

# STATISTIEK I

## R-CODES

R-code	Betekenis
<i>Dim</i> ()	Aantal rijen (= aantal personen id steekproef) en kolommen (= aantal variabelen) vd tabel
<i>Head</i>	1e 6 rijen + bijhorende kolommen
<i>Names</i> ()	Namen variabelen
\$	Waarde ve variabele uit de date (altijd gevolgd door de naam vd data)
<i>Table</i> ()	Absolute/relatieve frequentie
<i>Read.table</i> ()	Data lezen
<i>Pie</i> ()	Cirkeldiagram
<i>Barplot</i> ()	Staafdiagram
<i>Cut</i> ()	Klassen
<i>Breaks</i>	Grenzen tss de klassen
<i>Hist</i> ()	Histogram
<i>Cumsum</i> ()	Cumulatieve absolute frequentie
<i>Ecdf</i> ()	Cumulatieve frequentiecurve
<i>Plot</i> ()	
<i>Mean</i> ()	Gemiddelde
<i>Median</i> ()	Mediaan
<i>Min</i> ()	Minimum
<i>Max</i> ()	Maximum
<i>Aad</i> ()	Gemiddelde absolute afwijking
<i>Var</i> ()	Variantie
<i>Sd</i> ()	Standaarddeviatie
	Bv. $Sd = \text{sqrt}(4)$ : standaarddeviatie = $\sqrt{4} = 2$
<i>Sqrt</i> ()	Vierkantswortel (vd variantie om de standaarddeviatie te vinden)
<i>Quantile</i> ()	Percentielen/kwartielen
<i>IQR</i> ()	Interkwartielafstand
<i>Boxplot</i> ()	Boxplot
<i>Cov</i> ()	Covariantie
<i>Cor</i> ()	Correlatiecoëfficiënt
<i>Cor</i> () + optie "Kendall"	Kendall's tau

$D_{\text{binom}}(k, N, p)$	Kansdichtheid $P(X = k)$
$P_{\text{binom}}(k, N, p)$	Cumulatieve verdelingsfunctie ( $P \leq k$ )
$P_{\text{norm}}(x, \mu, \sigma)$	Kans $P(X \leq x)$ voor variabele $X \sim N(\mu, \sigma^2)$
$P_{\text{norm}}(\text{bekomen } Z)$	!! in R dus niet $\sigma^2$ maar $\sigma$ (door vierkantswortel te nemen) !! bij standaardnormale mogen $\mu, \sigma$ worden weggelaten
$D_{\text{norm}}(x, \mu, \sigma)$	Kansdichtheid $f_x(x)$
$P_{\text{chisq}}(y, k)$	Kansen $P(Y \leq y)$ voor variabele $Y \sim \chi_k^2$
$Q_{\text{chisq}}(f_y(y), k)$	= y
$P_t(t, k)$	Kansen $P(T \leq t)$ voor variabele $T \sim t_k$
$t.\text{test}()$	t-toets (g-toets maar in R $\rightarrow$ t)
$mu$	$H_0$
$t$	Toetsingsgrootte g
$Df()$	= degrees of freedom: aantal vrijheidsgraden = k = n - 1
$p\text{-value}$	= p-waarde = onderscheidingskans
$t.\text{test}(\text{alternative} = \text{"less"})$	Linkszijdige toets
$t.\text{test}(\text{alternative} = \text{"greater"})$	Rechtszijdige toets

#### Belangrijke commando's

1. q

= quantile: waarde

- We weten de kans
- We willen de waarde weten waarvoor ...% ligt eronder/boven

2. p

= probability: kans

- We weten de waarde
- We willen de kans weten die onder/boven een bep waarde ligt

#### FORMULES EN BIJHORENDE SYMBOLEN

!! steeds finaleresultaat ( $\neq$  tussenbewerkingen) afronden op 2 decimalen ( $<5$  naar  $\searrow$ ,  $\leq 5$  naar  $\nearrow$ )

