

Beschrijvende statistiek

Het berekenen en
interpreteren van
tabellen en statistieken

Bregje van Groningen
Connie de Boer

Vierde druk

Boom

**inclusief
website!**

Met behulp van onderstaande unieke activeringscode kunt u toegang krijgen tot www.beschrijvendestatistiek.nl voor extra materiaal. Deze code is persoonsgebonden en gekoppeld aan de 4e druk. Na activering van de code is de website 2 jaar toegankelijk. De code kan tot zes maanden na het verschijnen van een volgende druk geactiveerd worden.

9823-XL-53-LG

Omslagontwerp: Cunera Joosten, Amsterdam

Opmaak binnenwerk: Nu-nique grafische vormgeving, Goor

© 2016 Bregje van Groningen & Connie de Boer | Boom uitgevers Amsterdam

Behoudens de in of krachtens de Auteurswet gestelde uitzonderingen mag niets uit deze uitgave worden veelevoudigd, opgeslagen in een geautomatiseerd gegevensbestand, of openbaar gemaakt, in enige vorm of op enige wijze, hetzij elektronisch, mechanisch, door fotokopieën, opnamen of enige andere manier, zonder voorafgaande schriftelijke toestemming van de uitgever.

Voor zover het maken van reprografische veelevoudigingen uit deze uitgave is toegestaan op grond van artikel 16h Auteurswet dient men de daarvoor wettelijk verschuldigde vergoedingen te voldoen aan de Stichting Reprorecht (Postbus 3051, 2130 KB Hoofddorp, www.reprorecht.nl). Voor het overnemen van (een) gedeelte(n) uit deze uitgave in bloemlezingen, readers en andere compilatiewerken (art. 16 Auteurswet) kan men zich wenden tot de Stichting PRO (Stichting Publicatie- en Reproductierechten Organisatie, Postbus 3060, 2130 KB Hoofddorp, www.cedar.nl/pro).

No part of this book may be reproduced in any form, by print, photoprint, microfilm or any other means without written permission from the publisher.

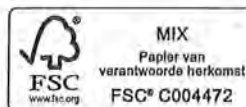
ISBN 9789089539823

ISBN 9789461278432 (e-book)

NUR 916

www.boomuitgeversamsterdam.nl

www.beschrijvendestatistiek.nl



Woord vooraf

Voor veel studenten is statistiek een moeilijk vak. Daarom vinden we het belangrijk dat de uitleg zo duidelijk mogelijk is. Dat stelt hoge eisen aan de leesbaarheid. Ter wille van de leesbaarheid hebben we nauwelijks verwijzingen naar literatuur in de tekst opgenomen. In de literatuurlijst achter in het boek staan de publicaties waarop wij onze kennis hebben gebaseerd. Deze boeken gaan op sommige statistieken veel uitgebreider in. De literatuurlijst kun je daardoor ook beschouwen als aanbevolen literatuur voor een verdere verdieping van de kennis.

Wij denken dat een goed begrip van statistiek pas ontstaat als je weet hoe je statistische gegevens berekent. Daarom is de uitleg mede gebaseerd op de formules die je gebruikt voor een kengetal of statistische analyse. In de praktijk zul je statistieken zelden met de hand, op basis van formules, uitrekenen. Daarvoor bestaan heel geschikte rekenprogramma's, zoals SPSS. Voor het interpreteren van de resultaten is het echter van groot belang dat je weet wat de betekenis is van de statistische gegevens. Dan is het nuttig deze een paar keer zelf te hebben uitgerekend. Daarna kun je ook oefenen met behulp van een computerprogramma. In aparte kaders is uitgelegd hoe je de besproken analyses door SPSS kunt laten uitvoeren. Deze combinatie van uitleg van berekeningen met de hand en berekeningen met SPSS geeft je de mogelijkheid zelf sommen en oefeningen te maken waarbij je de berekeningen met de hand controleert aan de hand van de resultaten van SPSS.

Eerdere versies van dit boek zijn gebruikt in het onderwijs over onderzoeksmethoden en statistiek aan propedeusestudenten Communicatiewetenschap aan de Universiteit van Amsterdam. De reacties en het commentaar van de studenten en de docenten van deze cursussen hebben geleid tot vele verbeteringen van de tekst.

