Was versteht man in der Umgangssprache unter einer "Statistik"? ine "Statistik" ist eine systematische Zusammenstellung von Zahlen und naten zur Beschreibung von Zuständen, Entwicklungen und Phänomen. Jacent zur deschriebung vor Lüssanien, Ertwikkunger und Prainanienselspiele Häufigkeitsverteilungen, Zeitreihenvergleich und -analyse, zusammenhangs- und Abhängigkeitsanalysen, statistische Kennzahlen zu Beschreibung von Verteilungen, Zusammenhängen. Worin unterscheiden sich "beschreibende" und "Schließende statistik" m Rahmen der "Beschreibende Statistik" sammelt man Daten bei allen ein der Statistik von der Statistik von der Weschreibende Statistik" sammelt man Daten bei allen ein der Statistik von der Statis

Intersuchungseinheiten, über die man Informationen erhalten will. Die untersuchungseinheiten, über die man Informationen ernätien will. Die seschreibende Extatstik hat zur Ziel, empirische Daten durch Tabellen, kennzahlen (auch: Maßzahlen oder Parameter) und Grafiken übersichtlich darzustellen und zu ordnen. Dies ist vor allem bei umfangreichem Datenmaterial sinnvoll, da dieses nicht leicht überblickt verden kann. In der schließenden Statistik" wählt man aus der Grundgesamtheit, über die man Informationen haben will, eine Grungesamtneit, über die man informationen naben will, eine Teilimenge = Stichprobe aus. (Grund für diese Vorgehensweise ist meistens die Größe der Grundgesamtheit). Die schließende Statistik wir zu einem wesentlichen Teil zum Beweis oder zur Widerlegung von vorhe aufgestellten Behauptungen, den Hypothesen, die sich auf derfinierte Grundgesamtheit bezieht, eingesetzt. Nur bei den Einheiten der Stichprobe erhebt man Daten, die man dann mit statistischen Methoder auswertet. Von den Stichprobenergebnissen versucht man, auf die Eigenschaften der Grundgesamtheit zu schließen. Was ist die wichtigste (in der Praxis aber sehr oft nicht gegebene)

Voraussetzung für eine "Random-Auswahl"?

iehe c): Liste/Datei aller Einheiten der Grundgesamtheit

### Bei einer Teilerhebung muss man Entscheidungen fällen über den mfang und das Auswahlverfahren. Welche Ents

st wichtiger?

Die Entscheidung über das Auswahlverfahren ist wesentlich wichtiger.
Eine Stichprobe "schlecht" ausgewählt kann noch so groß sein, sie liefert keine brauchbaren Erkenntnisse über die Grundgesamtheit. So ist inmittelbar einleuchtend, dass man bei der Grundgesamtheit Liniwohner in Deutschland" (rund 82 Mio.) und der daraus gezogenen Stichprobe "Studierende in Deutschland" (rund 2 Mio.) zwar einen sehr großen Stichprobenumfang hat, dass man aber nicht von der Stichprobe auf die Grundgesamtheit schließen kann. Diese Stichprobe ist "verzerrt" oder wie man oft auch sagt "nicht repräsentativ

c) Überlegen und Beschreiben Sie den Unterschied zwischen einer

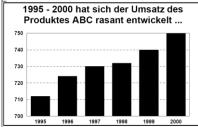
"willkürlichen Auswahl" und einer "zufälligen Auswahl".

Bei einer "willkürlichen Auswahl" (- Auswahl aufs Geratewohl oder 
\_convenience sample") gibt es keinen Auswahlplan. Die Interviewer sind frei in der Auswahl ihrer Interviewpartner. Daher suchen sie sich die Personen aus, die für sie am beguemsten zu erreichen sind. Das führt neist zu einer verzerrten Stichprobe. Bei einer "zufälligen Auswahl" ist meist zu einer verzerrten Stichprobe. Bei einer "zufäligen Auswah" ist die Auswah ist, d.h. jede Einheit der Grundgesamtheit (über die man Informationen erhalten will) muss mit gleicher Wahrscheinlichkeit in die Stichprobe gelangen können. Dies setzt vorau dass eine Liste/Datel aller Einhelten der Grundgesamtheit vorliegt. Die Interviewer sind nicht frei in der Auswahl ihrer Interviewpartner. Sie bekommen feste Zielpersonen vorgegeben. Nur bei der Zufallsauswahl (=Random-Auswahl) lässt sich mit Hilfe der Wahrscheinlichkeitsrechnung der Stichprobenfehler berechnen.

Beschreiben Sie die "Quota-Auswahl". Die "Quota-Auswahl" ist ein "bewusstes Auswahlverfahren". Über für die Interviewer verbindliche Quotenpläne wird erzwungen, dass für ausgewählte Quotierungsmerkmale (meist Geschlecht, Alter, Beruf) in de Stichprobe eine gleiche Verteilung wie in der Grundgesamtheit erreicht wird. Voraussetzung ist, dass man die Verteilung der wird. Vordabsetzung ist, dass indi die Verteilung Vordabsetzung ist, dass indi die Verteilung Quotierungsmerkmale in der Grundgesamtheit kennt. Amtliche Statistiken liefern hier in vielen Fällen die gewünschten Verteilungen Welche Art von Stichprobe ergibt sich bei einer so genannten TED-Umfrage im Fernsehen? (Bei der TED-Umfrage werden Fernsehzuschauer aufgefordert Fragen zu aktuellen Themen zu beantworten. beder der Antvortmöglichkeiten ist eine Telefonnummer zugeordnet, i dann – je nach persönlicher Meinung – gewöhlt werden soll.] Bei der TED-Unfrage im Fernsehn liegt eine willkürliche Ausswahl (des Bewohner eines Landes) vor. Es gibt keinertei Auswahlplan, jeder kann, wenn er Lust hat sich an der Umfrage beteiligen. Es gibt Personen, die grundsätzlich an solchen Umfrage niemals teilnehmen würden, andere

