

## Grammatik von Grafiken

Grafik = Daten + geometrische Elemente + Ästhetische Zuordnung  
+ Datentransformationen + Skalen + Koordinatensysteme  
+ Facettierung + Theme + {Grafik}

Geometrische Elemente = Punkte | Linien | Rechtecke | Boxplots | Dichtefkt. | ...

Ästhetische Zuordnung = Position & Farbe & Größe & Form

Datentransformationen = id | Mittelwerte | Anteile | ...

Skalen = Achsenabschnitte & Farbe & Legenden & Achsenbeschriftung & ...

Koordinatensysteme = kartesisch | logarithmisch | Polarkoord. | Kartenproj. | ...

Facettierung = small multiples | lattice plot | plot | ...

Theme = Font & Gitterlinien & Hintergrundfarben & Layout von Text & ...

## Lorenzkurve

Das Merkmal darf nur positive Werte annehmen.

$x_{(1)}, \dots, x_{(n)}$  sei die geordnete Stichprobe.

Die Lorenzkurve verbindet Punktepaare bestehend aus den Teilsommen von  $x_{(i)}$  (d.h.  $\sum_{i=1}^j x_{(i)}$ ) und dem relativen Anteil an Individuen, die diese Teilsumme besitzen.

## Berechnung

$u_{(0)} = 0, v_{(0)} = 0, j = 1, \dots, n$

$u_{(j)} := \frac{1}{n} \cdot (\text{Aufteilung } x\text{-Achse})$

$v_{(j)} := \frac{1}{n} \cdot x_{(j)} \cdot \frac{1}{\sum_{i=1}^n x_{(i)}}$   
(y-Werte)

$v_{(1)}, \dots, v_{(n)}$  ist monoton steigend

## Gini-Koeffizient

Maßzahl, die das Ausmaß der Konzentration beschreibt

Definiert als  $G = 2 \cdot F \in [0, \frac{n-1}{n}]$ , wobei  $F$  die Fläche zwischen

$y = x$  und der Lorenzkurve ist.  $G = 1 - \frac{1}{n} \cdot \sum_{i=1}^n (v_{(i-1)} + v_{(i)})$

$G = \frac{n-1}{n} - \frac{2}{n} \cdot \sum_{i=1}^{n-1} v_{(i)}$  Normiert:  $G^* = \frac{n}{n-1} \cdot G \in [0, 1]$

$G^* = 0$  bedeutet keine Konzentration (Gleichverteilung)

$G^* = 1$  bedeutet volle Konzentration (Monopol)

## Merke

Lorenzkurve muss

monoton steigend,

stetig, kleiner  $y = x$ ,

und eine Funktion sein

Startet in (0,0), endet

in (1,1)

## Faltungen

Sind  $X_1, X_2$  unabhängig, so gilt  $f_{X_1+X_2}(z) = \int_{\mathbb{R}} f_{X_1}(x_1) \cdot f_{X_2}(z-x_1) d\mu(x_1)$

! Grenzen bestimmen sich dadurch, dass  $x_1 \in \mathcal{X}_{X_1}$  und  $z-x_1 \in \mathcal{X}_{X_2}$  gelten muss.

## Beispiel

Seien  $X_1 \sim P(\lambda_1), X_2 \sim P(\lambda_2)$  und unabhängig. Wir betrachten  $\bar{Y} = X_1 + X_2$ .

$$P(Y=n) = P(X_1+X_2=n) = \sum_{k=0}^n P(X_1=k) P(X_2=n-k) = \sum_{k=0}^n \underbrace{\binom{n}{k} \cdot \lambda_1^k \cdot \lambda_2^{n-k}}_{= (\lambda_1 + \lambda_2)^n} \cdot \frac{e^{-\lambda_1 - \lambda_2}}{n!}$$

## Kontingenztafel und bedingte Häufigkeitsverteilung

$X \backslash Y$	$b_1$	...	$b_m$	
$a_1$	$h_{11}$	...	$h_{1m}$	$h_{1.}$
$\vdots$				
	absolute Häufigkeiten			
$\vdots$				
$a_k$	$h_{k1}$	...	$h_{km}$	$h_{k.}$
	$h_{.1}$	...	$h_{.m}$	$h_{..} = n$
	Randhäufigkeit von $b_j$ in $Y$			

Für relative Häufigkeiten ersetze  $h_{ij}$  durch  $f_{ij} = \frac{h_{ij}}{n}$ .

