

# BESCHRIJVENDE STATISTIEK EN KANSREKENEN

PETER  
GOOS



P. GOOS

BESCHRIJVENDE  
STATISTIEK  
EN KANSREKENEN

Acco Leuven / Den Haag

*Eerste druk*: 2013

Gepubliceerd door

Uitgeverij Acco, Blijde Inkomststraat 22, 3000 Leuven (België)

E-mail: [uitgeverij@acco.be](mailto:uitgeverij@acco.be) – Website: [www.uitgeverijacco.be](http://www.uitgeverijacco.be)

*Voor Nederland*

Acco Nederland, Westvlietweg 67 F, 2495 AA Den Haag, Nederland

E-mail: [info@uitgeverijacco.nl](mailto:info@uitgeverijacco.nl) – Website: [www.uitgeverijacco.nl](http://www.uitgeverijacco.nl)

Omslagontwerp: [www.frisco-ontwerpbureau.be](http://www.frisco-ontwerpbureau.be)

De uitgever heeft ernaar gestreefd de auteursrechten te regelen volgens de wettelijke bepalingen. Degenen die desondanks menen zekere rechten te kunnen doen gelden, kunnen zich alsnog tot de uitgever wenden.

© 2013 by Acco (Academische Coöperatieve Vennootschap cvba), Leuven (België)

Niets uit deze uitgave mag worden verveelvoudigd en/of openbaar gemaakt door middel van druk, fotokopie, microfilm of op welke andere wijze dan ook zonder voorafgaande schriftelijke toestemming van de uitgever.

No part of this book may be reproduced in any form, by mimeograph, film or any other means without permission in writing from the publisher.

# Woord vooraf

Dit boek is het resultaat van een grondige herwerking van de cursustekst “Beschrijvende Statistiek en Kansrekenen”, die ik ontwikkelde voor het opleidingsonderdeel “Statistiek met (bedrijfs)economische toepassingen 1” voor studenten Handelsingenieur, Handelsingenieur in de Beleidsinformatica en Toegepaste Economische Wetenschappen aan de Faculteit Toegepaste Economische Wetenschappen van de Universiteit Antwerpen. In vergelijking met eerder gepubliceerde versies van de tekst is een eerste grote wijziging dat een aanzienlijk aantal nieuwe onderwerpen aan bod komen, zoals de nominale en ordinale spreidingsindices, kurtosis of gepiektheid, de rangcorrelatiecoëfficiënt, en de gamma, de Weibull-, de lognormale en de beta kansdichtheden. Ook de multinomiale kansverdeling en de multivariate normale kansdichtheid worden nu behandeld. Daarnaast bevat de herwerkte versie van de tekst een aanzienlijk groter aantal figuren en tabellen, met als doel meer inzicht te verschaffen in de concepten en de formules die aangehaald worden. Ten slotte maakt dit boek gebruik van het statistisch softwarepakket JMP (spreek uit als “jump”), in tegenstelling tot eerdere versies van de tekst die steunden op het gebruik van Microsoft Excel.

Het grote voordeel van JMP is dat het gebruiksvriendelijk en krachtig is, uitstekende grafische mogelijkheden biedt (inclusief het maken van kaarten met statistische informatie), en geschikt is voor alle statistisch georiënteerde opleidingsonderdelen in de bachelorprogramma’s Handelsingenieur, Handelsingenieur in de Beleidsinformatica en Toegepaste Economische Wetenschappen. De bedoeling is om hetzelfde pakket doorheen het ganse programma te gebruiken. Aan veel onderwijsinstellingen is JMP beschikbaar voor lesgevers en studenten. Dit is ook zo aan de Universiteit Antwerpen: elke lesgever en student kan JMP downloaden en gebruiken in pc-classes, maar ook op zijn of haar eigen computer(s) thuis. Deze versie werkt gegarandeerd het ganse academiejaar. Je kunt ook te allen tijde een gratis proefversie van JMP downloaden ([www.jmp.com](http://www.jmp.com)). Je moet hiertoe wel eerst een korte registratieprocedure doorlopen. De proefversie werkt gedurende 30 dagen. Voor onderwijsinstellingen zijn campuslicenties van JMP overigens heel goedkoop.

Hoe gebruiksvriendelijk JMP ook is, het gebruik ervan zal enige oefening vergen. Om vertrouwd te raken met het softwarepakket kun je enkele *webcasts* bekijken op [www.jmp.com](http://www.jmp.com). De belangrijkste *webcast* is misschien wel de *on-demand webcast* met als titel “JMP for Students 1: Navigation and Use”. Je kunt deze *webcast* zien door te klikken op “News and Events”, vervolgens op “On-Demand Webcasts” en ten slotte op het tabblad “Academic” te klikken. Ook de *webcast* “JMP for Students 2: Basic Statistics” is het bekijken waard, omdat er diverse tips worden gegeven over het genereren van grafieken en het gebruik van de helpfunctie in JMP.



Het boek veronderstelt dat de studenten een grondige wiskundige voorkennis hebben overgehouden uit het middelbaar onderwijs, inclusief het gebruik van afgeleides, integralen en matrixalgebra. Een sterkte van het boek is dat alle wiskundige afleidingen gedetailleerd zijn weergegeven. Dit betekent dat alle benodigde tussenstappen in de wiskundige afleidingen zijn weergegeven, ook diegene die voor wiskundigen misschien triviaal zijn. De bedoeling hiervan is om de drempel voor minder wiskundig begaafde studenten niet hoger te maken dan strikt noodzakelijk.