grundsatzunt an soniern ohmage mennas termennen wurden, andere versuchen mehrmals ihre Meinung zu äußern. Geben Sie für die folgenden Merkmale jeweils an, welcher Merkmalstys vorliegt und auf welcher der 5 Skalen das Merkmal gemessen wird. Geschwindigkeit e. Fahrzeugs = quantitatives Merkmal, Verhältnisskala echtsform e. UN = qualitatives Merkmal. Nominalskala Preis e. Produktes = quantitatives Merkmal, Verhältnisskala Umsatzklasse des UN = Rangmerkmal, Rangskala Wohnort = qualitatives Merkmal, Nomalskala Mitarbeiterzahl = quantitatives Merkmal, Absolutskala Kundenzufriedenheit = Rangmerkmal, Rangskala / Ratingskala nsatz eines UN = quantitatives Merkmal, Verhältnisskala achsemesterzahl = quantitatives Merkmal, Absolutskala

eruf = qualitatives Merkmal, Nominalskala Steuerklasse = qualitatives Merkmal, Nominalskala Einkommensklasse = Rangmerkmal, Rangskala Geburtsjahr = quantitatives Merkmal, Intervallskala Nehmen Sie zum Schaubild (Abbildung) kritisch die Stellung. Überlegen Sie was darf bei einer Tabelle und einem Diagramm auf keinen Fall fehlen?



a) Bei der Abbildung fehlt die Maßeinheit, die Quellenangabe b) Aussage in Überschrift übertrieben. Jährlicher Anstieg nur etwa 1 %. Anstieg des Umsatzes über die Jahre bei ca 5,5 %. Skalierung der Ordinate achse täuscht, einen höheren Anstieg vor. Grund: Umsatz-Achse beginnt bei 700 nicht bei 0.

Im Folgenden ist die Zielsetzung verschiedener statistischer (= empirischer) Untersuchungen beschrieben. Geben Sie jeweils an, wer der was die Untersuchungseinheiten sind oder sein könnten und wie Sie die Grundgesamtheit unter sachlichen, räumlichen und zeitlichen Aspekten abgrenzen würden. Bitte beachten Sie, dass in den Aufgaber nicht zu allen Punkten Hinweise vorhanden sind. Hier geht es um die Operationalisierung (OP)der Aufgabe.

A) Hersteller von Schokolade michte Infos über Verbrauchsgewohnheiter n jugendlichen in Süddeutschland haben:

- Potentielle Untersuchungseinheit: Menschen
- Jugendliche: Möglich 10-18 Jahre, (definition nicht eindeutig)
- Räumliche abgrenzung: Süddeutschland (OP, wo genau Bundesland?) I. Zeitliche Abgrenzung: keine Infos, wahrscheinlich Jahresende

B) Glühbirnenfabrik untersucht bei Qualitätskontrolle Brenndauer und Funktionstüchtigkeit von Glühbirnen

- Potentielle Untersuchungseinheit: Glühbirnen
- Sachliche Abgrenzung: Neu produzierte Glühbirnen in dem UN. Räumliche Abgrenzung: No Infos, Beschränkung auf Hallen/Maschine
- Zeitliche Abgrenzung: No Infos, normal zeitlicher Intervall
- Teilerhebung: Brenndauer, durch Stichprobe Vollerhebung: Funktionstüchtigkeit

JN in Versandhandelsbranche stellt fest, dass aus neuen Bundesländern viele Beschwerden kommen. Welche Produkte besonders betroffen? Potentielle Untersuchungseinheiten: Reklamationen.

achliche Abgrenzung; schriftliche Reklamationen, die eingegangen sind Räumliche Abgrenzung: aus den "neuen Bundeländern" (OP eindeutig). Zeitliche Abgrenzung: 1. 1. – 31. 12. 2006

## Hausbrauerei Tauffenbach Bochum veranstaltet im Sommer ieden onntagvormittag einen "Frühschoppen mit Jazz-Musik".Betreiber will wissen, ob Gäste zufrieden mit Musik & Angebot Getränken / Speisen 1. Potentielle Untersuchungseinheiten: Menschen. 2. Sachliche Abgrenzung: Gäste der Hausbrauerei Tauffenbach.

- Räumliche Abgrenzung: z.B. die Räumlichkeiten der Hausbrauerei
- Zeitliche Abgrenzung: Sonntagvormittage im Sommer, an denen der Frühschoppen mit Jazz" stattfindet.

"r unschopper im Jazz- Saktimuse Eine Einzelhandelskette will die räumliche Anordnung des Warensorti ments "optimieren". Dazu will das UN untersuchen lassen, ob und in welchem Umfang es bei den Einkäufen ihrer Kunden so genannte Verbundeffekte" gibt, d.h., bestimmte Produkte aus dem Sortiment häufig zusammengekauft werden (z.B. Kaffee, Filtertüten und Gebäck).