NAAM	FORMULE	HOE?
RELATIEVE FREQUENTIE	$\frac{\text{Absolute frequentie}}{\text{Aantal steekproefelementen}}$	Kennen
HARMONISCH GEMIDDELDE	$\frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$ $x_i = \text{waarde van } i^{\text{e}} \text{ element}$	Zullen we niet echt hanteren
MEETKUNDIG OF GEOMETRISCH GEMIDDELDE	$\sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n}$	Zullen we niet echt hanteren
(REKENKUNDIG) STECKPROEFGEMIDDELDE $\bar{x}$  = GEMIDDELDE O.B.V. WAARDEN VAN VARIABELE	$\frac{(x_1 + x_2 + \dots + x_n)}{n}$ $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$	Begrijpen, niet reproduceren
GEMIDDELDE O.B.V. EEN FREQUENTIEVERDELING	$\bar{x} = \frac{1}{n} \sum_{i=1}^p f_i x_i^u$  Waarbij  - $x_i^u =$ unieke waarde van variabele X in steekproef  Bv. $x_1^u =$ vrouw en $x_2^u =$ man  - $f_i =$ absolute frequentie van deze waarde  - $p =$ aantal unieke waarden van variabele X in steekproef  !! uitkomst hiervan = uitkomst gem. o.b.v. waarden (logisch)	Kennen
KLASSENMIDDEN	$\frac{a + b}{2}$  !! klassenmiddens $]a,b[ = ]a,b[ = ]a,b[ = ]a,b[$	Begrijpend lezen

<b>GEMIDDELDE VAN GEGROEPEERDE DATA</b>	$\bar{x} = \frac{1}{n} \sum_{i=1}^p f_i \left( \frac{a_i + b_i}{2} \right)$ <p>Waarbij <math>f_i \left( \frac{a_i + b_i}{2} \right) = \text{absolute frequentie} \cdot \text{klassenmidden}</math></p> <p>!! uitkomst hiervan = <math>\pm</math> uitkomst gem. o.b.v. waarden of frequentieverdeling</p>	Begrijpend lezen
<b>MEDIAAN <math>m d_x</math></b>	<p>Bij oneven waarden: middelste waarde in <b>geordende rij</b></p> <p>Bij even waarden: middelste 2 waarden gedeeld door 2</p>	<b>Kennen</b>
<b>MODUS <math>m o</math></b>	<p>Waarde/klasse met de hoogste frequentie</p> <p>!! er kunnen meerdere waarden/klassen zijn = <b>modi</b></p> <ul style="list-style-type: none"> <li>- 1 modus = <b>unimodaal</b></li> <li>- 2 modi = <b>bimodaal</b></li> </ul>	<b>Kennen</b>
<b>VARIATIEBREEDTE <math>v x</math></b>	<p>Grootste – kleinste waarde</p> <p>OF</p> <p>Bovengrens laatste klasse – ondergrens eerste klasse (= 2 uitersten)</p> <p>!! kan nooit negatief zijn: max (altijd) &gt; min</p>	<b>Kennen</b>
<b>GEMIDDELDE ABSOLUTE AFWIJING <math>g a_x</math></b>	$\frac{ x_1 - \bar{x}  +  x_2 - \bar{x}  + \dots +  x_n - \bar{x} }{n}$ $g a_x = \frac{1}{n} \sum_{i=1}^n  x_i - \bar{x} $ <p>!! absolute waarde: anders altijd 0 (+ en - heffen elkaar op bij gem.)</p>	<b>Kennen</b>
<b>STEEKPROEFVARIANTIE <math>s n_x^2</math> of <math>s_x^2</math></b>	$s n_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ $s_x^2 = \frac{1}{n - 1} \sum_{i=1}^n (x_i - \bar{x})^2$	Formularium