De leerstof in dit boek wordt behandeld tijdens een 35-tal uur hoorcolleges, waarin de concepten uitgelegd worden, en 24 uur werkcolleges, waarin oefeningen gemaakt worden. Daarnaast worden ook sessies georganiseerd over het gebruik van het softwarepakket JMP.

Het boek onderscheidt zich van veel moderne tekstboeken over beschrijvende statistiek en kansrekenen doordat het een combinatie biedt van theoretische en wiskundige diepgang, gedetailleerde en begrijpbare uitleg bij de geïntroduceerde concepten, tal van praktische voorbeelden, en het gebruik van een gebruiksvriendelijk en toch bijzonder krachtig statistisch pakket.

Bij de uitgave van dit boek wil de auteur graag Leonids Aleksandrovs, Kris Annaert, Stefan Becuwe, Filip De Baerdemaeker, Roselinde Kessels, Ida Ruts, Bagus Sartono, Evelien Stoffels, Anja Struyf, Anil Haydar Topal, Katrien Van Driessen, Ellen Vandervieren, Kristel Van Rompaey, Diane Verbiest en Sara Weyns bedanken voor hun gedetailleerde opmerkingen en constructieve suggesties, en voor de technische ondersteuning bij het creëren van de figuren. De auteur wenst ook Bradley Jones, Volker Kraft, Brady Brady en Mia Stephens van de JMP-divisie in het SAS Institute te bedanken voor hun hulp met betrekking tot het softwarepakket JMP. Ten slotte wil de auteur ook professor dr. Willy Gochet van de KU Leuven bedanken voor de inspiratie die hij bood bij de ontwikkeling van dit boek.

# Inhoud

<b>1</b>	<b>Wat is statistiek?</b>	<b>13</b>
1.1	Waarom statistiek? . . . . .	13
1.2	Definitie van statistiek . . . . .	15
1.3	Voorbeelden . . . . .	16
1.4	Onderwerp van de statistiek . . . . .	17
1.5	Kansrekening . . . . .	19
1.6	Software . . . . .	21
<b>2</b>	<b>Data en hun voorstelling</b>	<b>23</b>
2.1	Soorten gegevens en meetschalen . . . . .	23
2.1.1	Categorische of kwalitatieve variabelen . . . . .	23
2.1.2	Kwantitatieve variabelen . . . . .	24
2.1.3	Hiërarchie van meetschalen . . . . .	25
2.1.4	Meetschalen in JMP . . . . .	26
2.2	De datamatrix . . . . .	26
2.3	Voorstellen van univariate kwalitatieve variabelen . . . . .	26
2.4	Voorstellen van univariate kwantitatieve variabelen . . . . .	31
2.4.1	Stam- en bladdiagram . . . . .	31
2.4.2	Naalddiagram voor univariate discrete kwantitatieve variabelen . . . . .	33
2.4.3	Histogrammen en frequentiepolygonen voor continue variabelen . . . . .	37
2.4.4	Empirische cumulatieve verdelingsfuncties . . . . .	43
2.5	Het voorstellen van bivariate variabelen . . . . .	44

2.5.1	Kwalitatieve variabelen . . . . .	44
2.5.2	Kwantitatieve variabelen . . . . .	49
2.6	Het voorstellen van tijdreeksen . . . . .	53
2.7	Het gebruik van kaarten . . . . .	54
2.8	Nog meer grafische mogelijkheden . . . . .	62
<b>3</b>	<b>Beschrijvende statistieken van steekproefgegevens</b>	<b>69</b>
3.1	Kengetallen van centrale ligging of locatie . . . . .	70
3.1.1	Mediaan . . . . .	70
3.1.2	Modus . . . . .	72
3.1.3	Rekenkundig gemiddelde . . . . .	73
3.1.4	Meetkundig of geometrisch gemiddelde . . . . .	77
3.2	Maatstaven van relatieve ligging . . . . .	78
3.2.1	Ordestatistiek, kwantiel, percentiel, deciel . . . . .	78
3.2.2	Kwartiel . . . . .	79
3.3	Kengetallen van spreiding . . . . .	80
3.3.1	Spreidingsbreedte . . . . .	80
3.3.2	Interkwartielbreedte . . . . .	80
3.3.3	Gemiddelde absolute afwijking . . . . .	81
3.3.4	Variantie . . . . .	81
3.3.5	Standaarddeviatie . . . . .	84
3.3.6	Variatiecoëfficiënt . . . . .	85
3.3.7	Spreidingsindices voor nominale en ordinale variabelen . . . . .	86
3.4	Kengetallen van scheefheid . . . . .	91
3.5	Gepiektheid of kurtosis . . . . .	94
3.6	Transformatie en standaardisatie van gegevens . . . . .	94
3.7	Boxplot . . . . .	95
3.8	Bivariate variabelen . . . . .	99
3.8.1	Covariantie . . . . .	100
3.8.2	Correlatie . . . . .	104