- Potentielle Untersuchungseinheiten: Kassenbons / Warenkörbe Sachliche Abgrenzung: Kassenbons der Einzelhandelskette Räumliche Abgrenzung: Keine Infos. Geschäfte einz. Regionen / Fillain Zeitliche Abgrenzung: no infos, Sinnv: Beschränk. Auf best. Jahreszeit

Überlegen Sie bitte, was bei der Einteilung der Umsatzklassen in Tabel Uberiegen sie bitze, was bei der inteilung der umsatzkiassen in Tabeilet 1 und bei dem Diagramm in Abbildung 1 (Modul 1 Folien 13 und 14), bis unter" (Abkürzung "b.u.") bedeutet. Warum schreibt man statt "bis unter" nicht einfach "bis", z.B. 100 bis 200 Tsd. €? in mathematischer Schreibweise bedeutet die Umsatzklasse "100 bis unter 200 Tsd. €", dass alle Großhändler mit einem Umsatz x, mit 100 =

< 200 Tsd. €, in diese Klasse fallen. D.h. wenn der Umsatz genau 100 Tsd. € beträgt, dann fällt er in diese Umsatzklasse, wenn der Umsatz abe isat. Euerbagt, uann laite in indises Onisatzakes, wenn uer Unisatzake genau 200 Tsd. € beträgt, dann fällt er nicht mehr in diese Klasse, sonder n die nächste Klasse. Liegt der Umsatz allerdings ganz knapp unter 200 Γsd. €, z. B. bei 199 999,99 €, dann liegt er in der Klasse "100 bis unter 200 Tsd. €"

Wenn man wie folgt klassieren würde: "100 bis 200 Tsd. €", "200 bis 300 Tsd. €", "200 bis 300 Tsd. €", uzon könnte ein Umsatz von genau 200 Tsd. € nicht eindeutig einer Umsatzklasse zugeordnet werden. Eine eindeutige Zuordnung zu genau einer Klasse ist aber bei der Klassierung von Daten unbedingt notwendig.

welche zwei Probleme hat man bei der Klassierung von Daten?
Übersichtlichkeit – Informationsverlust für klassierte Daten können keine exakten statistischen Kennzahlen (z.B. Mittelwerte) berechnet werden. Näherungswerte können nur unter bestimmten Annahmen (z.B. Gleichverteilung in den Klassen) berechnet werden.

Was versteht man unter "offenen Randklassen"?

iei "offenen Randklassen" ist entweder die Klassenuntergrenze nicht ingegeben (untere offene Randklasse, z.B. "b.u. 100 kg") oder die lassenobergrenze (obere offen Randklasse, z.B. "200 kg und schwerer Bei offenen Randklassen kann die Klassenmitte und die Klassenbreite nicht berechnet werden

Bestimmen Sie für die Klasse "150 b.u. 180 cm" die Klassenbreite und die

Klassenmitte. Klassenbreite = 180 – 150 = 30 cm Klassenmitte = (150 + 180)/2 = 165 Von 2006 bis 2007 ist der Umsatz eines Unternehmens um 50% gesunken. Im darauffolgenden Jahr 2008 war das Unternehmen erfolgreicher: der Umsatz stieg von 2007 bis 2008 um 70% erfolgreicher: der Umsatz stieg von 2007 bis 2008 um 70%.
Welche der folgenden Aussagen ist richtig? (kurze Begründung
Rechnung) Der Umsatz ist im Zeitraum von 2006 bis 2008
a) um 10 % gestiegen - Falsch
b) um 15 % gestiegen - falsch
c) um 15 % gesunken – richtig Bsp: 100/2 = 50, 50/100\*70 = 35.

50 + 35 = 85. 100 + x = 85, x = - 15

Für 200 Unternehmen liegt für das Jahr 2006 die folgende Umsatzverteilung vor. Vervollständigen Sie die klassierte Häufigkeitsverteilung. a) Bestimmen Sie den Modus, den Median und das arithmetische Mittel b) Bestimmen Sie die Streuungsparameter Spannweite w, Varianz s2, andardabweichung s, Variationskoeffizient v.

Klasse Nr. i	Umsatzklasse (Mio €)	Anzahl Unternehmen h <sub>i</sub>	Anteil f <sub>i</sub> (%)	Hi	F; (%)	Klassenmitte m <sub>i</sub>
1	0 b.u. 1	60	30%	60	30%	0,5
2	1 b.u.2	80	40%	140	70%	1,5
3	2 b.u.5	40	20%	180	90%	3,5
4	5 b.u.10	10	5%	190	95%	7,5
5	10 b.u.20	10	5%	200	100%	15
Σ		200	100%			

a)  $\overline{x}D$ = 1.5 Mio € (Mitte der Klasse 2 mit größtem hi)

infallsklasse k

Elliansiasses T (S. Bsp. 2 Folie 15 Modul 4)  $\overline{x}Z=1+(2-1)0,5-0,30,4=1+0,20,4=1,5$  ( $Mio \in$ )  $\overline{x}=\Sigma mi*fiki=1=0,5\cdot0,3+1,5\cdot0,4+3,5\cdot0,2+7,5\cdot0,05+15\cdot0,05=0,15+0,6+0,7+0$ 375+0,75+2,575 (Mio €)

b) w = 20 - 0 = 20

s2=Σ(mi-x)2\*fiki=1=(0,5-2,575)2·0,3+(1,5-2,575)2·0,4+(3,5-2,575)2·0,i+(7,5-2,575)2·0,05+(15-2,575)2·0,05=10,86 s=Vs2=V10,86=3,29 ν=sx=3,29/2,575=1,28