Bedingte Häufigkeitsverteilung:  $f(Y=b_j | X=a_i)$

$$f_Y(b_j | a_i) = \frac{h_{ij}}{h_{i.}}, \dots, f_Y(b_m | a_i) = \frac{h_{im}}{h_{i.}}$$

## Erwartete absolute/relative Häufigkeit

Wenn  $X, Y$  unabhängig sind, erwarten wir: absolute Häufigkeit:  $\hat{h}_{ij} = \frac{h_{i.} \cdot h_{.j}}{n}$

relative Häufigkeit:  $\hat{f}_{ij} = \frac{f_{i.} \cdot f_{.j}}{1}$

## Mosaikplot

		1	2
Y			
X	1		
	2		

Breite vom Kasten  $Y=j$  ist  $\frac{h_{.j}}{n}$

Höhe vom Kasten  $ij$  ist  $\frac{h_{ij}}{h_{i.}}$

## (Korrigierter) Kontingenzkoeffizient

Normierung von  $\chi^2$ : Kontingenzkoeffizient  $K := \sqrt{\frac{\chi^2}{n + \chi^2}}$ ,  $K \in [0, \sqrt{\frac{\min\{k, m\}-1}{\min\{k, m\}}}]$

Korrigierter Kontingenzkoeffizient  $K^* := \frac{K}{\sqrt{\frac{\min\{k, m\}-1}{\min\{k, m\}}}}$ ,  $K^* \in [0, 1]$

## Sensitivität und Spezifität

$Y \in \{0, 1\}$  (Zielgröße),  $X$  mindestens ordinalskaliert

$Y = 1$ : "positiver" Fall,  $Y = 0$ : "negativer" Fall

Sei  $\hat{y}_i$  die Prognose für  $y_i$  auf Basis von  $x_i$ . Es soll gelten

$\hat{y}_i = 1 \Leftrightarrow x_i \geq c$

	$y_i = 0$	$y_i = 1$	
Vorhersage $\hat{y}_i = 0$	wahr negativ	falsch negativ	# negative Vorhersagen
Vorhersage $\hat{y}_i = 1$	falsch positiv	wahr positiv	# positive Vorhersagen
	# negative	# positive	

## ROC-Kurve

Die ROC-Kurve zeigt die Zuverlässigkeit der Vorhersagen

für alle möglichen Schwellenwerte  $c$  an.

Verbindet die Punkte  $(FPR(c), TPR(c)) \forall c \in [x_{(1)}, x_{(n)}]$

Für  $c < x_{(1)} \Rightarrow \hat{y}_i = 1 \forall i \Rightarrow (FPR(c), TPR(c)) = (0, 1)$

Für  $c > x_{(n)} \Rightarrow \hat{y}_i = 0 \forall i \Rightarrow (FPR(c), TPR(c)) = (1, 0)$

## Bedingte Odds / Odds ratio

Bedingte Odds:  $y(1,2|X=a_i) = \frac{h_{i1}}{h_{i2}}$

Relative Chancen (Odds ratio):  $y(1,2|X=a_i) = \frac{y(1,2|X=a_i)}{y(1,2|X=a_j)} = \frac{h_{i1} \cdot h_{j2}}{h_{i2} \cdot h_{j1}}$

! Symmetrisches Maß und Risikofaktor

$y > 1$ : Odds in  $X=a_i$  höher als in  $X=a_j$

## Odds-Update

$$\frac{P[B|A]}{P[B^c|A]} = \frac{P[A|B]}{P[A|B^c]} \cdot \frac{P[B]}{P[B^c]}$$

a-posteriori-Verhältnis = neue Information / a-priori-Verhältnis

## $\chi^2$ -Koeffizient (geht auf allen Skalen, aber nur sinnvoll für nominal)

Zusammenhangsmaß zum Quantifizieren von "Abstand" zwischen beobachteten und erwarteten Häufigkeiten.