3.8.3	Rangcorrelatie . . . . .	105
3.9	Complementariteit van kengetallen en grafische voorstellingen . . . . .	111
3.10	Beschrijvende statistiek met behulp van JMP . . . . .	112
<b>4</b>	<b>Kansrekenen</b>	<b>119</b>
4.1	Kansexperimenten . . . . .	121
4.2	Definitie van kans . . . . .	123
4.3	Rekenregels . . . . .	126
4.4	Voorwaardelijke kans . . . . .	127
4.5	Onafhankelijke en afhankelijke gebeurtenissen . . . . .	132
4.6	Totale kans en de regel van Bayes . . . . .	136
4.7	Het simuleren van kansexperimenten . . . . .	140
<b>5</b>	<b>Bijkomende aspecten van kansrekening</b>	<b>143</b>
5.1	Combinatieleer . . . . .	143
5.1.1	Optelregel . . . . .	143
5.1.2	Vermenigvuldigingsprincipe . . . . .	144
5.1.3	Permutaties . . . . .	145
5.1.4	Combinaties . . . . .	145
5.2	Aantal mogelijke volgordes . . . . .	146
5.2.1	Twee verschillende objecten . . . . .	147
5.2.2	Meer dan twee verschillende objecten . . . . .	148
5.3	Toepassingen van kansrekenen . . . . .	148
5.3.1	Reeksen van onafhankelijke kansexperimenten . . . . .	148
5.3.2	Euromillions . . . . .	150
<b>6</b>	<b>Univariate kansvariabelen</b>	<b>153</b>
6.1	Kansvariabelen en verdelingsfuncties . . . . .	153
6.2	Discrete kansvariabelen en kansverdelingen . . . . .	155
6.3	Continue kansvariabelen en kansdichtheden . . . . .	157
6.4	Functies van kansvariabelen . . . . .	167



6.4.1	Functies van een discrete kansvariabele . . . . .	167
6.4.2	Functies van een continue kansvariabele . . . . .	167
6.5	Families van kansverdelingen en kansdichtheden . . . . .	169
6.6	Simulatie van kansvariabelen . . . . .	170
<b>7</b>	<b>Kengetallen van populaties en processen</b>	<b>175</b>
7.1	Verwachte waarde van een kansvariabele . . . . .	175
7.2	Verwachte waarde van een functie van een kansvariabele . . . . .	178
7.3	Speciale gevallen . . . . .	179
7.4	Variantie en standaarddeviatie van een kansvariabele . . . . .	180
7.5	Andere kengetallen . . . . .	183
7.6	Momentgenererende functie . . . . .	185
<b>8</b>	<b>Belangrijke discrete kansverdelingen</b>	<b>189</b>
8.1	De uniforme verdeling . . . . .	189
8.2	De Bernoulli-verdeling . . . . .	191
8.3	De binomiale verdeling . . . . .	192
8.3.1	Kansverdeling . . . . .	192
8.3.2	Verwachte waarde en variantie . . . . .	198
8.4	De hypergeometrische verdeling . . . . .	200
8.5	De Poisson-verdeling . . . . .	204
8.6	De geometrische verdeling . . . . .	211
8.7	De negatief binomiale verdeling . . . . .	214
8.8	Kansen en kansverdelingen in JMP . . . . .	216
8.8.1	Tabellen met kansverdelingen en cumulatieve verdelingsfuncties . . . . .	217
8.8.2	Grafische voorstellingen . . . . .	218
8.9	Het simuleren van discrete kansvariabelen met JMP . . . . .	224
<b>9</b>	<b>Belangrijke continue kansdichtheden</b>	<b>231</b>
9.1	De continue uniforme dichtheid . . . . .	232
9.2	De exponentiële dichtheid . . . . .	233

9.2.1	Definitie en kengetallen . . . . .	234
9.2.2	Enkele interessante eigenschappen . . . . .	235
9.3	De gamma dichtheid . . . . .	239
9.4	De Weibull-dichtheid . . . . .	240
9.5	De beta dichtheid . . . . .	242
9.6	Andere dichtheden . . . . .	242
9.7	Grafische voorstellingen en kansberekeningen in JMP . . . . .	245
9.8	Het simuleren van continue kansvariabelen in JMP . . . . .	249
<b>10</b>	<b>De normale verdeling</b>	<b>251</b>
10.1	De dichtheid . . . . .	252
10.2	Berekening van kansen voor normaal verdeelde variabelen . . . . .	256
10.2.1	Standaardnormaal verdeelde variabelen . . . . .	256
10.2.2	Normaal verdeelde variabelen . . . . .	257
10.2.3	JMP . . . . .	258
10.2.4	Voorbeelden . . . . .	259
10.3	Lognormale kansdichtheid . . . . .	265
<b>11</b>	<b>Multivariate kansvariabelen</b>	<b>271</b>
11.1	Inleidende begrippen . . . . .	271
11.2	Gezamenlijke (discrete) kansverdeling . . . . .	273
11.3	Marginale of onvoorwaardelijke (discrete) kansverdeling . . . . .	274
11.4	Voorwaardelijke (discrete) kansverdeling . . . . .	276
11.5	Voorbeelden met discrete bivariate kansvariabelen . . . . .	277
11.6	De multinomiale kansverdeling . . . . .	284
11.7	Gezamenlijke (continue) kansdichtheid . . . . .	286
11.8	Marginale of onvoorwaardelijke (continue) kansdichtheid . . . . .	293
11.9	Voorwaardelijke (continue) kansdichtheid . . . . .	297
<b>12</b>	<b>Functies van meerdere kansvariabelen</b>	<b>299</b>
12.1	Een functie van meerdere kansvariabelen . . . . .	299