Im Rahmen einer Marktforschungsstudie wurden n = 12 Personen u.a. gefragt nach den drei Merkmalen Geschlecht G (w = weiblich, m geragt nach den drei merkmalen descrilectht G (W = Welbildn, m = männlich), Alter A (Alter in Jahren) und Markenpräferenz M (A = Produkt A, B = Produkt B). Die Erhebung ergab die folgenden 12
Befragungsergebnisse (Beobachtungswertekombinationen):
(w, 37, A), (m, 65, A), (w, 26, A), (m, 37, B), (w, 21, B), (m, 29, A), (w, 52, B), (m, 43, A), (w, 48, A), (m, 58, B), (w, 24, A), (m, 58, B).
Lessbeispiel: Die 1. Person ist weiblich, 37 Jahre alt und bevorzugt rodukt A.

n Buduk A. a) Erstellen Sie die zwei folgenden zweidimensionalen Kreuztabellen: Geschlecht x Markenpräferenz klassiertes Alter x Markenpräferenz (2 Altersklassen: 1. Klasse: bis unter 40 Jahre, 2. Klasse: 40 Jahre und älter) b) Versuchen Sie eine dreidimensionale Kreuztabelle zu erstellen für die . drei Merkmale: Geschlecht x klassiertes Alter x Markenpräferenz c) Welche der drei Merkmale kann man als unabhängige bzw. abhängige

rekrhale betrachten?
) Sie erheben bei einer Grundgesamtheit Daten für 3 Merkmale. Was efert mehr Information: die dreidimensionale Häufigkeitsverteilung ode alle drei möglichen ein- und zweidimensionalen Häufigkeitsverteilungen

M/A	bis unter 40 Jahre alt	40 Jahre und älter	Σ	
Produkt A	4	3	7	58,3%
Produkt B	2	3	5	41,7%
Σ	6	6	12	100%
	50%	50%		

b)						
	weib	lich	mänr	nlich		
M/GA	bis unter 40	40 Jahre und älter	bis unter 40	40 Jahre	2	
	Jahre alt		Jahre alt	und älter		
Produkt A	3	1	1	2	7	58,3 9
Produkt B	1	1	1	2	5	41,79
Σ	4	2	2	4	12	100 %
	33,3%	16,7%	16,7%	33,3%		

c) Geschlecht G und Alter A sind unabhängige Merkmale und Markenwahl ist das abhängige Merkmal. Die Antwort ist hier eindeutig. Denn Geschlecht und Alter können nicht vom Merkmal Markenwahl abhängige Merkmale sein; die Markenwahl kann das Geschlecht bzw. da Alter nicht beeinflussen. Es sist aber sinnvoll zu untersuchen, ob die Markenwahl bei Frauen und Männern unterschiedlich oder gleich ist bzw ob die Markenwahl bei einer Produktgruppe altersabhängig ist. Denn es st möglich, dass Frauen und Männer ein unterschiedliches Kaufverhalter haben. Genauso ist möglich bzw. sogar wahrscheinlich, dass jüngere onsumenten andere Marken wählen als ältere Konsument d) Die dreidimensionale Häufigkeitsverteilung liefert mehr Information a ille möglichen ein- und zweidimensionalen Häufigkeitsverteilungen usammen. Nur aus ihr kann man erfahren, welchen Zusammenhang es wischen **allen drei Merkmalen** gibt.

Ein Merkmal X mit m Merkmalsausprägungen (x1, ..., xm) wird bei n Untersuchungseinheiten gemessen. Beantworten Sie für die ntsprechende Häufigkeitsverteilung die folgenden Fragen:

- a) Wie groß ist  $\Sigma h(ximi=1)$  Lösung:  $\Sigma h(ximi=1)=n$ b) Wie groß ist  $\Sigma f(ximi=1)$  Lösung:  $\Sigma f(ximi=1)=1$  bzw. 100%
- ) Welchen Wert hat H(xm), F(xm)? Lös. H(xm)=n, F(xm)=1 bzw. 100%

# Vervollständigen Sie die folgende Häufigkeitstabelle (relative Häufigkeiten sind ganze Zahlen!): Lösung der Aufgabe 5

i	Merkmalsausprägung x <sub>i</sub>	h(x <sub>i</sub> )	f(x <sub>i</sub> ) (%)	H(x <sub>i</sub> )	F(x <sub>i</sub> ) (%)
1	3	240	15%	240	15%
2	4	400	25%	640	40%
3	5	320	20%	960	60%
4	6	640	40%	1.600	100%
Σ		1.600	100%		-

ervollständigen Sie die folgende klassierte Häufigkeitstabelle (relative äufigkeiten mit einer Nachkommastelle!) Lösung der Aufsabe 6

F<sub>1</sub> (%) (Mio €) (%) 0 b.u. 40 16,7% 40 b.u. 100 50.0% 200 66.7% 100 b.u. 200 16.7% 250 83.4%

Aufgabe 9 (Klausuraufgabe WS18/19 mit 9 Punkten): Um die Entwicklung der Telefonkosten der letzten 6 Monate des om die Eritwickung der Teierbrinssten der Erzter in Windrade des vergangenen Jahres zu analysieren, wird Claudia von ihrem Vater beauftragt, die mittleren Telefonkosten sowie deren Streuung zu berechnen. Die Telefonkosten (in €) sind in der folgenden Tabelle ıfgeführt. Berechnung mit 2 Nachkommastellen!