De volgende personen willen wij expliciet noemen als dank voor hun inhoudelijke en tekstuele bijdragen: Tiede Bijlsma, Sanneke Schouwstra, Rob de Lange, Wouter de Nooy, Mieke Sillekens, Reza Kartosen, Floris Müller, Johannes von Engelhardt, Rob Erven, Carel van Wijk, Nadine Bol, Jeroen Jonkman, Marianne Ouwehand, Gert van Driel en Rhianne Hoek. Verder hebben verschillende docenten en studenten van de cursus 'Methoden van Communicatie Onderzoek en Statistiek' commentaar geleverd op eerdere versies van deze tekst.

Nieuw in de vierde druk

In de vierde druk gaan we nog meer in op het belang van de operationalisatie, het kiezen van de juiste meetniveaus en de mogelijke consequenties van keuzes die gemaakt (moeten) worden in het analyseproces. Voorbeelden zijn aangepast aan (op het moment van schrijven) meer actuele zaken.

In hoofdstuk 3 is een aantal paragrafen toegevoegd waarin we onder andere ingaan op de normale verdeling als kansverdeling en *z*-scores bij het berekenen van kansen. Omdat het een boek over *beschrijvende* statistiek is, zullen we dat niet koppelen aan het toetsen van hypothesen en significantie.

Daarnaast is een nieuw hoofdstuk toegevoegd (hoofdstuk 4) waarin het ver- en bewerken van data in SPSS centraal staat, en de consequenties daarvan voor het meetniveau van de variabelen en de analyses die mogelijk zijn. In hoofdstuk 8 wordt uitgebreider dan in de derde druk ingegaan op het berekenen en interpreteren van een variantieanalyse. Hoofdstuk 9 en 10 zijn samengevoegd en uitgebreid. Het heet nu *Schaalconstructie*, en we bespreken daarin zowel de theoretische als statistische overwegingen bij het maken van een valide en betrouwbare schaal. In bijna alle hoofdstukken wordt meer stilgestaan bij de voorwaarden van bepaalde analyses en de consequenties wanneer niet aan de voorwaarden wordt voldaan.

Tot slot is de vierde druk aangepast aan de ontwikkelingen in de SPSS-software. Alle voorbeelden waarin SPSS wordt gebruikt, zijn nu uitgevoerd met SPSS 23.

Bregje van Groningen
Connie de Boer

Inhoud

Woord vooraf	5
Inleiding	13
1 Basiselementen	17
1.1 Datamatrix	17
1.2 Frequentietabellen	20
1.2.1 Hoe ziet dit eruit in SPSS?	21
1.2.2 Missing values	23
1.2.3 Grafieken	23
1.3 Kruistabellen	25
1.4 Variabelen	30
1.5 Meetniveaus	32
1.5.1 Nominaal meetniveau	33
1.5.2 Ordinaal meetniveau	33
1.5.3 Interval meetniveau	34
1.5.4 Ratio meetniveau	35
1.5.5 Criteria	36
1.6 Waarden van variabelen	36
1.6.1 Continue en discrete meetschalen	36
1.7 Univariate, bivariate en multivariate analyses	37
1.7.1 Univariate analyses	38
1.7.2 Bivariate analyses	38
1.7.3 Multivariate analyses	39
2 Centrummaten	41
2.1 Modus	41
2.2 Mediaan	42
2.3 (Rekenkundig) gemiddelde	45
2.4 Keuze tussen centrummaten	49
2.5 Samenvatting	49
3 Spreiding	51
3.1 Kwartielen	52
3.1.1 Boxplot	53
3.2 Variantie	54
3.3 Standaarddeviatie	58
3.4 Centrum- en spreidingsmaten in SPSS	61
3.5 Standaardiseren (z-scores)	63