**STEEKPROEFVARIANTIE O.B.V.  
FREQUENTIEVERDELING**

$$sn_x^2 = \frac{1}{n} \sum_{i=1}^p f_i (x_i^u - \bar{x})^2$$

Formularium

$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n f_i (x_i - \bar{x})^2$$

Weten

Waarbij

- $x_i^u$  = unieke waarde vd variabele X in steekproef

Bv.  $x_1^u$  = vrouw en  $x_2^u$  = man

- $f_i$  = absolute frequentie van deze waarde

**STEEKPROEFSTANDAARDDEVIATIE  $sn_x$**

Vierkantswortel vd variantie

**Kennen**

$$sn_x = \sqrt{sn_x^2}$$

$$s_x = \sqrt{s_x^2}$$

**PERCENTIEL  $P_k$** 

$$\frac{F(P_k)}{n} = \frac{k}{100}$$

**Kennen**

Waarbij

- $P_k$  = het k-de percentiel

- $\frac{F(P_k)}{n}$  = cumulatieve relatieve frequentie

Bv. Voor het 10<sup>e</sup> percentiel:

- 1)  $k = 10$
- 2) Cumulatieve relatieve frequentie =  $k/100 = 10\%$
- 3) 10% vd waarden zijn hetzelfde of kleiner

Bijzondere percentielen:

Kwartiel Symbool      % vd waarden die gelijk zijn of eronder liggen

1 <sup>e</sup> kwartiel	$P_{25}$	25%
2 <sup>e</sup> kwartiel	$P_{50}$	50%
3 <sup>e</sup> kwartiel	$P_{75}$	75%
4 <sup>e</sup> kwartiel	$P_{100}$	100%

!!  $m d_x = P_{50}$ **INTERKWARTIELAFSTAND  $Q$** Verschil 3<sup>e</sup> en 1<sup>e</sup> kwartiel**Weten**

$$P_{75} - P_{25}$$

**INTERKWARTIELINTERVAL**

$$[P_{75}, P_{25}]$$

**Weten**

!! bevat 50% van alle waarden

**SPREIDINGSMAAT D**

$$d = \frac{1 - \frac{f_{mo}}{n}}{1 - \frac{1}{p}}$$

Formularium

Waarbij

- $p$  = aantal unieke waarden
- $f_{mo}$  = aantal keer dat de  $m$
- $n$  = aantal steekproefelementen

!! van 0 (geen spreiding) t.e.m. 1 (maximale spreiding)

**OUTLIERS**

$$P_{25} - 1,5 \cdot Q$$

**Weten**

$$P_{75} + 1,5 \cdot Q$$

**COVARIANTIE  $cov_{xy}$** 

$$cov_{XY} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Formularium

<p><b>CORRELATIECOËFFICIËNT <math>r_{xy}</math></b></p>	$r_{XY} = \frac{cov_{XY}}{s_X s_Y}$ <p>1. Covariantie berekenen</p> $cov_{XY} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$ <p>2. Standaarddeviatie berekenen voor x</p> $\sqrt{s_x^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$ <p>3. Standaarddeviatie berekenen voor y</p> $\sqrt{s_y^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$ <p>4. <math>\frac{\text{covariantie}}{\text{standaarddeviatie } x \cdot \text{standaarddevia}}</math></p> <p>!! correlatiecoëfficiënt heeft altijd zelfde teken als covariantie</p>	<p><b>Weten</b></p>
<p><b>CONCORDANT PAAR</b></p>	$\frac{y_j - y_i}{x_j - x_i} > 0$ <p>= positieve hellingsgraad</p>	<p><b>Kennen</b></p>
<p><b>DISCORDANT PAAR</b></p>	$\frac{y_j - y_i}{x_j - x_i} < 0$ <p>= negatieve hellingsgraad</p>	<p><b>Kennen</b></p>
<p><b>KENDALL'S TAU <math>\tau</math></b></p>	$\tau = \frac{2(\#concordante\ paren - \#discordante\ paren)}{n(n-1)}$	<p>Formularium</p>
<p><b>REGRESSIELIJN</b></p>	$Y = b_0 + b_1 X$ <p>Waarbij</p> <ul style="list-style-type: none"> <li>- <math>b_1</math> = regressiecoëfficiënt: helling vd rechte</li> <li>- <math>b_0</math> = intercept: snijpunt met y-as</li> </ul>	<p>Formularium</p>