12.2	Verwachte waarde van functies van meerdere kansvariabelen . . . . .	300
12.3	Voorwaardelijke verwachte waarden . . . . .	305
12.4	Kansverdeling van functies van kansvariabelen . . . . .	306
12.4.1	Discrete kansvariabelen . . . . .	306
12.4.2	Continue kansvariabelen . . . . .	307
12.5	Functies van onafhankelijke Poisson, normaal en lognormaal verdeelde kansvariabelen . . . . .	312
<b>13</b>	<b>Covariantie, correlatie en variantie van lineaire functies</b>	<b>317</b>
13.1	Covariantie en correlatie . . . . .	317
13.2	Variantie van een lineaire functie van twee kansvariabelen . . . . .	321
13.3	Variantie van een lineaire functie van meer dan twee kansvariabelen . . . . .	323
13.4	Variantie van een lineaire combinatie van onafhankelijke kansvariabelen . . . . .	324
13.4.1	Twee onafhankelijke kansvariabelen . . . . .	324
13.4.2	Meerdere onderling onafhankelijke kansvariabelen . . . . .	324
13.5	Lineaire combinatie van normaal verdeelde kansvariabelen . . . . .	325
13.6	Bivariate en multivariate normale kansdichtheid . . . . .	326
13.6.1	Bivariate normale kansdichtheid . . . . .	326
13.6.2	Grafische voorstellingen . . . . .	327
13.6.3	Onafhankelijkheid, marginale en voorwaardelijke dichtheden . . . . .	330
13.6.4	Algemene multivariate normale kansdichtheid . . . . .	335
<b>14</b>	<b>De centrale limietstelling</b>	<b>337</b>
14.1	Kansdichtheid van het steekproefgemiddelde uit een normaal verdeelde populatie	337
14.2	Kansverdeling of -dichtheid van het steekproefgemiddelde uit een niet-normaal verdeelde populatie . . . . .	338
14.2.1	Centrale limietstelling . . . . .	338
14.2.2	Illustratie van de centrale limietstelling . . . . .	340
14.3	Toepassing . . . . .	345
14.4	Normale benadering van de binomiale verdeling . . . . .	346
	<b>Appendix A. Het Griekse alfabet</b>	<b>349</b>

## INHOUD

Appendix B. Binomiale verdeling	351
Appendix C. Poisson-verdeling	357
Appendix D. Exponentiële verdeling	361
Appendix E. Standaardnormale verdeling	363





# Hoofdstuk 1

## Wat is statistiek?

*The world is ready for the truth; the modern age is here; every year another report appears that examines poverty by means of statistical research rather than romantic claptrap. (uit The Crimson Petal and the White, Michel Faber, p. 334)*

In dit inleidende hoofdstuk wordt een algemene omschrijving gegeven van de onderwerpen statistiek en kansrekenen. Enkele voorbeelden illustreren het doel en de toepassingsmogelijkheden van beide disciplines, alsook het verschil ertussen. Aangezien binnen de bedrijfswereld, de industrie, het management en de economie de statistiek meer toepassingen vindt dan de kansrekening, krijgt statistiek in de opleidingen Handelsingenieur, Toegepaste Economische Wetenschappen en Sociaal-Economische Wetenschappen veruit de meeste aandacht. Desondanks wordt tijdens deze opleidingen ook de nodige aandacht besteed aan de kansrekening. Beide disciplines kunnen immers niet van elkaar losgekoppeld worden. In dit boek komt zowel kansrekenen als statistiek aan bod.

### 1.1 Waarom statistiek?

Statistiek is sinds jaar en dag een vak, niet zelden een gevreesd vak, in een groot aantal studierichtingen aan universiteiten en hogescholen. De reden hiervoor is dat heel wat mensen tijdens hun beroepswerkzaamheden vroeg of laat met problemen van gegevensanalyse geconfronteerd worden. Een degelijke statistische achtergrond laat hen niet enkel toe om de gegevens te analyseren en op basis van de analyse concrete beslissingen te nemen, maar het geeft hen ook een voorsprong bij het verzamelen van gegevens.

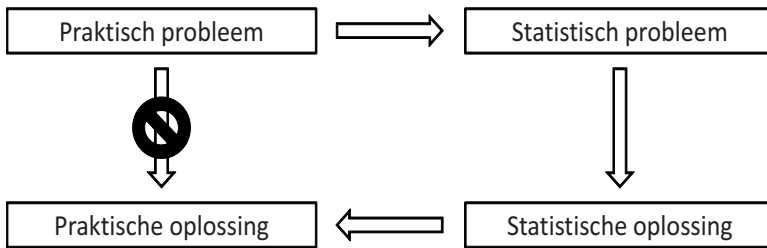
Dat statistiek desondanks door de meeste studenten niet onmiddellijk als nuttig wordt ervaren is veelal te wijten aan het feit dat zij bij het volgen van een lessenspakket statistiek nog niet in aanraking gekomen zijn met allerlei praktische beslissingsproblemen waarmee managers, economen, ingenieurs en onderzoekers dagelijks geconfronteerd worden. Typisch gaan studenten pas bij het uitvoeren van hun thesiswerkzaamheden beseffen dat statistiek ook voor hen

nuttig is. De vele voorbeelden in dit boek zijn bedoeld om deze bewustwording een aantal jaar te vervroegen.

Doorgaans worden bij het introduceren van een cursus statistiek een resem citaten uit de kast gehaald in een poging de studenten te motiveren. Een klassiek voorbeeld is “*Statistical thinking will one day be as necessary for efficient citizenship as the ability to read and write*” van de Britse filosoof Herbert George Wells (1866-1946). Meer recent is de uitspraak van de Amerikaanse kwaliteitsgoeroe W. Edwards Deming, aan wie een groot gedeelte van de ronduit spectaculaire economische heropleving van Japan na de Tweede Wereldoorlog wordt toegeschreven. Hij beweerde dat “*Statistics is too important to be left to statisticians. The goal is to have many statistically-skilled workers: engineers, scientists, managers, ...*” Hal Varian, *chief economist* bij Google beweert het volgende: “*I keep saying that the most sexy job in the next 10 years will be statistician. And I’m not kidding.*” Dichter bij huis beweert Willy Buysse, topman bij SN Brussels Airlines, dat veel te weinig beslissingen genomen worden op basis van gegevens. Recentelijk werd zijn jarenlange ijver — eerst binnen het ter ziele gegane SABENA en later bij de opvolger SN Brussels Airlines — voor de oprichting van een onderzoeksafdeling, waarin statistische en andere kwantitatieve methoden aangewend worden om allerlei problemen aan te pakken, beloond.