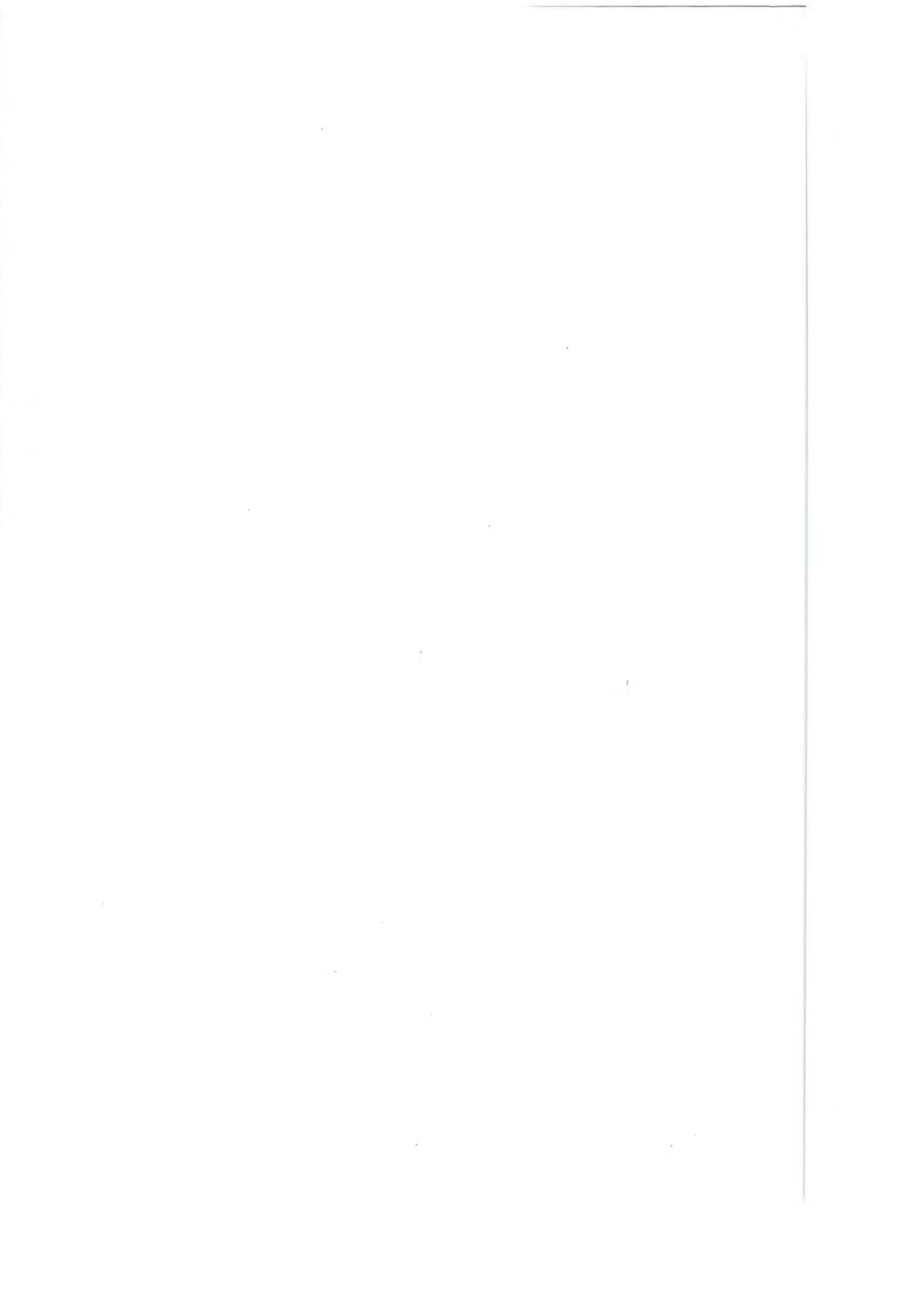
3.6	Normale en scheve verdelingen	66
3.6.1	Normale verdeling	67
3.6.2	Scheve verdelingen	72
3.7	Samenvatting	74
4	Bewerken van je data	77
4.1	Syntax	77
4.2	Missing values	80
4.3	Compute	83
4.4	Hercoderen	87
4.5	Select Cases	91
4.6	Samenvatting	96
5	Associatiematen op nominaal niveau	99
5.1	Wat zijn associatiematen?	99
5.1.1	Meetniveau van de variabelen	99
5.1.2	Symmetrische en asymmetrische relaties	100
5.1.3	Samenhang in kruistabellen	100
5.2	Cramers V	102
5.2.1	Interpretatie	102
5.2.2	Berekening	105
5.3	Phi	108
5.3.1	Interpretatie	108
5.3.2	Berekening	109
5.4	Goodman en Kruskals tau	110
5.4.1	Berekening	111
5.4.2	Interpretatie	116
5.5	Lambda	119
5.5.1	Berekening	121
5.5.2	Interpretatie	123
5.6	Voorwaarden bij het maken van een kruistabel	124
5.7	Samenvatting	125
6	Associatiematen op ordinaal niveau	127
6.1	Samenhang in kruistabellen met ordinaal meetniveau	127
6.2	Gamma	129
6.2.1	Interpretatie	130
6.2.2	Berekening	132
6.3	Somers' d	139
6.3.1	Interpretatie	139
6.3.2	Berekening	141
6.4	Kendalls tau-b	145
6.4.1	Interpretatie	146
6.4.2	Berekening	148
6.4.3	Kendalls tau-b in een correlatiematrix	149