Monat	Juli	August	September	Oktober	November	Dezember
Kosten (€)	31,44	30,18	31,04	33,60	38,16	132,40

a) (5 Punkte) Berechnen Sie das arithmetische Mittel sowie die Streuung

=16\*(31,44+30,18+31,04+33,6+38,16+132,4)=296,826=49,47

s2=16\*(31,442+30,182+31,042+33,62+38,162+132,42)-49,472=22977,69 6-2447.28=1382.34 s=Vs2=V1382.34=37.18

b) **(4 Punkte)** Claudia, die im Dezember häufig bei teuren Hotli angerufen hat, ist entsetzt über den hohen Mittelwert und befürchtet Taschengeldentzug durch ihren Vater. Helfen Sie Claudia aus der Patsche, indem Sie ein alternatives Lageparameter, zu Claudias Gunsten, vorschlagen. Begründen Sie Ihren Vorschlag kurz (maximal drei Sätze) und berechnen Sie den Wert Ihres vorgeschlagenen Lagemaßes.

b) (4 Punkte) "Besserer" Vorschlag: Median – weniger empfindlich gegenüber Ausreißern. n gerade, Reihe erst sortieren!

xZ=12(31.44+33.60)=32.52

lm Rahmen einer Bürgerbefragung wurden n = 1.200 Personen u.a. efragt nach ihrem Alter (A) und ihrer Meinung zur Ausweitung der Fußgängerzone in der Innenstadt (F). Das Merkmal Alter wurde wie olgt klassiert: "bis unter 40 Jahre" und "40 Jahre und älter". Auf die rage nach der Ausweitung der Fußgängerzone konnte nur mit "ja" nein" geantwortet werden. Die Befragungsergebnisse wurden uusgezählt. Es ergaben sich die in der folgenden Kreuztabelle aufgeührten absoluten Häufigkeiten für die Merkmalsausprägungskombin ) Tragen Sie in die obige Kreuztabelle ein (absolut und relativ in % mit mmastelle):

1) die heiden Randverteilung

(2) die relativen Spaltenhäufigkeiten in % (3) die relativen Zeilenhäufigkeiten in %

(4) die relativen Häufigkeiten der Merkmalsausprägungskombi. in % Lösung der Aufgabe 7

	F/A⇒	40-	40+	
	th.	bis unter 40 Jahre alt	40 Jahre und älter	Σ
		396	204	(1)
		(2) 60,0%	(2) 37,8%	600
g g	ja	(3) 66,0%	(3) 34,0%	50,0%
ung ierz		(4) 33,0%	(4) 17,0%	
Ausweitung der Fußgängerzone		264	336	(1)
ag a		(2) 40,0%	(2) 62,2%	600
4 E	nein	(3) 44,0%	(3) 56,0%	50,0%
		(4) 22,0%	(4) 28,0%	
		(1)	(1)	
	Σ	660	540	1.200
		55,0%	45,0%	100,0%

## b) Beantworten Sie mit Hilfe der obigen Kreuztabelle die folgender

(3) Wie viel % der Befragten sind bis unter 40 Jahre alt und befürworten die Ausweitung der Fußgängerzone? 33,0.%. (4) Wie viel % der Befragten, die 40 Jahre und älter sind, sind gegen die Ausweitung der Fußgängerzone? .62,2%.

a) Sie lesen in einer Studie über die Altersverteilung in einer Gruppe, dass XZ= 32 Jahre und X = 40 Jahre ist. Welche Schlüsse können Sie

Jass 242 - 34 mile ruin 2 - 40 mile 31. Welche Schmasse komen ind daraus über die Altersverteilung ziehen? b) Sie lesen in einer Studie über die Einkommensverteilung einer Berufsgruper 247 = 30.000 €, ₹ = 40.000 €, Q1 = 25.000 €, Q3 = 45.000 € Welche Informationen erhalten Sie aus diesen 4 statistischen Kennzahlen über die Einkommensverteilung? Erhalten Sie auch Informationen über die Streuung der Verteilung? LÖSUNG:

a) Wegen XZ = 32 < 40 = X, liegt hier eine rechtsschiefe Altersverteilung vor. Es gibt Ausreißer im oberen Altersbereich, die das arithmetische Mittel nach oben ziehen, den Median aber nicht beeinflussen.

b) Die Einkommensverteilung ist rechtsschief, da  $\overline{x}Z < \overline{x}$ . Es gibt Ausreiß im oberen Einkommensbereich. 25% der Personen verdienen weniger and 75% mehr als 25,000 €, 75% verdienen weniger und 25% mehr als 9.000 €. Ein Einkommen von 30.000 € wird von 50% über- und von den anderen 50% unterschritten. Der Interquartilsabstand IQR = Q3 - Q1 = 20.000 € gibt die Spannweite bei den mittleren 50% an und informiert omit über die Streuung des "mittleren" Teils der Einkommensverteilung