<b>REGRESSIECOËFFICIËNT BIJ PERFECT LINEAIR VERBAND</b>	$b_1 = \frac{y_j - y_i}{x_j - x_i}$	<b>Kennen</b>
<b>INTERCEPT BIJ PERFECT LINEAIR VERBAND</b>	$b_0 = y_i - b_1 x_i$	<b>Kennen</b>
<b>REGRESSIECOËFFICIËNT BIJ NIET PERFECT LINEAIR VERBAND</b>	$b_1 = r_{XY} \frac{s_Y}{s_X}$  !! $b_1$ zal altijd zelfde teken hebben als $r_{xy}$	Formularium
<b>INTERCEPT BIJ NIET PERFECT LINEAIR VERBAND</b>	$b_0 = \bar{y} - b_1 \bar{x}$	Formularium
<b>VANAF HIER: OP POPULATIE NIVEAU</b>		
<b>LIMIET VAN DE RELatieve FREQUENTIE (WANNEER STEEKPROEF = <math>\infty</math>)</b>	$P(X = x_i) = \lim_{n \rightarrow \infty} \frac{f_i}{n}$	Begrijpen, niet reproduceren
<b>CUMULATIEVE VERDELINGSFUNCTIE <math>F_X(x)</math> BIJ DISCRETE EN CONTINUE VARIABELEN</b>	Univariate: $F_X(x) = P(X \leq x)$ Bivariate: $F_{X,Y}(x, y) = P(X \leq x \text{ en } Y \leq y)$	Begrijpen, niet reproduceren
<b>KANS BIJ CONTINUE VARIABELEN</b>	$P(X = x) = 0$  !! daarom beroep doen op dichtheidsfunctie	Begrijpen, niet reproduceren
<b>DICHTHEIDSFUNCTIE OF KANSDICHTHEID</b>	$f_x(x) = \lim_{b \rightarrow 0} \frac{F_x(x+b) - F_x(x)}{b}$	Begrijpen, niet reproduceren
<b>INTEGRATIE DICHTHEIDSFUNCTIE</b>	$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} f_x(x) dx$  $P(X \leq x) = \int_{-\infty}^x f_x(x) dx$  $P(X > x) = \int_x^{+\infty} f_x(x) dx$	Begrijpen, niet reproduceren
<b>KANS BIJ CONTINUE VARIABELEN O.B.V. EIGENSCHAP</b>	$P(x_1 \leq X \leq x_2) = P(X \leq x_2) - P(X \leq x_1) = F_x(x_2) - F_x(x_1)$	<b>Kennen</b>
<b>VOLLEDIGE OPPERVLAKTE ONDER DICHTHEIDSFUNCTIE</b>	$\int_{-\infty}^{+\infty} f_x(x) dx = 1$	Begrijpen, niet reproduceren

**POPULATIE GEMIDDELDE OF  
VERWACHTINGSWAARDE  $E(X)$ ,  $\mu_x$  of  $\mu$**

Discreet:  $E(X) = \sum_{i=1}^p P(X = x_i)x_i$

Continu:  $E(X) = \int_{-\infty}^{+\infty} fx(x)dx$

Formularium

Begrijpen,  
niet  
reproduceren

**POPULATIE VARIANTIE  $V(X)$ ,  $\sigma_x^2$  of  $\sigma^2$**

Discreet:  $V(X) = \sum_{i=1}^p P(X = x_i)(x_i - E(x))^2$

Continu :  $V(X) = \int_{-\infty}^{+\infty} fx(x)(x - E(X))^2 dx$

Formularium

Begrijpen,  
niet  
reproduceren

**POPULATIESTANDAARDDEVIATIE  $\sigma_x$  of  $\sigma$**

$$\sigma_x = \sqrt{V(X)}$$

**Weten**

**UNIVARIATE KANSVERDELING BIJ  
DISCRETE VARIABELEN**

$$P(X = x_i) = \sum_{j=1}^q P(X = x_i \text{ en } Y = y_j)$$

$$P(Y = y_i) = \sum_{i=1}^p P(X = x_i \text{ en } Y = y_i)$$

Waarbij

- p = aantal mogelijke waarden dat X kan aannemen
- q = aantal mogelijke waarden dat Y kan aannemen