Een andere motivering voor het aanleren van statistische methoden aan economen is te vinden in het zogenaamde Six Sigma verbeterprogramma. De bedoeling van dit programma is om binnen zowel dienstverlenende als industriële ondernemingen concrete problemen met een grote financiële impact op te lossen en op die manier het aantal fouten en defecten te herleiden tot 3.4 per miljoen operaties. De aanpak is gebaseerd op statistische methoden, zoals voorgesteld in Figuur 1.1. De figuur geeft aan dat de traditionele werkwijze om praktische problemen onmiddellijk met praktische oplossingen te verhelpen verlaten wordt. Deze aanpak is immers vaak gebaseerd op gegis, nattevingerwerk en trial-and-error, met als gevolg dat een definitieve oplossing niet zelden lang op zich laat wachten. Het Six Sigma verbeterprogramma promoot een meer doordachte, wetenschappelijke aanpak van problemen. Eerst worden gegevens verzameld in de zogenaamde meetfase. Met behulp van statistische methoden worden de gegevens vervolgens onderzocht, wat dikwijls leidt tot interessante inzichten en aanbevelingen ter verbetering van bestaande producten, diensten of processen. De Six Sigma aanpak steunt ook op het gebruik van statistische procescontrole en het uitvoeren van experimenten. Op die manier helpt de statistiek de best mogelijke oplossing te vinden voor allerlei praktische problemen.

Om tot een geslaagd huwelijk te komen tussen enerzijds praktijkmensen en anderzijds statistici wordt van beide kanten vanzelfsprekend enige openheid gevergd. Van de kant van ingenieurs, economen en managers is een degelijke basiskennis van de basisprincipes en -technieken van de statistiek vereist. De statistiek is bijgevolg een onmisbaar facet in de waaier van vaardigheden die van een polyvalente medewerker verwacht worden. Dit verklaart waarom niet enkel in eerste en tweede bachelor Handelsingenieur, Handelsingenieur in de Beleidsinformatica, Toegepaste Economische Wetenschappen en Sociaal-Economische Wetenschappen statistiek gedoceerd wordt, maar ook in de latere jaren van de opleidingen.



**Figuur 1.1:** Gebruik van statistische methoden om problemen op te lossen.

Ten slotte is een grondige opleiding in de statistiek ook een *conditio sine qua non* voor studenten Beleidseconomie. Ook zij zullen in hun beroepsloopbaan met tal van datasets geconfronteerd worden. Een grondige studie van deze gegevens is onmogelijk zonder statistische achtergrond. Voor hen vormen de cursussen statistiek in de opleiding een opstapje naar econometrische onderzoeksmethoden, die daarbij doorgaans gehanteerd worden.

## 1.2 Definitie van statistiek

Het woord statistiek klinkt iedereen wellicht vertrouwd in de oren. Een statistiek verwijst doorgaans naar numerieke informatie, bijvoorbeeld informatie omtrent

- de bevolking van een land: geboorte- en sterftcijfers, immigraties en emigraties, ... (deze informatie noemt men bevolkingsstatistieken);
- de economie: tewerkstellings- en werkloosheidscijfers, investeringen, prijzen, bruto nationaal product (BNP), ... (deze statistieken noemt men economische statistieken); of
- een bedrijf of sector: verkoopcijfers, resultatenrekening, groei, aanwervingen, afvloeiingen, ... (deze cijfers noemt men bedrijfsstatistieken).

Meer formeel kan statistiek gedefinieerd worden als het geheel van methodologieën voor het verzamelen, voorstellen, analyseren en interpreteren van data of gegevens. Hieruit blijkt dat de statistiek een heel algemene hulpwetenschap is, waarvoor in nagenoeg elke werkomgeving een belangrijke rol is weggelegd. Toepassingen van de statistiek in de geneeskunde, de economie, de scheikunde en het bedrijfsmanagement zijn legio, maar ook in de literatuurwetenschap, geschiedenis, politieke wetenschappen, criminologie en zelfs musicologie wordt statistiek gebruikt.

Data of gegevens zijn in de moderne maatschappij massaal aanwezig:

- Computerbestanden in bedrijven bevatten verkoopcijfers, kostprijzen en klantgegevens (zoals adres, bestelde hoeveelheden en frequentie van bestellen).



- De financiële pagina's van kranten bevatten aandelenkoersen, grondstofprijzen en wisselkoersen.
- De federale en regionale overheden publiceren regelmatig data over de bevolking, handel en industrie.
- Het internet is een bron van tal van datasets.

Ondernemingen verzamelen uiteraard ook zelf actief gegevens. Dit gebeurt onder meer door het uitvoeren van experimenten (bijvoorbeeld om nieuwe producten te ontwerpen), in het kader van statistische procescontrole, of door het meten van allerlei eigenschappen van producten, diensten en processen. Kwaliteitsafdelingen van ondernemingen trachten door het voortdurend analyseren van gegevens producten of diensten af te leveren met zo weinig mogelijk defecten en een maximale betrouwbaarheid. Bovendien wordt gepoogd om het bedrijfsproces zodanig te organiseren dat de afvalberg minimaal is, er weinig tot geen inspectie van afgewerkte producten nodig is, en dat de producten en diensten met een minimum aan kosten aan de klantenvereisten voldoen.

