x^2 = \left\{ \sum_{i=1}^k [(x'_i)^2 \cdot f_i] \right\} - \bar{x}^2$	Alternative Zeichen der Varianz sind $s_x^2 = s^2 = \sigma^2$
		$s^2 = \frac{N_1 s_1^2 + N_2 s_2^2}{N_1 + N_2} + \frac{N_1 (\bar{x}_1 - \bar{x})^2 + N_2 (\bar{x}_2 - \bar{x})^2}{N_1 + N_2}$	Varianz der Grundgesamtheit. Gleichungsbeispiel bei der Annahme, dass es zwei Teilmengen gibt und die jeweils die Varianzen und Mittelwerte bekannt sind.
s_x	sx	$s_x = \sqrt{s_x^2}$	Standardabweichung mittlere Abweichung vom Mittelwert. Nachteil: Bei großen Merkmalsmengen nimmt die Schwankungsbreite zu. Ein Vergleich zwischen Messreihen mit großen und kleinen Verteilungen ist daher ggf. nicht mehr sinnvoll. Statt dessen: Variationskoeffizient.
v	v	$v = \frac{s_x}{\bar{x}}$	Hinweis: Bei $s_x = 0$ gibt es einen eindeutigen Hinweis auf Konzentration. Ansonsten nicht. Variationskoeffizient = Auf den Mittelwert bezogenes relatives Streuungsmaß, sofern nur positive Werte auftreten.

Tabelle 0.7.: Streuungsparameter

Math	TR	Formel	Erläuterung
p_I	pi	Gruppe	$p_i = \frac{x_i \cdot h_i}{N \cdot \bar{x}}$
		Klasse	$p_i = \frac{x'_i \cdot h_i}{N \cdot \bar{x}}$
P_i	spi	$P_i = \sum_{j=1}^i p_j$	Das Konzentrationsmaß beschreibt die relative Merkmalssumme. Die Lorenzkurve veranschaulicht das Konzentrationsmaß grafisch.
$L(F_i)$	-	-	Die Fläche zwischen Gleichverteilung und Lorenzkurve wird als Lorenzfläche bezeichnet und ist ein weiteres Konzentrationsmaß. Je größer die Fläche, desto größer die Konzentration. Beispiel zur Lorenzkurve siehe 0.5 auf Seite 14
			Lorenzkurve: Welchen Anteil haben Merkmalsträger an Merkmalen. Etwa 0.5 = 50% der Autohersteller (F_i) haben Anteil von 0.25 = 25% P_i der Produktion
G	g	$G = \frac{0.5 - A(L)}{0.5}$ im Bereich $0 \leq G \leq 1$	Stat. Maß zur Darstellung der Ungleichverteilung. Der Gini-Koeffizient misst die Höhe der relativen Konzentration über das Verhältnis der Lorenzfläche zur Fläche bei maximaler Konzentration. Es kann unterschiedliche Lorenzflächen bei identischem G geben. <i>Eigenschaft:</i> Werden alle x_i um denselben Prozentsatz erhöht oder gesenkt, dann bleibt der Gini-Koeffizient unverändert. Werden alle x_i um einen additiven Zuschlag erhöht, dann wird der Gini-Koeffizient kleiner. Wird ein x_i -teibetrag von einem größeren zu einem kleineren x_i transferiert, so wird der Gini-Koeffizient kleiner. <i>Beispiel:</i> Der Ginikoeffizient für die Einkommensverteilung liegt in Deutschland bei 0,274 (2003), in Frankreich bei 0,327 (1995), in Großbritannien bei 0,360 (1999), in Japan bei 0,249 (1993) und in den USA bei 0,408
$A(L)$	al	$A(L) = \sum_{i=1}^k \frac{P_{i-1} + P_i}{2} \cdot f_i$ mit $P_0 = 0$	Fläche unter der Trapezkurve Hinweis: Sind großgeschriebene P's also spi.