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^m \frac{(h_{ij} - \hat{h}_{ij})^2}{\hat{h}_{ij}} = n \cdot \sum_{i=1}^k \sum_{j=1}^m \frac{(f_{ij} - \hat{f}_{ij})^2}{\hat{f}_{ij}}$$

## Satz

Für eine Kontingenztafel der Form

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \text{ gilt } \chi^2 = \frac{n \cdot (a \cdot d - b \cdot c)^2}{(a+b)(a+c)(c+d)(b+d)}$$

! Misst nur Stärke des Zusammenhangs von  $X$  und  $Y$ .

Nicht die Richtung, wie bei  $\gamma$

## Schiefe

Ist  $X^3$  quasi-int. bar, dann heißt

$$g(X) := \frac{E[(X - E[X])^3]}{\sqrt{\text{Var}[X]^3}} \quad \text{Schiefe von } X.$$

$g(X) = 0$ : symmetrisch  
 $g(X) > 0$ : linkssteil  
 $g(X) < 0$ : rechtssteil

## Kurtosis

Ist  $X^4$  quasi-int. bar, dann heißt

$$K(X) := \frac{E[(X - E[X])^4]}{\text{Var}[X]^2} \quad \text{Kurtosis von } X$$

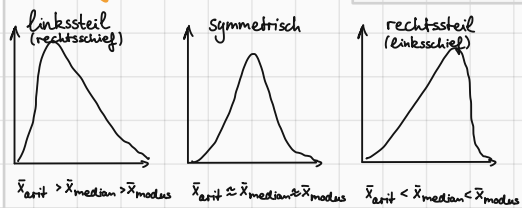
$K^*(X) = K(X) - 3$  heißt Exzess-Kurtosis.

$K^*(X) = 0$ : mesokurtisch

$K^*(X) > 0$ : leptokurtisch (viele extreme Werte)

$K^*(X) < 0$ : platykurtisch (wenig extreme Werte)

## Verteilung charakterisieren



## Verteilungen

Verteilung	Verteilungsfunktion
$U(a, b)$	$F_X(x) = \frac{x-a}{b-a} \cdot \mathbb{1}_{[a,b)}(x) + \mathbb{1}_{[b,\infty)}(x)$
$N(\mu, \sigma^2)$	—
Cauchy	$F_X(x) = \frac{1}{\pi} + \frac{\arctan(\frac{x-\mu}{\sigma})}{\pi}$
Gamma $\text{Ga}(a, b)$ $\Gamma(a, b)$	$F_X(x) = \frac{1}{\Gamma(a)} \int_0^x t^{a-1} \cdot e^{-t/b} \cdot \frac{1}{b} dt$ $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$ $\Gamma(a, b) = \frac{\Gamma(a) \cdot \Gamma(b)}{\Gamma(a+b)} = \int_0^1 x^{a-1} (1-x)^{b-1} dx$ $\Gamma(a, b) = \frac{\Gamma(a) \cdot \Gamma(b)}{\Gamma(a+b)}$ $\Gamma(a, b) = \frac{\Gamma(a) \cdot \Gamma(b)}{\Gamma(a+b)}$ $\Gamma(a, b) = \frac{\Gamma(a) \cdot \Gamma(b)}{\Gamma(a+b)}$
Beta( $\alpha, \beta$ )	$F_X(x) = \frac{1}{B(\alpha, \beta)} \int_0^x t^{\alpha-1} (1-t)^{\beta-1} dt$ $B(\alpha, \beta) = \frac{\Gamma(\alpha) \cdot \Gamma(\beta)}{\Gamma(\alpha+\beta)}$
$\text{Exp}(\lambda)$	$F_X(x) = (1 - e^{-\lambda x}) \cdot \mathbb{1}_{[0, \infty)}(x)$
$\text{Bin}(n, \pi)$	$F_X(x) = \sum_{k=0}^x \binom{n}{k} \cdot \pi^k \cdot (1-\pi)^{n-k} \cdot \mathbb{1}_{[0, n]}(x) + \mathbb{1}_{(n, \infty)}(x)$
$\text{Poi}(\lambda)$	$F_X(x) = \sum_{k=0}^x P(X=k) = \frac{\Gamma(Lx+1, \lambda)}{[x]!}$ $\Gamma(x, \lambda) = \int_x^\infty t^{x-1} e^{-t} dt$
$\text{Geom}(\pi)$	$F_X(x) = 1 - (1-\pi)^{[x]}$ $\mathbb{P}(\text{genau } x \text{ Versuche für ersten Erfolg einer Bernoulli-ZV})$
$\text{Geom}(\pi)$	$F_X(x) = 1 - (1-\pi)^{[x]+1}$ $\mathbb{P}(\text{genau } x \text{ Fehlversuche vor erstem Erfolg einer Bernoulli-ZV})$