Aufgabe 7 (Klausuraufgabe WS17/18 mit 18 Punkten): In der folgenden Tabeile ist die Verteilung der männlichen Teilnehmer bei einer Umfrage auf Altersklassen dargestellt. Dabei wurde zwischen Personen mit und ohne Migrationshintergrund unterschieden

a) Berechnen Sie approximativ die Altersquartile (Q1, Q2, Q3) für die peiden Gruppen. Alle Quartile auf ganze Zahlen runden. TIPP: vervollständigen Sie zuerst die Häufigkeitstabelle. o) Zeichnen Sie die Boxplots für beide Gruppen in einem Diagramm. Bitte

lenken Sie an die "Lesbarkeit" des Boxplots (Achsenbeschriftung und Legende nicht vergessen!) c) Wie groß ist der Anteil der Personen im Alter zwischen 15 und 75 Jahre

für jede Gruppe? (gemeint ist das Intervall [15;75))

Lösung	der	Aufga	be 7

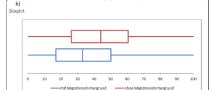
Klassen-Nr.	Alter (von – b.u.) in Jahren	mit Migrations- hintergrund f(x) in %	ohne Migrations- hintergrund f(x) in %	mit Migrations- hintergrund kumm. F(x) in %	ohne Migrations- hintergrund kumm. F(x) in %
1	b.u. 15	22	12	22	12
2	15 - 35	31	23	53	35
3	35 – 55	29	33	82	68
4	55 – 75	15	25	97	93
5	75 und älter	3	7	100	100

a) Median = Q2 →Median ist auch Quartil (s. Folie 38 Modul 4) erechnung von Q1, Q2, Q3 wie Median für klassierte Daten! (s. Folie 13 15, 16 Modul 4)

 $\overline{x}Z := xk-1+(xk-xk-1)*0,5-Fk-1fk$  mit Migrationshintergrund: für Q1 (25%) k=2, QmM=15+(35-15)\*(0,25-0,22)/0,31=17 für Q2 (50%) k=2, Q2mM=15+(35-15)\*(0,5-0,22)/0,31=33 für Q3 (75%) k=3, Q3mM=35+(55-35)\*(0,75-0,53)/0,29=50

ohne Migrationshintergrund: für Q1 (25%) k=2, Q1oM=15+(35-15)\*(0,25-

0,12)/0,23=26 für Q2 (50%) k=3, Q2oM=35+(55-35)\*(**0,5**-0,35)/0,33=44 für Q3 (75%) k=4, Q3oM=55+(75-55)\*(**0,75**-0,68)/0,25=61



c) Mit M = 31 + 29 + 15 = 75% (ODER: 97 – 22 = 75%) →Anteile Personen us den Klassen 2, 3, 4 addieren ODER subtrahieren Anteil Perso aus den Klassen 2, 3, 4 augreten Oben 1111. Klasse 1 von kumuliertem Anteil Personen der Klasse 4. Ohne M = 23 + 33 + 25 = 81% (ODER: 93 - 12 = 81%)

Im Rahmen einer Marktforschungsstudie wurden n = 2.000 Personen u. a. gefragt nach den drei Merkmalen Geschlecht G, Alter A (klassiert in 2 Klassen: "b.u. 40 Jahre" und "40 Jahre und älter") und Markenpräferen // (A = Produkt A, B = Produkt B). Die Befragungsergebnisse sind in der Digenden dreidimensionalen Häufigkeitstabelle zusammengefasst:

	wei	blich	män		
M/GA	bis unter 40 Jahre alt	40 Jahre und älter	bis unter 40 Jahre alt	40 Jahre und älter	Σ
Produkt A	400	150	200	450	1200
Produkt B	100	450	200	50	800
Σ	500	600	400	500	2000

Beantworten Sie auf der Basis der obigen Tabelle die folgenden Fragen mit 1 Nachkommastelle):

init I vachkolmische); a) Wieviel % der jüngeren Frauen (bis unter 40 J) bevorzugen Produkt B? Lösung 100 / 500 = 20,0 % b) Wieviel % der älter. Männer (40 Jahre & älter) bevorzugen Produkt A?

50/500 = 90,0 %

c) Wie viel % der Frauen bevorzugen Produkt A?

50/1100 = 50.0 %

I) Wie viel % der A-Käufer sind Frauen? 50/1200 = 45.8 %

e) Wie viel % der Befragten sind Frauen, unter 40 Jahre & A-Käuferinnen? 400/2000 = 20,0 %

f) Wie viel % der Frauen, die 40 Jahre & älter sind, bevorzugen Produkt A

g) Wie viel % der Befragten sind Männer? 900/2000 = 45 % h) Welchen Wert hat f(B | b.u. 40 J.)

300/900 = 33,3 %

i) Wie viel Prozent der A-Käufer, die 40 Jahre und älter sind, sind Frauen?

j) Wie viel %der Befragten sind A-Käufer, 40 Jahre und älter & Männer? 450/2000 = 22,5 %

Für die 20 Unternehmen einer Branche wurden im Jahr 2005 die folgenden Umsätze (in Mio. €) ermittelt: 35, 150, 190, 20, 8, 74, 44, 89, 25, 12, 17, 5, 22, 10, 13, 150, 47, 65, 49, 55 Erstellen Sie für diese 13.12. 1.7. 12. 10. 13. 130. 17. 13. 13. 13. Etseinen Jet in uiese tattstische Reihe in der folgenden Tabelle eine klassierte Häufig-eitsverteilung für die gegebenen Klassen (mit absoluten und relativen in % mit 1 Nachkommastelle) Häufigkeiten und Summenhäufigkeiten).