Tabelle 0.8.: Relative Konzentration

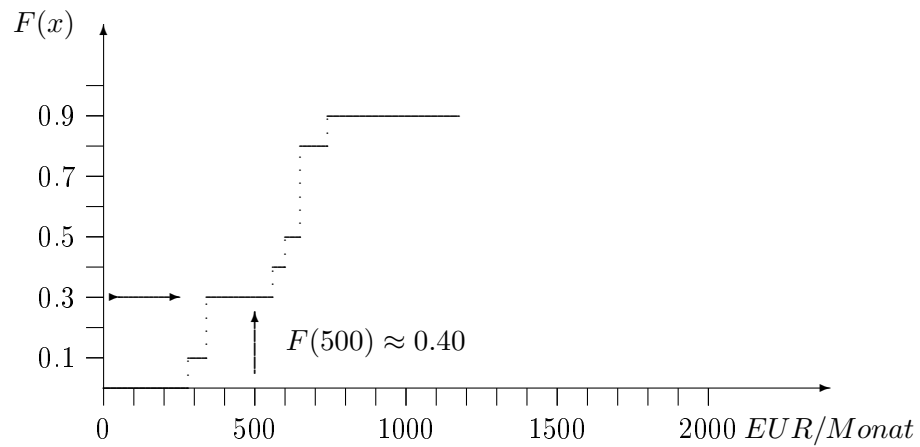


Abbildung 0.1.: Funktion relativer Sumenhäufigkeit $F(x)$ bei gruppierten Daten

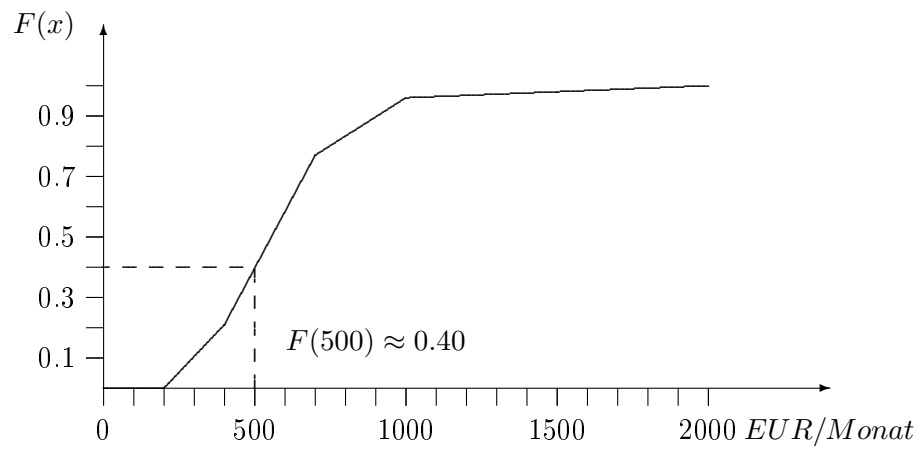


Abbildung 0.2.: Funktion $F(x)$ relativer Summenhäufigkeit bei klass. Daten

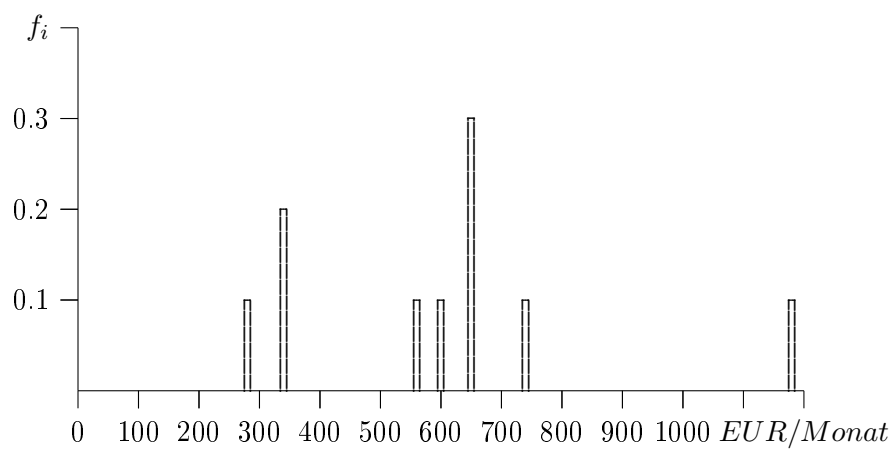


Abbildung 0.3.: Darstellung rel Häufigkeit von Gruppen: Stabdiagramm

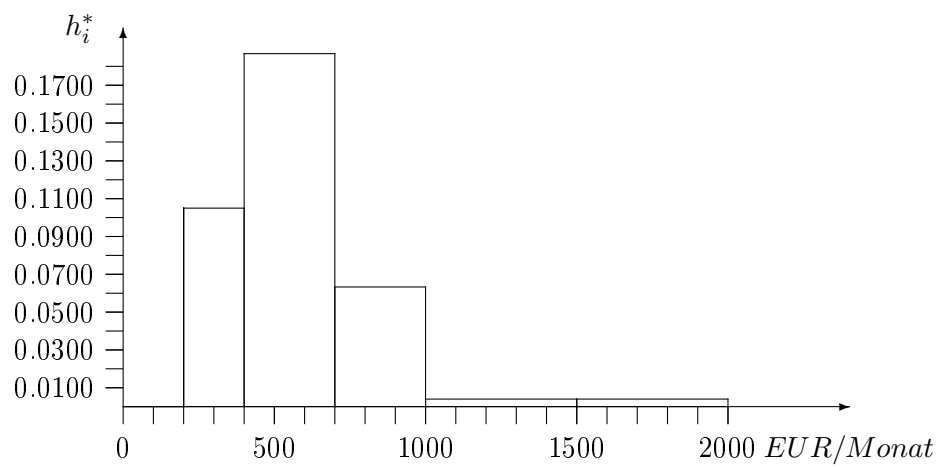


Abbildung 0.4.: Darstellung rel. Häufigkeit von Klassen: Histogramm

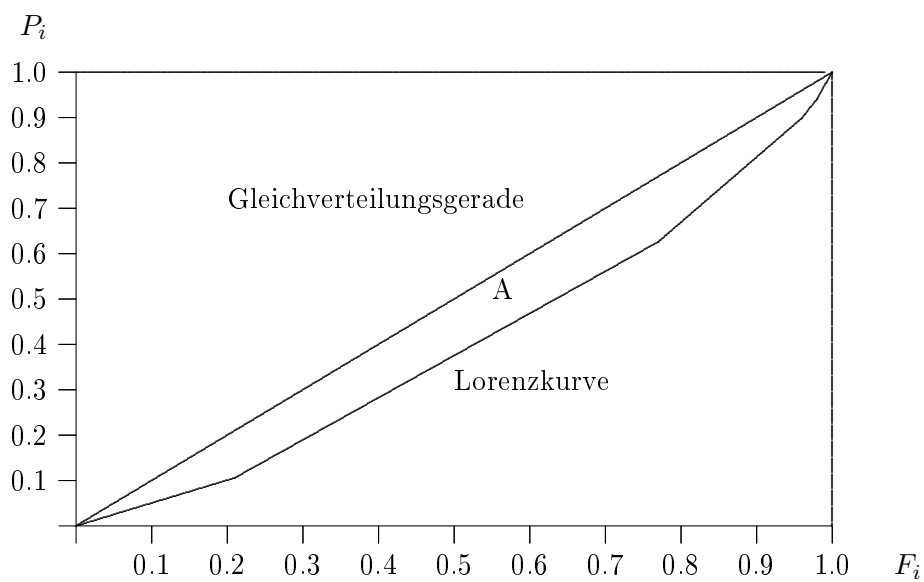


Abbildung 0.5.: Lorenzkurve

Teil III.

Formblätter

x	x	Klasse oder Gruppe einer statistischen Zählung. Variable kann Zeichen haben wie 1, i , k die für das 1-te, i -te oder letzte Gruppe/Klasse stehen.
x_d	xd	Modalwert, der Wert mit der häufigsten Merkmalsausprägung
x_z	xz	Median, Mitte aller Merkmalsausprägungen, d.h. nach oben und unten gleich viele Merkmalsausprägungen
x_p	xp	Quantile überschreiten einen gewissen Anteil von Merkmalsausprägungen <i>nicht</i>
x'_i		Klassenmitte der i -ten Klasse
x_i^u x_i^o		untere bzw. obere Grenze der i -ten Klasse
h	h	Anzahl von Einheiten innerhalb einer Gruppe oder Klasse. Tiefgestellte Zeichen gleiche Bedeutung wie bei x Die Summe aller h ist die statistische Masse
H_i	shi	absolute Summenhäufigkeit, wie h_i aber aufsteigend addiert. Der größte Wert= N
f_i	fi	relative Häufigkeit. Summe aller $f_i = 1$ Entspricht dem prozentualen Anteil an der statistischen Masse.
F_i	sfi	relative Summenhäufigkeit. Wie f_i aber aufsummiert. Der größte Wert = 1
Δx_i	dx _i	Klassenbreite der i -ten Klasse
s_i	si	relative Summenhäufigkeit einer Klasse
N	n	Statistische Masse, also die Menge aller Merkmalsausprägungen.

