

# Anmerkungen zur Vorlesung Statistik vom 11.10.2016

## Inhalt:

univariate deskriptive Statistik:

- Häufigkeiten und Darstellungsmöglichkeiten (Plots)
- Lageparameter: Modus, arithmetisches Mittel, Median, Quantile
- empirische Verteilungsfunktion

## Aufgabensammlung Statistik:

	<b>Aufgaben zu R Grundlagen</b>	<b>3</b>
04.10.2016 Hausaufgabe A1 - A7	Aufgabe 1: RStudio und erste Versuche . . .	4
	Aufgabe 2: Zuweisungen und Variablen . . .	6
	Aufgabe 3: Vektoren . . . . .	8
	Aufgabe 4: Mehrere Merkmale: Data Frames	10
	Aufgabe 5: Skalenniveaus und Data Frames .	12
	Aufgabe 6: Datenimport aus Textdateien . . .	13
	Aufgabe 7: R-Skripten als Logbuch . . . . .	14
11.10.2016 Hausaufgabe A8 - A13 A10: ohne Boxplot	Aufgabe 8: Deskriptives mit R . . . . .	15
	Aufgabe 9: Einfache Grafiken in R . . . . .	17
	Aufgabe 10: Emp. Vtlgs.f. Quantil Boxplot .	18
	<b>Aufgaben zur deskriptiven Statistik</b>	<b>19</b>
	Aufgabe 11: Häufigkeit1b . . . . .	19
	Aufgabe 12: Lageparameter . . . . .	20
	Aufgabe 13: Lageparameter . . . . .	21

Notation:

$X$  (Großbuchstaben) : Merkmalsvariable (später auch Zufallsvariable)  
(typisch auch:  $Y, Z$ )

$x$  : Realisation(en) einer Merkmalsvariable

$i$  ( $i = 1, \dots, n$ ) : Index der Realisationen, steht für einen bestimmten Merkmalsträger

$n$  : Anzahl der Merkmalsträger

$(x_1, x_2, \dots, x_n)$  : Urliste

weil Urlisten oft unübersichtlich  $\rightarrow$  Häufigkeitsverteilungen

$a_j$  : Merkmalsausprägung der Merk. var. Index  $j = 1, \dots, J$

$h(a_j)$  : absolute Häufigkeiten (# des Auftretens)

$H(a_j) = \sum_{k=1}^j h(a_k)$  : kumulierte absolute Häufigkeiten

$f(a_j) = \frac{h(a_j)}{n}$  : relative Häufigkeit

$F(a_j) = \sum_{k=1}^j f(a_k)$  : kumulierte relative Häufigkeit

Anmerkung:  $H(a_j)$  und  $F(a_j)$  ( $\hat{=}$  Häufigkeitsverteilungen) sind nur sinnvoll, falls Merkmal  $X$  geordnet werden kann.

Falls  $X$  metrisch  $\wedge$  empirische Verteilungsfunktion

( $\hat{=}$  Erweiterung der kumulierten relativen Häufigkeit auf alle Zwischenwerte)

$\downarrow$  d.h.