6.5	Spearman's rho	151
6.5.1	Interpretatie	152
6.5.2	Berekening	154
6.6	Samenvatting	158
7	Tabelsplitsing	159
7.1	Interpretatie	160
7.1.1	Spurieuze samenhang	161
7.1.2	Specificatie	166
7.1.3	Versluiting	169
7.2	Samenvatting	171
8	Associatiematen op interval- en rationiveau	173
8.1	Pearson's correlatiecoëfficiënt	173
8.1.1	Grafische weergave	174
8.1.2	Interpretatie	177
8.1.3	Berekening	179
8.1.4	Partiële correlaties	184
8.2	Enkelvoudige regressie	187
8.2.1	Berekening	188
8.2.2	Interpretatie	195
8.3	Meervoudige regressieanalyse	199
8.3.1	Dummyvariabelen	199
8.3.2	Interpretatie meervoudige regressieanalyse	201
8.3.3	Schijnsamenhang in een meervoudige regressie	204
8.3.4	Regressie- en correlatieanalyses in wetenschappelijke tijdschriften	205
8.4	Samenvatting	208
9	Associatiematen: tot slot	211
9.1	Eta en eta-kwadraat	211
9.1.1	Interpretatie	211
9.1.2	Berekening	219
9.1.3	Interactie-effecten bij variantieanalyse	221
9.2	Het kiezen van een associatiemaat	228
9.2.1	Formulering van uitspraken op basis van je onderzoek	229
9.2.2	Kenmerken van de associatiematen	230
9.3	Samenvatting	233
10	Schaalconstructie	235
10.1	Validiteit van een meting	235
10.1.1	Latente en manifeste variabelen	236
10.2	Betrouwbaarheid van een meting	238

10.3	Schaalconstructie	238
10.3.1	Factoranalyse	239
10.3.2	Betrouwbaarheidsanalyse: interne consistentie	247
10.3.3	Maken en beschrijven van de schaal	249
10.4	Meerdere factoren	251
10.4.1	Het vinden van de factoren	255
10.4.2	Het interpreteren van de factoren en de noodzaak van rotatie	256
10.5	Gebruik en presentatie van de resultaten	260
10.6	Overige vormen van betrouwbaarheid	261
10.6.1	Stabiliteit en equivalentie	262
10.7	Samenvatting	264
	Formuleblad Beschrijvende statistiek	265
	Literatuur	273
	Register	275
	Over de auteurs	279

Kaders

Kader 1.1	Invoeren van gegevens en het maken van een frequentietabel	22
Kader 1.2	Het maken van grafieken	24
Kader 1.3	Kruistabellen maken	29
Kader 1.4	Gedrag van mensen is bepalend voor koffiekeuze	30
Kader 2.1	Centrummaten	44
Kader 3.1	Centrum- en spreidingsmaten	61
Kader 3.2	Berekenen van z-scores	65
Kader 4.1	Missing maken van waarden	82
Kader 4.2	Nieuwe variabele maken door middel van Compute	84
Kader 4.3	Herocoderen van variabelen	88
Kader 5.1	Berekenen van nominale associatiematen	121
Kader 6.1	Berekenen van gamma, Somers' d en Kendalls tau-b	141
Kader 6.2	Correlatiematrix met Kendalls tau-b	150
Kader 7.1	Tabelsplitsing	165
Kader 8.1	Het maken van een spreidingsdiagram	176
Kader 8.2	Het berekenen van de correlatie	184
Kader 8.3	Het berekenen van partiële correlaties	186
Kader 8.4	Het uitvoeren van een regressieanalyse	197
Kader 9.1	Vergelijken van gemiddelden tussen afzonderlijke groepen	217
Kader 9.2	Berekenen van interactie-effecten via GLM	222
Kader 10.1	Uitvoeren van een factoranalyse	246
Kader 10.2	Berekenen van Cronbachs alfa	250



Inleiding

In dit boek staat de beschrijvende statistiek centraal. Maar wat is beschrijvende statistiek? Sociale wetenschappers hebben als doel kennis te genereren over de sociale werkelijkheid. Dat kan nieuwe kennis zijn, maar ook kennis waaruit blijkt dat wat we 'wisten' niet of niet helemaal klopt. Kennis kun je verkrijgen door het uitvoeren van een onderzoek. Wanneer je een onderzoek hebt uitgevoerd, wil je de resultaten van het onderzoek zo duidelijk mogelijk weergeven. Dat kan op een korte en overzichtelijke manier gebeuren door middel van kengetallen, tabellen of grafieken. Dit is precies waar het in de beschrijvende statistiek om gaat: het samenvattend beschrijven van de kenmerken van een groep onderzoekseenheden.