**STATISTISCHE ONAFHANKELIJKHEID**

Discreet:

**Kennen**

2 discrete variabelen X en Y zijn onafhankelijk indien:

$$P(X = x_i \text{ en } Y = y_j) = P(X = x_i)P(Y = y_j)$$

geldt voor alle mogelijke combinaties i en j

1. Alle kansen neerschrijven
2. Marginale verdeling berekenen
3. Marginale verdeling vermenigvuldigen
4. Deze aan elkaar stellen

Niet kunnen

Continu:

2 continue variabelen X en Y zijn onafhankelijk indien:

$$P(X \leq x \text{ en } Y \leq y) = P(X \leq x)P(Y \leq y)$$

Voor alle mogelijke waarden x en y

**POPULATIE COVARIANTIE  $COV(X, Y)$**

Discreet:

Formularium

$$COV(X, Y) = \sum_{i=1}^p \sum_{j=1}^q P(X = x_i \text{ en } Y = y_j)(y_j - E(Y))$$

Continu:

Begrijpen,  
niet  
reduceren

$$COV(X, Y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{x,y}(x,y)(x - E(X))(y - E(Y)) dx dy$$

**POPULATIE CORRELATIECOËFFICIËNT  $\rho_{XY}$**

Discreet:  $\rho_{XY} = \frac{COV(X, Y)}{\sigma_X \sigma_Y}$

**Kennen**

Continu:  $\rho_{XY} = \frac{COV(X, Y)}{\sigma_X \sigma_Y}$

Waarbij

- $\sigma_X$  = standaarddeviatie van x
- $\sigma_Y$  = standaarddeviatie van y

**EIGENSCHAP POPULATIE GEMIDDELDE OF VERWACHTINGSWAARDE**

$E(a) = a$

**Weten**

Indien a = constante

Bv. Iedereen 2000 euro inkomen  $\rightarrow$  gem. = 2000

**EIGENSCHAP POPULATIE VARIANTIE**

$V(a) = 0$

**Weten**

Bv. Iedereen zelfde inkomen: geen spreiding en variantie = 0

**BIJZONDERE VERDELINGEN****BINOMIALE KANSVERDELING**

$X \sim \text{Binom}(N, p)$

$$P(X = k) = \frac{N!}{k!(N-k)!} p^k (1-p)^{N-k}$$

Formularium

**Weten**

Waarbij

- N = max. aantal successen
- k = aantal gewenste successen
- p = de kans op een succes
- ! = faculteit (via GRM uitrekenen)

Bv.  $4! = 4 \cdot 3 \cdot 2 \cdot 1$

Enkel wanneer:

- N = vast
- p blijft ongewijzigd

Op grafiek

- Bij kleine kans op succes: scheef naar rechts  
  
Logisch want meeste hebben lage score en enkel de uitzonderingen een hoge
- Bij helft kans op succes: symmetrisch
- Bij grote kans op succes: scheef naar links  
  
Logisch want meeste hebben hoge score en enkel de uitzonderingen een lage

**POPULATIE GEMIDDELDE OF VERWACHTINGSWAARDE VAN EEN BINOMIALE VARIABELE**

$E(X) = N \cdot p$

Formularium

**POPULATIE VARIANTIE VAN EEN BINOMIALE VARIABELE**

$V(X) = N \cdot p (1 - p)$

Formularium

**DICHTHEIDSFUNCTIE NORMALE VERDELING**

$$X \sim N(\mu, \sigma^2)$$

$$f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Formularium

**Weten**

Waarbij

- $\pi = 3,14\dots$
- $e = 2,71$
- $\mu =$  populatie gemiddelde =  $E(X)$
- $\sigma^2 =$  populatie variantie =  $V(X)$
- $\sigma =$  populatie standaarddeviatie