Table .9.: Überblick Variablen

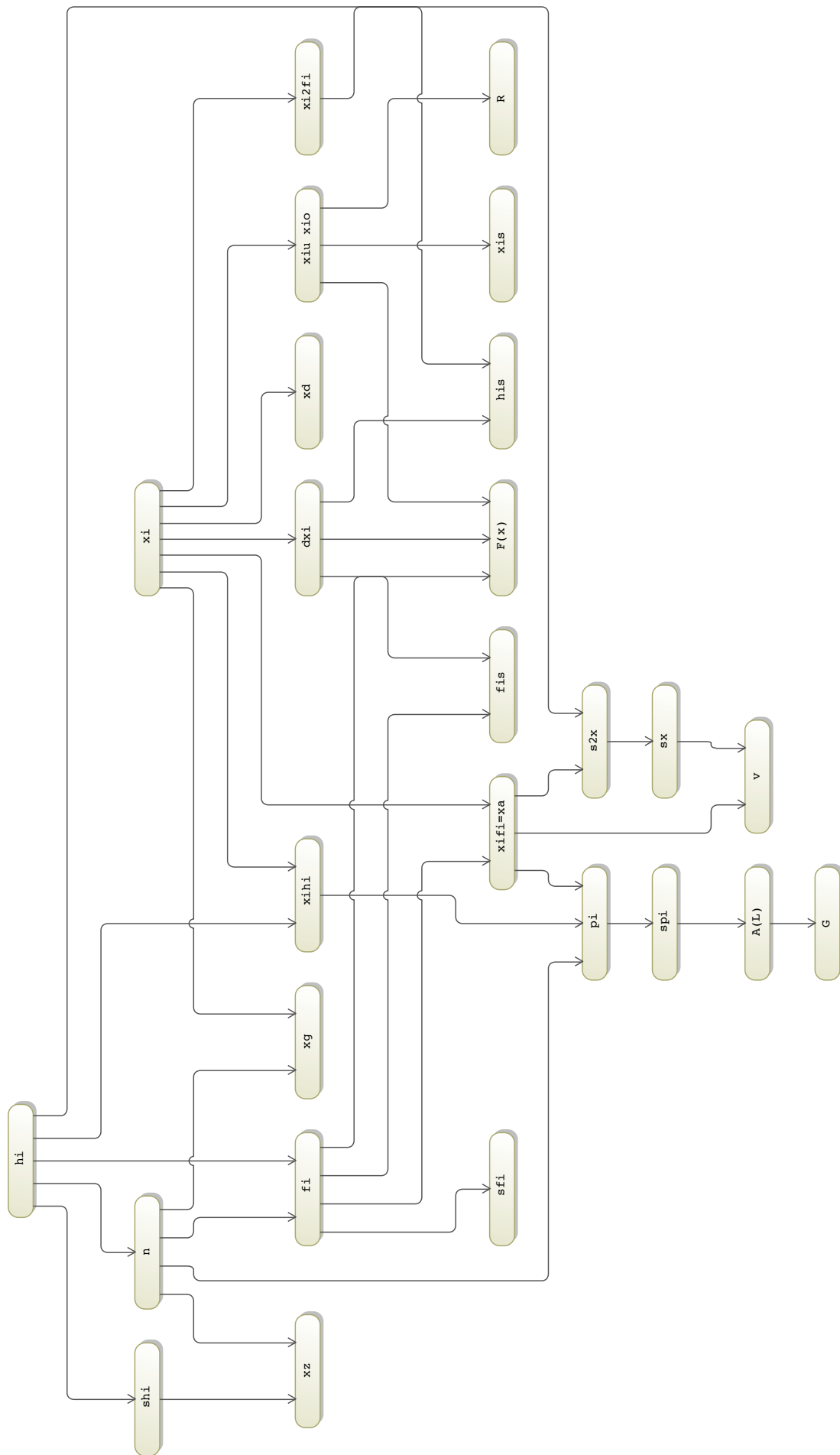


Figure .6.: Zusammenhänge von Variablen

Gruppe	abs. Häufig	abs. Summen- häufig	rel. Häufig	rel. Summen- häufig	$x_i \cdot f_i$	$(x_i)^2 \cdot f_i$	Konz- koeff.	Konz- maß	Fläche unter Lorenzkurve
x_i	h_i	H_i	f_i	F_i	$x_i \cdot h_i$		p_i	P_i	$A(L)$
\sum	$N =$	-	$= 1$	-	$\bar{x} =$	$=$	-	-	$=$

$$G = \frac{0.5 - A(L)}{0.5} = \frac{0.5 -}{0.5} =$$

$$G = \frac{0.5 - A(L)}{0.5} = \frac{0.5 -}{0.5} =$$