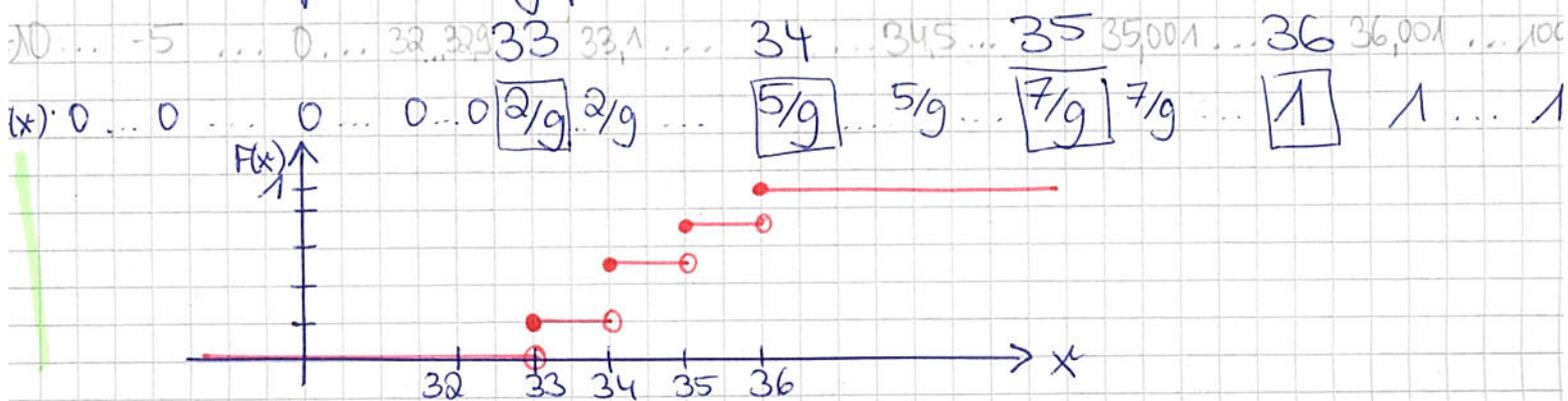
wir betrachten nicht mehr nur  $F(a_j)$ ,  $j = 1, \dots, J$  sondern

$F(x) \quad \forall x \in \mathbb{R}$

4.1 sortierte "Urliste":

	33	33	34	34	34	35	35	36	36
$\downarrow$	1		2			3		4	
$a_j$	33		34			35		36	
$h(a_j)$	2		3			2		2	$\Sigma = g = n$
$f(a_j)$	$2/g$		$3/g$			$2/g$		$2/g$	
$F(a_j)$	$2/g$		$5/g$			$7/g$		$9/g = 1$	

$\downarrow$  empirische Verteilfkt.  $F(x)$



4.1 Wesentliche Eigenschaften der empirischen Verteilfkt.:

- 1)  $F(x) = 0$  für  $x < a_1$  ( $\hat{=}$  kleinste Merkmalsausprägung)
- 2)  $F(x) = 1$  für  $x \geq a_g$  ( $\hat{=}$  größte "—" )
- 3) "Treppenfunktion" mit Sprungstellen an den Merkmalsausprägungen  $a_j$  ( $\rightarrow$  # Sprungstellen =  $g$ )
- 4) mit Sprunghöhe  $f(a_j)$

! Je mehr Sprungstellen, desto glatter wird die Funktion.

Beispiele: ecdf für Variable "Alter" mit  $\mathbb{R}$

FL2

[?] Für welchen Wert von  $x$  erreicht  $F(x)$  einen bestimmten Wert  $p$  ( $p \in [0, 1]$ )  
↳ Wahrscheinlichkeit

Beispiel (Fortsetzung):

"Sortierte" Urliste: 33 33 34 34 34 35 35 36 36  
 $X_{[1]}$   $X_{[2]}$   $X_{[3]}$   $X_{[4]}$   $X_{[5]}$   $X_{[6]}$   $X_{[7]}$   $X_{[8]}$   $X_{[9]}$

also:  $n = 9$

Ges.:  $\tilde{X}_{0.4}$  und  $\tilde{X}_{2/9}$

1)  $\tilde{X}_{0.4} = \frac{0.4 \cdot 9}{p} = 3.6 \notin \mathbb{N}_0 \Rightarrow \tilde{X}_{0.4} = X_{\lceil n \cdot p \rceil} = X_{\lceil 3.6 \rceil} = X_{[4]} = 34$   
4. Wert der sort. Urliste  
*Aufrundungsfunktion*