Een voorbeeld van een kengetal is het rekenkundig gemiddelde. Door de gemiddelde leeftijd van een groep uit te rekenen geef je een samenvattende beschrijving van een specifiek kenmerk van die groep, namelijk hun leeftijd. Je kunt dit kenmerk ook beschrijven door een tabel of een grafiek te gebruiken, zoals in tabel 1 is gedaan.

Tabel 1 Beschrijving van leeftijd in tabel

Leeftijd	Absolute frequentie	Percentage
18	5	25
19	6	30
20	1	5
21	2	10
22	2	10
24	4	20
Totaal	20	100

Om kengetallen, tabellen en grafieken te produceren moet een onderzoeker een aantal handelingen verrichten. Ten eerste moet hij de gegevens verzamelen en die verzamelde data onderbrengen in een datamatrix. Daarna moet hij de data verwerken (kengetallen bepalen en tabellen en grafieken maken). Deze gegevens vermeldt hij vervolgens in een onderzoeksverslag, waarbij het belangrijk is dat hij de statistieken op de juiste manier weergeeft en interpreteert. Deze fasen, de datamatrix, de data-analyse en de interpretatie van statistieken, worden in dit boek besproken.

In de komende hoofdstukken zal steeds over het algemeen eerst de theoretische achtergrond worden besproken. Het gaat daarbij om de vraag waarom je de analyses nodig hebt. Dit wordt geïllustreerd aan de hand van voorbeelden en waar mogelijk met een grafische weergave. De analyses zijn uitgevoerd met het computerprogramma SPSS, een programma dat in de sociale wetenschap veel wordt gebruikt. In deze vierde druk is gebruikgemaakt van SPSS 23 (Windows-versie) en Windows 10. Het is mogelijk dat de kaders en tabellen er wat betreft lay-out iets anders uitzien wanneer een andere versie van SPSS of Windows-versie, of een Mac is gebruikt. De inhoud is uiteraard steeds hetzelfde. We zullen in aparte kaders uitleggen hoe je met behulp van SPSS de besproken analyses zelf kunt uitvoeren.

In een wetenschappelijke tekst mag je nooit SPSS-tabellen opnemen, ook niet in de bijlagen. Je maakt als onderzoeker zelf tabellen waarin alleen de voor het onderzoek relevante gegevens staan. In dit boek wordt van deze regel afgeweken, omdat we willen laten zien hoe je een SPSS-output kunt lezen en interpreteren. In dit boek zijn de data, wanneer niet zelf verzonnen, afkomstig van databestanden die worden gebruikt door collega's van de afdeling Communicatiewetenschap, scripties van Masterstudenten, en enquêtes die zijn afgenomen door studenten van de module 'Methoden en Technieken voor Communicatieonderzoek'. Veel van deze data zijn bewerkt door de auteurs om de voorbeelden in het boek zo duidelijk mogelijk te maken.

Naast de interpretatie vanuit SPSS zal de handmatige berekening van verschillende analyses worden uitgelegd. Het zelf uitrekenen van statistieken op basis van de formules maakt duidelijk welke betekenis je eraan kunt toekennen. Dit zal het interpreteren van de gegenereerde kengetallen, tabellen en grafieken vergemakkelijken. Deze aspecten – nut en toepassing, voorbeelden, interpretatie en handmatige berekening – zullen in de komende hoofdstukken aan de orde komen bij de behandeling van de verschillende statistieken. Om zo dicht mogelijk in de buurt te komen van de waarden zoals ze in SPSS gepresenteerd zijn, rekenen wij altijd met drie decimalen achter de komma. In teksten over het interpreteren van de maten zullen we, zoals dat ook in de meeste wetenschappelijke artikelen wordt gedaan, ons beperken tot twee decimalen achter de komma.