Klasse Nr.	Umsatzklasse (Mio €)	Anzahl Filialen h	Anteil f <sub>i</sub> (%)	Hi	F; (%)
1	0 b.u. 20	6	30,0%	6	30,0%
2	20 b.u.50	7	35,0%	13	65,0%
3	50 b.u.100	4	20,0%	17	85,0%
4	100 b.u.200	3	15,0%	20	100,0%
Σ		20	100.0%		-

Auf die Frage "Wie viel Stück des Produktes ABC haben Sie im letzten uie ringe\_wie ver zuck des rioutekes zuc inzuerinse im einzeten ant gekauft?" gab es bei der Hauptuntersuchung unterschiedliche worten zur Zahl der gekauften Stücke von ABC. Die statistische he wurde zusammengefasst in der folgenden Häufigkeitsverteilung: a) Bestimmen Sie den Modus, den Median und das arithmetische Mittel.

b) Bestimmen Sie die Streuungsparameter Spannweite w, Varianz s2, riationskoeffizient v

c) Welche Aussagen können Sie auf der Basis der Werte der Lageparameter über die Form der Verteilung machen?

Lösung	der	Aufgabe	5

i	gekaufte Stückzahl von ABC x <sub>i</sub>	Anzahl der Nennungen $h(x_i)$	Anteil f(x <sub>i</sub> ) (%)	H(x <sub>i</sub> )	F(x <sub>i</sub> ) (%)	$(x_i - \overline{x})^2$	$(x_i - \overline{x})^2 h(x_i)$
1	0	100	20%	100	20%	1,44	144
2	1	300	60%	400	80%	0,04	12
3	2	50	10%	450	90%	0,64	32
4	3	20	4%	470	94%	3,24	64,8
5	4	20	4%	490	98%	7,84	156,8
6	5	5	1%	495	99%	14,44	72,2
7	7	5	1%	500	100%	33,64	168,2
Σ		500	100%				650

a)  $\overline{x}D=1$   $\overline{x}Z=1$ 

=0·100+1·300+2·50+3·20+4·20+5·5+7·5500=600500=1,2

o) w = xmax - xmin = 7-0=7

. Variante (Formel (2), s.Folie 20 Modul 5)

 $2=(\Sigma xi2*f(xi))-\Sigma 2ji=1=(0+1\cdot0,6+4\cdot0,1+9\cdot0,04+16\cdot0,04+25\cdot0,01+49\cdot0,01\cdot1,22=0,6+0,4+0,36+0,64+0,25+0,49-1,44=1,3$ 

. Variante (Formel (1), s. Folie 20 Modul 5)  $s2=1n\Sigma(xi-\overline{x})2*h(xi)=650500ji=1=1,3$   $s=\sqrt{s}2=\sqrt{1,3}=1,14$   $v=s\overline{x}=1,14/1,2=0,95$ 

c) Die Verteilung ist rechtsschief wegen  $\overline{x}Z = 1 < 1,2 = \overline{x}$ 

Die Tabelle zeigt die Zahl der Eheschließungen bzw. die Zahl der hescheidungen je 10.000 Ehen in Deutschland

Jahr	Eheschließungen	Ehescheidungen je 10.000 Ehen	
2001	389.000	198,2	
2000	418.550	194	
1999	430.674	187,7	
1998	417.420	191,4	
1997	422.776	181,2	
1996	427.297	161	
1995	430.534	153,8	
1994	440.244	150	
1993	442.605	135,6	
1992	452.428	104,8	
1991	454.291	104	

a) Bestimme den Modalwert, den Zentralwert und das arithmetische Mittel sowohl von den Eheschließungen als auch von den Ehescheidungen. Was fällt dir an den Ergebnissen auf?

b) Bestimme die Varianz und die Standardabweichung sowohl von den Eheschließungen als auch von den Ehescheidungen.

c) Welche Veranschaulichungsmöglichkeiten für solch einen tabellarischen Zusammenhang hast du bereits kennen gelernt? Wähle zwei davon aus und realisiere sie! Welche Visualisierungsform ist in iesem Fall besonders geeignet bzw. ungeeignet und warum?

a) Worüber informiert .

... die Standardabweichung? Informiert über die (absolute) Streuung einer Verteilung. Die Standardabweichung hat die gleiche Maßeinheit wi

die Beobachtungswerte und die Mittelwerte. b. .. das Quartil (3? informiert darüber, welcher Wert von 75% der Beobachtungswerte einer Verteilung unterschritten wird und von 25% der Beobachtungswerte überschritten wird.

.. der Variationskoeffizient? informiert über die relative Streuung einer

... der Interquartilsabstand? informiert über die Spannweite de

u... ue interquat utsassaturi innimiter uoer uie spanimente der mittleren 50% der Beobachtungswerte einer Verteilung, e... das 5%-Quantill? informiert über den Wert, der von 5% der Beobachtungswerte einer Verteilung unterschritten wird und von 95% der Beobachtungswerte überschritten wird.

... der Median? informiert über den Wert, der von der Hälfte der Beobachtungswerte überschritten und von der anderen Hälfte unterschritten wird.