Onderzoeksbureaus verzamelen gegevens via enquêtes per telefoon, per post of door straatinterviews. Dergelijke enquêtes worden opgezet om informatie in te winnen over het winkelgedrag van consumenten, over het kiesgedrag van de bevolking of over de publieke opinie betreffende maatschappelijke problemen.

Statistiek laat ons toe om gegevens of data te verwerken tot bruikbare informatie. De rol die de statistiek hierin speelt kan het best geïllustreerd worden aan de hand van enkele voorbeelden.

## 1.3 Voorbeelden

### Voorbeeld 1.1

Een luchtvaartmaatschappij voerde een onderzoek uit naar het gedrag van haar passagiers op intercontinentale vluchten en registreerde daarom

- het aantal passagiers met reservatie dat niet opdaagt (de zogenaamde *no-shows*);
- het gewicht aan bagage dat de passagiers meenemen (vaak geldt een limiet van 20 kilogram); en
- de tijd die de passagiers aankomen vóór het officiële vertrek van de vlucht (voor intercontinentale vluchten wordt de passagiers gevraagd minimaal twee uur voor vertrek aanwezig te zijn).

De maatschappij registreerde gedurende enkele maanden deze gegevens en maakte hierbij een onderscheid tussen de passagiers die *economy* vliegen en de passagiers in *business class*. De gegevens moesten vervolgens geanalyseerd worden met de bedoeling maatregelen te nemen. Een voorbeeld hiervan kan zijn meer overboekingen toe te laten, dit wil zeggen meer reservaties opnemen dan er plaatsen zijn op het vliegtuig, en/of strenger optreden tegen passagiers die te veel bagage meenemen.

### Voorbeeld 1.2

Bij de productie van koffie is de luchtvochtigheid in de productiehal van cruciaal belang voor de kwaliteit van het eindproduct. De vochtigheidsgraad wordt onder controle gehouden door een systeem dat niet feilloos werkt. Daarom worden dagelijks meerdere metingen van het vochtgehalte verricht om na te gaan of het binnen de perken blijft. Deze aanpak valt onder de noemer statistische procescontrole.

### Voorbeeld 1.3

Een vulmachine voor flessen heeft doorgaans meerdere vulkoppen, zodat in één beweging meteen een aantal flessen gevuld kunnen worden. Bij een dergelijk vulproces worden elk uur typisch een aantal gevulde flessen zorgvuldig nagewogen om na te gaan of elke vulkop precies de gewenste hoeveelheid in de flessen deponeert. Een andere interessante onderzoeksvraag in deze context is of er verschillen zijn tussen de metingen uitgevoerd door verschillende analisten.

### Voorbeeld 1.4

Grootwarenhuizen verzamelen dankzij de klantenkaarten massa's gegevens. Zaken die typisch geregistreerd worden zijn

- het gespendeerde bedrag per winkelbeurt, al dan niet opgesplitst in categorieën (voeding, kleding, ...);
- het aantal verkochte artikelen;
- de betalingswijze (contant, debetkaart, kredietkaart en maaltijdcheque).

Onderzoekers maken gebruik van statistische methoden om deze gigantische hoeveelheid informatie samen te vatten en op een overzichtelijke manier voor te stellen.

### Voorbeeld 1.5

Financiële analisten zijn onder meer geïnteresseerd in de graad van risico van het beleggen in een bepaald aandeel. Daartoe houden zij gedurende jaren de maandelijkse rendementen van dat aandeel bij. Hierbij wordt niet enkel met koerswijzigingen, maar ook met uitgekeerde dividenden rekening gehouden. Bovendien worden de maandelijkse rendementen van de globale markt, bijvoorbeeld de BEL20-index, bijgehouden. Indien het aandeel gemiddeld procentueel sterker stijgt of daalt dan de markt, dan noemt men het aandeel risicovol. In het andere geval spreekt men van een aandeel met weinig risico. Via statistische methoden kunnen verbanden tussen de rendementen van het aandeel en de globale markt onderzocht worden.

## 1.4 Onderwerp van de statistiek

In de voorbeelden uit de vorige sectie worden telkens een of meerdere vragen onderzocht over een **populatie** van objecten of elementen of over een **proces** dat objecten of elementen genereert.

De gegevens over de populatie of het proces worden bekomen door een of meerdere eigenschappen of karakteristieken van hun elementen te registreren. Deze eigenschappen of karakteristieken worden **variabelen** genoemd. Deze naam geeft aan dat de waarde van de eigenschap varieert van element tot element. Daarom wordt de statistiek soms de studie van de variabiliteit genoemd.

Het is meestal onpraktisch om alle elementen uit een populatie of gegenereerd door een proces in een studie op te nemen. In die gevallen werkt men slechts met een deel van de elementen: de **steekproef**. Het is niet altijd eenvoudig om steekproefgegevens op een correcte manier te verzamelen. Aan het verzamelen van gegevens moet bij elk statistisch onderzoek dan ook de nodige aandacht besteed worden. In deze context wordt soms de afkorting GIGO<sup>1</sup> gebruikt. Dit staat voor *garbage in, garbage out* en slaat op het feit dat de meest geavanceerde statistische methoden weinig tot geen betrouwbare informatie kunnen halen uit gegevens van slechte kwaliteit.