In de beschrijvende statistiek gebruik je de gegevens van een groep onderzoekseenheden. Die groep kan een *steekproef* of een *populatie* zijn. De gebruikte statistieken zeggen dan iets over de samenstelling van die steekproef of over die populatie, afhankelijk van welke gegevens je gebruikt. Als het om een steekproef gaat, zegt de beschrijvende statistiek dus iets over de steekproef en niet over de populatie waaruit die steekproef afkomstig is. *Beschrijvende* statistiek onderscheidt zich van de *inferentiële* statistiek doordat je je bij de beschrijvende statistiek beperkt tot uitspraken over een groep onderzoekseenheden waarvan je de gegevens hebt. Wanneer je niet de hele populatie hebt onderzocht, kun je over de populatie niets met zekerheid zeggen. Op basis van de steekproefgegevens kun je echter wel met een bepaalde mate van waarschijnlijkheid uitspraken

doen over de populatie. Dit gebeurt in de *inferentiële statistiek*. Op basis van steekproefgegevens maak je dan een schatting van een populatiewaarde.

De statistieken of kengetallen uit de beschrijvende en de inferentiële statistiek geef je aan met symbolen, met letters. Zo duid je het gemiddelde in de populatie aan met een μ (mu), en een gemiddelde in een steekproef met een \bar{x} (x streep) of in een wetenschappelijk artikel met M (mean). De standaarddeviatie in een populatie geef je aan met een σ (sigma). Als je de standaarddeviatie van de steekproef (s) gebruikt als schatter van de standaarddeviatie in de populatie, ben je met inferentiële statistiek bezig. De formules van σ (standaarddeviatie van de populatie) en s (standaarddeviatie van de steekproef) vertonen een klein verschil. Bij σ wordt door n (aantal onderzoekseenheden) gedeeld, en bij s door $n - 1$ (aantal vrijheidsgraden¹). Als je door n zou delen, dan schat je de standaarddeviatie in de populatie iets te laag. Ditzelfde geldt niet alleen voor de standaarddeviatie, maar ook voor andere kengetallen die in dit boek worden behandeld.

Bij de behandeling van de beschrijvende statistiek zou het logisch zijn om te kiezen voor de kengetallen waarbij door n wordt gedeeld. De statistieken worden immers niet gebruikt als schatter voor de populatieparameters. Hier hebben we niet voor gekozen, omdat we naast het met de hand uitrekenen van statistische gegevens aandacht besteden aan SPSS. Je kunt de met de hand uitgerekenen statistieken controleren door de gegevens in een SPSS-bestand in te voeren. Het is ook mogelijk om op basis van een SPSS-bestand zelf oefeningen te maken voor het met de hand uitrekenen van statistieken. Aangezien je SPSS ook gebruikt voor de inferentiële statistiek, zou er – zeker als het gaat om een klein aantal onderzoekseenheden – een verschil kunnen bestaan tussen de output van SPSS en de zelf uitgerekenen gegevens als je de formules uit de beschrijvende statistiek gebruikt. Daarom hebben we ervoor gekozen de formules te gebruiken die horen bij de inferentiële statistiek.

In hoofdstuk 1 staan de basiselementen van de beschrijvende statistiek centraal. De centrummaten (modus, mediaan en gemiddelde) en spreidingsmaten worden respectievelijk in hoofdstuk 2 en 3 beschreven. In hoofdstuk 3 wordt ook de normale verdeling als kansverdeling besproken. Hoofdstuk 4 geeft informatie over hoe je in SPSS je data kunt bewerken. De hoofdstukken 5, 6, 8 en 9 gaan over de associatiematen op nominaal, ordinaal en interval- en rationiveau. In hoofdstuk 7 wordt aandacht besteed aan tabelsplitsing. Hoofdstuk 10 gaat over schaalconstructie en behandelt factoranalyse en betrouwbaarheidsanalyses.

Website



Bij dit boek hoort een website met aanvullend materiaal: www.beschrijvendestatistiek.nl. Op deze website staan per hoofdstuk opdrachten. Studenten kunnen zo oefenen met de lesstof van het hoofdstuk. Op pagina 4 van dit boek staat een persoonlijke inlogcode voor toegang tot de website.

Noot

- 1 Dit is het aantal keren dat een waarde 'vrijelijk kan variëren'. Als je van een groep mensen de gemiddelde leeftijd weet, en je de leeftijden weet van alle individuen uit die groep op één na ($n-1$), staat de leeftijd van die laatste persoon vast. Die kan niet meer variëren. De leeftijd van deze laatste persoon kun je precies uitrekenen op basis van de leeftijden van de anderen uit de groep en de gemiddelde leeftijd.