... der Modus? informiert darüber, welche Merkmalsausprägung(en) in der Häufigkeitsverteilung am häufigsten vorkommen

h. ... die Spannweite?

nformiert über den maximalen Abstand der Merkmalsausprägungen, die n der Häufigkeitsverteilung vorkommen. Sie ergibt sich aus der Diffi wischen dem größten und kleinsten Merkmalswert der Verteilung. ... das Quartil Q1? informiert darüber, welcher Wert von 25% der

Beobachtungswerte einer Verteilung unterschritten wird und von 75% der Beobachtungswerte überschritten wird.

b) Was haben die statistischen Parameter Varianz, Standardabweichung und Variationskoeffizient gemeinsam? Wodurch unterscheiden sich die

Lösung: Varianz, Standardabweichung und Variationskoeffizient sind Streuungsparameter. Mit ihnen lässt sich die Streuung einer Häufigkeitsverteilung charakterisieren. Varianz und Standardabweichung informieren über die absolute Streuung, der Variationskoeffizient über die **relative** Streuung. In einer Statistik über die Einkommen (Jahreseinkommen) vor

leitenden Angestellten im Rechnungswesen lesen Sie: 1. Quartil Q1: 75.000 €, 3. Quartil Q3: 150.000 €, Median **xZ**: 100.000 €, Mittelwert **X**: 140.000 €. Welche der folgenden Aussagen über die nkommensverteilung sind richtig? (Zutreffendes ankreuzen!)

	RICHTIG	FALSCH
(1) Die Einkommensverteilung ist linksschief.		X
(2) 25% der Befragten verdienen weniger als 75.000 €.	X	
(3) 50% der Befragten verdienen zwischen 75.000 und 150.000 €.	X	
(4) 75% der Befragten verdienen mehr als 150.000 €.		X
(5) Die Einkommensverteilung ist symmetrisch.		X
(6) 50% der Befragten verdienen weniger als 140.000 €.		X

Erstellen Sie auf der Basis der folgenden Angaben eine ndimensionale, unklassierte Häufigkeitsverteilung: In einer Stadt haben vier Taxi-Unternehmen jeweils drei Wagen, ein Taxi-Unternehmen hat 26 Wagen. Die übrigen Taxi-Unternehmen in der Stadt sind kleiner und haben weniger Wagen: acht Taxi-Untermehmen haben jeweils nur einen Wagen, sieben haben jeweils zwei Wagen. Fragen auf der Basis der erstellten Tabelle beantworten Fragen (mit 1 Nachkommastelle): a) Wie viele Taxi-Unternehmen gibt es in dieser Stadt?

20 Taxi Unternehmen in dieser Stadt.

b) Wie viele Wagen bieten in dieser Stadt ihre Leistungen an? 60 Wagen stehen in der Stadt zur Verfügung. 1\*8+2\*7+3\*4+26\*1=60

c) Bestimmen / berechnen Sie für die obige Verteilung

 $\overline{x}D$ =1 (da Anz. Unternehmen mit 1 Wagen am höchsten b. den Median

b.  $\overline{x}Z$ =2 (s. Folie 11 Modul 4) Reihe: 11111112222223333326

c. das arithmetische Mittel x=1.8+2.7+3.4+26.120=6020=3

d) Machen Sie Aussagen über die Schiefe der Verteilung rechtsschiefe Verteilung

e) Berechnen Sie die Standardabweichung

s2=1/20\*((1-3)2\*8+(2-3)2\*7+(3-3)2\*4+(26-3)2\*1)=120\*(32+7+0+529): 56820=28.4 s=vs2=v28.4=5.33

f) Wie viel % der Taxi-Unternehmen haben mehr als einen Wagen

g) Wie viel % der Taxi-Unternehmen haben weniger als 3 Wage 40+35=75%

Lösung der Aufgabe 3

i	Anzahl Wagen	Anzahl Taxi- Unternehmen h(x <sub>i</sub> )	f(x <sub>i</sub> ) (%)	H(x <sub>i</sub> )	F(x <sub>i</sub> ) (%)
1	1	8	40,0%	8	40,0%
2	2	7	35,0%	15	75,0%
3	3	4	20,0%	19	95,0%
4	26	1	5,0%	20	100,0%
Σ		20	100,0%		
	•				•

sung der Aufgabe 8

 Der Modalwert der Eheschließungen in Deutschland und nescheidungen je Ehe existiert nicht, denn jeder Wert kommt nur ein inziges Mal vor. Kein Wert der Datenreihe ist ein Modalwert!

Zentralwert=Median → ist der 6. Element der sortierten Reihe für Eheschließungen ₹Z=430.534 für Ehescheidungen ₹Z=1.610.000 *oder* 161 *je* 10.000 *Ehen* 

TEheschließungen=389.000 +···+ 454.29111=4.725.81911=429.619.909

Das arithmetische Mittel entspricht ungefähr dem Zentralwert der Datenmenge sowohl bei den Eheschließungen als auch bei den Ehescheidungen, d.h. dass die Entwicklung der Zahlen gleichmäßig ist und keine Ausreißer aufweist.

b) Berechnung der Varianz über die absolute Häufigkeit sEheschlieSungen2=111 $\Sigma(xi-\overline{x})$ 211i=1=111((389.000-429.620)2+···+(4 54.291-429.620)2)=307.808.635

c) Darstellungsformen: Graph, Säulendiagramm, Stabdiagramm, Balkendiagramm, Kreisdiagramm. Besonders geeignete Darstellungsformen: Graph, Balkendiagramm, Säulendiagramm und Stabdiagramm; ungeeignete Darstellungsform ist Kreisdiagramm.