#### Voorbeeld 1.6

Bij het peilen naar het kiesgedrag bij Belgische verkiezingen is de populatie erg duidelijk te omschrijven: alle kiesgerechtigde burgers van het land. Variabelen die in deze context geregistreerd kunnen worden zijn geslacht, beroep, politieke overtuiging en leeftijd.

#### Voorbeeld 1.7

Het opgooien van een dobbelsteen is een proces. Een mogelijke steekproef bestaat erin de dobbelsteen vijftig keer op te gooien. Variabelen die bij elke worp of waarneming geregistreerd kunnen worden zijn het aantal gegooide ogen of het al dan niet even zijn van het gegooide aantal ogen.

In de Voorbeelden 1.2 en 1.3 kunnen alle tijdstippen waarop het productieproces in werking is als de populatie beschouwd worden. Op een beperkt of eindig aantal tijdstippen kunnen metingen of waarnemingen omtrent de variabelen luchtvochtigheid (Voorbeeld 1.2) of gewicht (Voorbeeld 1.3) verricht worden. De metingen vormen samen de steekproef. Voor de financiële analist uit Voorbeeld 1.5 wordt de steekproef gevormd door een beperkte reeks rendementen en marktindices. In Voorbeeld 1.4 is de populatie waarnaar de interesse van de onderzoeker uitgaat de verzameling van alle klanten van het grootwarenhuis. Een mogelijke steekproef bestaat uit alle klanten die het warenhuis in de maand mei bezocht hebben en van hun klantenkaart gebruikgemaakt hebben.

De in een steekproef verzamelde gegevens kunnen op allerlei manieren overzichtelijk voorgesteld worden met behulp van tabellen en grafieken. Daarnaast laat ook het berekenen van een aantal kenmerkende waarden of statistieken, zoals bijvoorbeeld een gemiddelde, het toe om een overzichtelijk beeld van een verzameling gegevens samen te stellen. Het voorstellen van steekproefgegevens valt onder de noemer **beschrijvende of descriptieve statistiek**. Dit onderdeel

---

<sup>1</sup>Deze afkorting parodieert de afkortingen FIFO (first in first out) en LIFO (last in first out), die in accounting gebruikt worden bij het boeken van artikelen uit voorraad.

komt aan bod in de Hoofdstukken 2 en 3.

Het beschrijven van de steekproefgegevens is in veel gevallen slechts een eerste stap in een onderzoek. Een tweede onderdeel omvat het analyseren en interpreteren van de steekproefgegevens. Deze analyse en interpretatie zijn vereist om een antwoord te vinden op een aantal vooraf gestelde vragen over de populatie of het proces, om gestelde hypothesen te testen, of om de waarde of kwaliteit van een voorgesteld statistisch model te toetsen. De antwoorden en conclusies die hierbij bekomen worden, worden veralgemeend naar de populatie of het proces. Deze veralgemening wordt **inferentie** genoemd, wat meteen de benaming **inferentiële statistiek** verklaart. Andere vaak gebruikte benamingen zijn wiskundige statistiek, **verklarende statistiek** en steekproeftheorie.

De veralgemening van conclusies betreffende een reeks steekproefgegevens naar een gehele populatie of naar een proces is meteen de zwakke plek van de statistiek: op basis van steekproefgegevens kunnen nooit met zekerheid uitspraken gedaan worden over de populatie of het proces in kwestie. Aan de gedane uitspraken kan wel een betrouwbaarheid meegegeven worden indien bij het verzamelen van de steekproefgegevens statistisch verantwoorde methoden gebruikt werden. Deze graad van betrouwbaarheid wordt uitgedrukt met behulp van een kans, zodat een basiskennis van de kansrekening vereist is om statistische methoden te kunnen begrijpen en te kunnen toepassen.

### 1.5 Kansrekening

Nog meer dan het begrip statistiek klinken de woorden kans en waarschijnlijkheid vertrouwd in de oren. Zo heeft iedereen intuïtief een goed idee van de betekenis van een kans van  $1/4$  bij deelname aan een gokspel. Een dergelijke kans kan dan ook door vrijwel iedereen gebruikt worden om af te wegen of men al dan niet aan het spel zal deelnemen. De berekening van zo'n kans kan evenwel voor grotere moeilijkheden zorgen.

De kansrekening bestudeert processen of experimenten waarbij de uitkomst onzeker is. De begrippen proces en experiment dienen hier in hun breedste betekenis geïnterpreteerd te worden. Voorbeelden zijn het gooien van een dobbelsteen, de prijs van een aandeel bij het sluiten van de Nasdaq, de hypotheaire rentevoet, de vraag naar laptopcomputers van een merk, het percentage defecte producten in een productielijn gedurende een bepaalde periode, het aantal bezoekers op een website of het lukraak trekken van een winnaar uit alle deelnemers aan een tombola.

Het verschil tussen de kansrekening en de statistiek bestaat erin dat de kansrekening populaties en processen rechtstreeks bestudeert, terwijl de statistiek dit doet via steekproefgegevens. De kansrekening vertrekt hierbij telkens van een aantal veronderstellingen, aannames of assumpties omtrent de populatie of het proces. Een aantal voorbeelden verduidelijken dit.