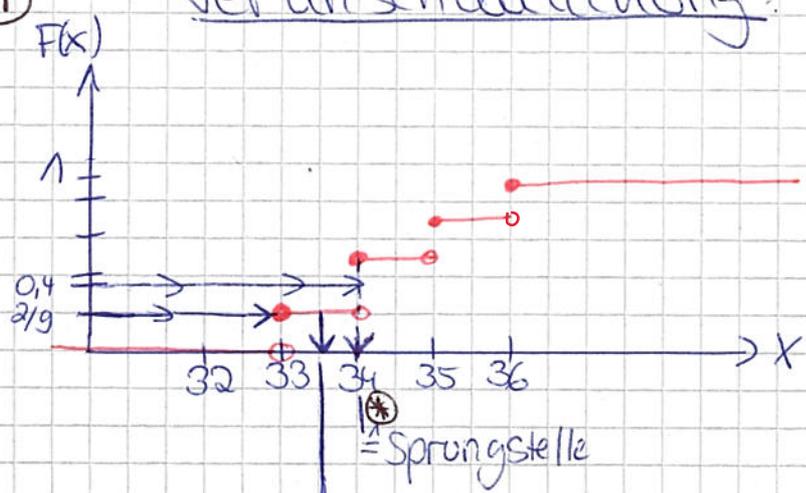
2)  $\tilde{X}_{2/9} = \frac{2/9 \cdot 9}{p} = 2 \in \mathbb{N}_0 \Rightarrow \tilde{X}_{2/9} = \frac{1}{2} (X_{[n \cdot p]} + X_{[n \cdot p + 1]}) = \frac{1}{2} (X_{[2]} + X_{[3]}) = \frac{1}{2} \cdot (33 + 34) = 33.5$

②

R. Umsetzung!

④

① Veranschaulichung:



⊛  
Zwischenwert / Mittelwert  
zwischen 2 Sprungstellen  
⊛⊛

⊛ Quantileswert zu  $p=0,4$ , also  $\hat{x}_{0,4}$

⊛⊛ Quantileswert zu  $p=2/9$ , also  $\hat{x}_{2/9}$

**F43-46**

Visualisierungsmöglichkeiten mit R

**43** Balkendiagramm: für absolute Häufigkeiten (Säulen)  
 für relative Häufigkeiten  
Attribut  $y_{lim} = c(\dots)$

**44** Kreisdiagramm:  $w_j =$  Winkel für Merkmalsausprägung  $a_j$   
 $w_j = \frac{360^\circ}{n} \cdot f(a_j)$   
 Funktion: addmargins()

**45** Balkendiagramme, Klassen getrennt oder gestapelt:

**46** klassierte Daten  $\rightarrow$  Histogramme:

Besondere Eigenschaft: Fläche ist proport. zur Häufigk.

$c$  ist frei wählbar,  
 z.B.  $c = 1$  oder  $c = 1/12$

Möglichkeit Alter und Alter darstellen / nächster Schritt: Dichtfunktion sel. Merkmal

**F4S11**

- nun: Lage- und Streuungsparameter & visuelle Darst.
- Modus
  - Median
  - arithm. Mittel
  - Spannweite SP
  - MSE (mittlere quad. Abw.)  
 $\frac{MQA}{s^2}$
  - Standardabweichung  $s$
  - Variationskoeffizient  $v$
  - Boxplot

**F5D**

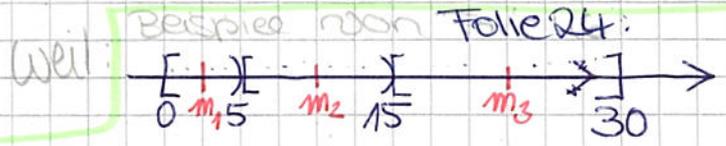
**Modus**: = Modalwert  $x_{Mod} =$  die häufigste Ausprägung  $a_j$

**Median**: = der mittlere Wert

$x_{Med} = \tilde{x}_{0.5}$

**arithmetisches Mittel**:  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{j=1}^J a_j \cdot h(a_j) = \sum_{j=1}^J a_j \cdot f(a_j)$

bei klassierten Daten:  
 ( $\hat{=}$  Schätzwert!)



$\rightarrow$  Überschätzung:  $\bar{x}$  ist zu hoch  
 $\rightarrow$  Unterschätzung:  $\bar{x}$  ist zu niedrig

**F52**

- ③ Frau =  $X_{\text{Mod}}$  mit größter Anzahl
- ① `na.exclude()` falls Merkmalsvariable fehlende Werte enthält
- ② für metrische Merkmale  
    `summary(.)` statt separat `median()` / `mean()` möglich
- ④ `summary(.)` für Faktoren  $\hat{=}$  Häufigkeitstabelle