## Konvergenz

$$X_n \xrightarrow{f} X : \mathbb{P}(\lim_{n \rightarrow \infty} f(X_n) = f(X)) = 1$$

$$X_n \xrightarrow{f} X : \forall \varepsilon > 0 : \lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0$$

$$X_n \xrightarrow{f} X : \lim_{n \rightarrow \infty} E[|X_n - X|^p] = 0$$

$$X_n \xrightarrow{d} X : \forall x \in \mathbb{R} : \lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

## Datenformen

Querschnittsdaten: Eine Beobachtung pro UE.

Zeitreihe: Mehrere Beobachtungen einer UE.

Längsschnittdaten: Mehrere Beobachtungen mehrerer UE.

## Wahrscheinlichkeitsmaß

$(\Omega, \mathcal{F})$  sei ein Maßraum und  $\mu: \mathcal{F} \rightarrow \mathbb{R}$

i)  $\mu(\emptyset) = 0$  ii)  $\forall A \in \mathcal{F} : \mu(A) \geq 0$

iii)  $\sigma$ -Additivität:  $\mu(\bigcup_{i=1}^\infty A_i) = \sum_{i=1}^\infty \mu(A_i)$

iv) Normiertheit:  $\mu(\Omega) = 1$

## $\sigma$ -Algebra

$\mathcal{F}$  heißt  $\sigma$ -Algebra, wenn

1)  $\Omega \in \mathcal{F}$  2)  $A \in \mathcal{F} \Rightarrow \bar{A} \in \mathcal{F}$

3)  $A_i \in \mathcal{F}, i \in \mathbb{I} \Rightarrow \bigcup_{i \in \mathbb{I}} A_i \in \mathcal{F}$

## Hypergeometrische Verteilung

Wahrscheinlichkeit für  $x$  Erfolge beim "ziehen"

von  $n$  Elementen aus einer Stichprobe der Größe

$N$ , wobei es insgesamt  $M$  günstige Elemente gibt.

$X$  ist hypergeometrisch verteilt  $X \sim H(N, M, n)$ ,

wenn  $P_{N, M, n}(X=x) = \frac{\binom{M}{x} \cdot \binom{N-M}{n-x}}{\binom{N}{n}}$ , für

$x \in \{\max\{0, n-(N-M)\}, \dots, \min\{n, M\}\}$

Dann ist  $F_X(x) = \sum_{k=0}^x P_{N, M, n}(X=k) \cdot \mathbb{1}_{[0, n]}$

## Farbskalentypen

- Qualitativ: (eher) nur für nominales Skalenniveau.
- Sequentiell: mindestens ordinales Skalenniveau.
- Divergent: mindestens ordinales Skalenniveau mit "neutralem" mittleren Wert.

## Sinnvolle Zusammenhangsmaße

$X, Y$  nominal  $\rightarrow \chi^2$ , Odds ratio, Kontingenzkoeffizient

$X, Y$  ordinal  $\rightarrow$  Rangkorr. nach Spearman  $r_{sp}$

$X, Y$  metrisch  $\rightarrow$  Korrelation Pearson  $r_{op}$

## Bedingte Momente

$$E(X|Z=z) = \begin{cases} \sum_{x \in T_X} x \cdot P(X=x|Z=z), & X \text{ diskret} \\ \int_{\mathbb{R}} x \cdot f_{X|Z}(x|Z=z) dx, & X \text{ stetig} \end{cases}$$

$E(X|Z)$  ist ein ZV  
 $E(X|Z=z) = E(X|Z(z))$  ist ein Zahlenwert

$$\text{Var}(X|Z=z) = E(X - E(X|Z=z))^2 | Z=z) = \begin{cases} \sum_{x \in T_X} (x - E(X|Z=z))^2 \cdot P(X=x|Z=z), & X \text{ diskret} \\ \int_{\mathbb{R}} (x - E(X|Z=z))^2 \cdot f_{X|Z}(x|Z=z) dx, & X \text{ stetig} \end{cases}$$