**Voorbeeld 1.8**

Indien het bestudeerde proces het opgooien van een dobbelsteen is, dan kunnen we met behulp van de kansrekening de kans proberen te berekenen om minstens 20 keer een zes te gooien wanneer we de dobbelsteen 100 keer opgooien. Deze berekening is enkel mogelijk indien een belangrijke veronderstelling gemaakt wordt omtrent de gebruikte dobbelsteen, namelijk dat de dobbelsteen eerlijk is, of, met andere woorden, dat de dobbelsteen volledig homogeen en totaal symmetrisch is, zodat het even waarschijnlijk is om een één te gooien als een twee, een drie, een vier, een vijf of een zes.

Een mogelijk statistisch onderzoek over de dobbelsteen zou kunnen zijn: het nagaan van de eerlijkheid van de dobbelsteen. Hiertoe kan de dobbelsteen een (groot) aantal keer opgegooid worden om op die manier de nodige steekproefgegevens te verzamelen. Vervolgens kan een statistisch onderbouwde conclusie getrokken worden omtrent de hypothese dat de dobbelsteen eerlijk is.

**Voorbeeld 1.9**

Bij een industrieel vulproces kan, gegeven een aantal instellingen van de vulmachine en een aantal veronderstellingen inzake de nauwkeurigheid van de machine, de kans berekend worden dat een fles te weinig gevuld zal zijn. Een andere mogelijkheid is om de kans te berekenen dat er bij een levering van 1000 kratten ten hoogste 5% van de flessen zal zijn met een te kleine inhoud.

Een statistisch onderzoek van hetzelfde vulproces bestaat typisch uit het regelmatig wegen van een aantal flessen (steekproef) om na te gaan of de gemiddelde inhoud van de flessen niet te groot of te klein is, en of de variabiliteit of veranderlijkheid in de inhoud van fles tot fles niet te groot is.

**Voorbeeld 1.10**

De kansrekening maakt bij het bestuderen van het kiesgedrag van de Belgische bevolking de veronderstelling dat 30% voor partij A zal stemmen, 25% voor partij B, 20% voor partij C, en 25% blanco of voor een aantal kleinere partijen. De kansrekening kan in dat geval berekenen dat van elke 500 kiezers er gemiddeld 150 voor partij A zullen opteren, 125 voor partij B, 100 voor partij C en 125 blanco zullen stemmen of voor andere partijen kiezen.

De statistiek zal daarentegen een statistische voorspelling maken aan de hand van een steekproef van bijvoorbeeld 2000 kiezers. Bij deze voorspelling kan ook een foutenmarge gegeven worden.

Belangrijk is dus dat bij statistiek gewerkt wordt met een beperkte hoeveelheid steekproefinformatie en dat uitspraken over populaties en processen daarom fout kunnen zijn. Dit is de zwakte van statistiek. Idealiter zijn de kansen op fouten natuurlijk klein. De kans op fouten kan kleiner gemaakt worden door op oordeelkundige wijze veel kwaliteitsvolle gegevens te verzamelen.

Kansberekening heeft ook een zwakke plek: de veronderstellingen omtrent het bestudeerde proces of de bestudeerde populatie kunnen fout zijn, waardoor de conclusies ongeldig worden.

## 1.6 Software

Bij kansrekening en statistiek zijn vaak heel wat berekeningen nodig. Ook is het belangrijk om overzichtstabellen te maken van alle gegevens in een steekproef, of om grafische voorstellingen te maken van de gegevens. Het gebruik van een computer en van gespecialiseerde statistische software is dan noodzakelijk. Zoals aangegeven in het woord vooraf gebruiken we in dit boek het statistisch softwarepakket JMP.

[Dit boek is online te koop \(klik hier\)](#)

Dit boek levert eerst en vooral een toegankelijk en diepgaand overzicht van de belangrijkste statistische kengetallen voor nominale, ordinale en kwantitatieve gegevens. Omdat een figuur meer zegt dan 1000 woorden, besteedt het boek ook veel aandacht aan de grafische voorstelling van de gegevens. Niet enkel de oude vertrouwde histogrammen en taartdiagrammen komen aan bod, maar ook moderne grafieken, zoals de mozaïekplot en de bubble plot, worden uitvoerig besproken.

Het tweede deel van het boek handelt over kansrekening. Alle belangrijke discrete kansverdelingen (zoals de binomiale en Poisson verdelingen) en continue kansdichtheden (zoals de exponentiële, normale en lognormale kansdichtheden) worden in detail besproken, en hun gebruik wordt geïllustreerd aan de hand van leuke voorbeelden.

Doorheen het boek wordt gebruik gemaakt van het gebruiksvriendelijke, interactieve statistische pakket JMP voor het uitvoeren van berekeningen, het bepalen van kansen en het creëren van figuren. De benodigde tussenstappen worden in detail beschreven, zodat de lezer niet alleen de aangeleerde basisconcepten begrijpt, maar er ook nog mee aan de slag kan. Dit boek maakt van het overzichtelijk voorstellen van grote aantallen gegevens kinderspel.

PETER GOOS is gewoon hoogleraar aan de Faculteit Toegepaste Economische Wetenschappen van de Universiteit Antwerpen en de Faculteit Bio-Ingenieurswetenschappen van de KU Leuven, en hoogleraar aan de Erasmus School of Economics van de Erasmus Universiteit Rotterdam, waar hij kansrekening en statistiek doceert aan studenten toegepaste economische wetenschappen, sociaal-economische wetenschappen en handelsingenieur en bio-ingenieur, en aan studenten wiskunde en econometrie. Peter Goos is auteur van het boek *Kansen en Verwachtingen* (Acco, 2009) en ontving de Shewell Award en de Lloyd S. Nelson Award van de American Society for Quality en de Statistics in Chemistry Award van de American Statistical Association.



9 789033 493232