## Satz vom iterierten Erwartungswert

Für beliebige ZV  $X, Y$  und Funktion  $f$  gilt:

$$E[E(f(X)|Z)] = E(f(X))$$

## Satz von der totalen Varianz

Für beliebige ZV  $X, Y$  gilt:

$$\text{Var}(X) = E(\text{Var}(X|Z)) + \text{Var}(E(X|Z))$$

↑ Erwartete bedingte Varianz    ↑ Varianz des bedingten Erwartungswertes

## Markov- und Chebyshev-Ungleichungen

Sei  $X: \Omega \rightarrow \mathbb{R}$  eine reelle ZV. Dann gilt

$$\forall \varepsilon > 0: \mathbb{P}(|X| > \varepsilon) \leq \frac{1}{\varepsilon^2} \cdot E[X^2]$$

Markov-Ungleichung ( $n=1$ ):  $\mathbb{P}(X > \varepsilon) \leq \frac{E[X]}{\varepsilon}$

Chebyshev-Ungleichung ( $n=2$ ):  $\mathbb{P}(|X - E[X]| > \varepsilon) \leq \frac{1}{\varepsilon^2} \cdot \text{Var}(X)$

## Jensen-Ungleichung

Sei  $X$  eine int. bare ZV und

$g: \mathbb{R} \rightarrow \mathbb{R}$  konvex. Dann gilt

$$E[g(X)] \geq g(E[X]).$$

Ist  $f: \mathbb{R} \rightarrow \mathbb{R}$  konkav, so gilt  $f(E[X]) \geq E[f(X)]$

## Zentraler Grenzwertsatz

Seien  $(X_n)_{n \in \mathbb{N}}$  i.i.d. ZVen mit  $E[X_n] = \mu$  und  $\text{Var}[X_n] = \sigma^2 < \infty$ . Sei  $S_n = \sum_{i=1}^n X_i$ . Dann gilt:

$$1) \lim_{n \rightarrow \infty} \mathbb{P}\left[\frac{S_n - n\mu}{\sqrt{n}\sigma} \leq x\right] = \mathcal{N}(0, 1) \quad \text{bzw.} \quad \frac{S_n - n\mu}{\sqrt{n}\sigma} \xrightarrow{D} \mathcal{N}(0, 1) \quad \text{bzw.} \quad S_n \sim \mathcal{N}(n\mu, n\sigma^2) \quad \mathbb{P}(S_n \leq x) = \Phi\left(\frac{x - n\mu}{\sqrt{n}\sigma}\right)$$

$$2) \lim_{n \rightarrow \infty} \mathbb{P}\left[\frac{S_n - n\mu}{\sigma/\sqrt{n}} \leq x\right] = \mathcal{N}(0, 1) \quad \text{bzw.} \quad \frac{S_n - n\mu}{\sigma/\sqrt{n}} \xrightarrow{D} \mathcal{N}(0, 1) \quad \text{bzw.} \quad \bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) \quad \mathbb{P}(\bar{X}_n \leq x) = \Phi\left(\frac{x - \mu}{\frac{\sigma}{\sqrt{n}}}\right)$$

## Charakteristische / Momenterzeugende Funktion

Sei  $X$  eine ZV. Die Funktion  $\varphi_X: \mathbb{R} \rightarrow \mathbb{C}$  mit

$$\varphi_X(t) := E[\exp(itX)] = \int \exp(itx) dP = \int f_X(x) \cdot \exp(itx) d\mu$$

heißt charakteristische Fkt. von  $X$ . Die Funktion  $M: \mathbb{D} \rightarrow \mathbb{R}$

$$M(t) := E[\exp(tX)] = \int f_X(x) \cdot \exp(tx) d\mu \quad \text{heißt momenterzeug. Fkt.}$$

## Dichte transformationssatz

Sei  $X: \Omega \rightarrow \mathbb{R}$  eine ZV mit stetiger Verteilungsfkt.  $F_X$  und Dichte

$f_X(x) = \frac{\partial F_X(x)}{\partial x}$  bzgl. des  $\lambda$ -Maßes. Sei  $g: \mathbb{R} \rightarrow \mathbb{R}$  bijektiv und stetig diff. bar mit

$g'(x) \neq 0$ . Dann hat  $g \circ X$  die Dichte  $f_{g \circ X}(y) = f_X(g^{-1}(y)) \cdot |g^{-1}'(y)|$