**Weten**

Op grafiek:

- Hoogste punt (top) = gemiddelde
- Grote variantie = laag + breed
- Kleine variantie = hoog + smal
- Symmetrisch
- Enkel positieve waarden

**INTEGRAAL VAN DE KANS VAN EEN NORMALE VERDELING**

$$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

Begrijpen,  
niet  
reduceren

**PRINCIPES STANDAARDNORMALE  
VERDELING  $X \sim N(0,1)$**

$$\mu = 0$$

**Weten**

$$\sigma^2 = 1$$

2 belangrijke eigenschappen

1. Symmetrisch rond 0

Bijgevolg:

$$P(X > x) = P(X < -x)$$

2. Totale opp. = 1

Bijgevolg:

$$P(X \leq -x) = 1 - P(X \leq x)$$

$$P(X \geq -x) = 1 - P(X \geq x)$$

$$P(X < x) = 1 - P(X > x)$$

$$P(X > x) = 1 - P(X < x)$$

Sidenote:

- Bij - teken veranderen haakjes van kant
- $< = \leq$  en  $> = \geq$  (want continue variabelen)

**STANDAARDISEREN VAN X BIJ EEN  
NORMALE VERDELING**

Omdat  $\mu \neq$  altijd 0 en  $\sigma^2 \neq$  altijd 1  $\rightarrow$   
standaardiseren:

1. 
$$P(X \leq x) = P\left(\frac{X - \mu}{\sigma} \leq \frac{x - \mu}{\sigma}\right)$$

Formularium

**Weten**

2.  $Z = \frac{X - \mu}{\sigma} \rightarrow$  dit vervangen in  
bovenstaande formule

3. 
$$P(X \leq x) = P\left(Z \leq \frac{x - \mu}{\sigma}\right)$$

Hierdoor: nieuwe variabele  $Z$  die de  
standaardnormale verdeling wél volgt  $\rightarrow$   
 $Z \sim N(0,1)$

<b>VARIABELE VAN DE CHI-KWADRAAT VERDELING</b> $Y \sim \chi_k^2$	$Y = x_1^2 + x_2^2 + \dots + x_k^2$ <p>Formularium</p> <p><b>Weten</b></p> <p>Waarbij</p> <ul style="list-style-type: none"> <li>- k = aantal vrijheidsgraden (en ook populatie gemiddelde)</li> <li>- <math>\chi_k \sim N(0,1)</math></li> </ul> <p>Op grafiek:</p> <ul style="list-style-type: none"> <li>- Hoogste punt (top) = k en dus ook populatie gemiddelde</li> <li>- Scheve verdeling: <u>a</u>symmetrisch</li> <li>- Enkel positieve waarden</li> <li>- Totale opp./kans = 1</li> </ul>	
<b>POPULATIE GEMIDDELDE BIJ EEN CHI-KWARDRAAT VERDELING</b>	$E(Y) = k$	Formularium
<b>POPULATIE VARIANTIE BIJ EEN CHI-KWARDRAAT VERDELING</b>	$V(Y) = 2k$	Formularium
<b>VARIABELE VAN DE STUDENT-T VERDELING</b> $T \sim t_k$	$T = \frac{X}{\sqrt{\frac{1}{k}Y}}$ <p>k = aantal vrijheidsgraden</p> <p>Op grafiek:</p> <ul style="list-style-type: none"> <li>- Indien <math>k \rightarrow \infty (= t_\infty)</math>: valt exact samen met standaardnormale</li> <li>- Hoogste punt (top) = populatie gemiddelde = 0</li> <li>- Symmetrisch</li> </ul>	Formularium  Interpretatie ervan niet kennen  <b>Weten</b>
<b>POPULATIEGEMIDDELDE BIJ EEN T-VERDELING</b>	$E(T) = 0$	Formularium
<b>VARIANTIE BIJ EEN T-VERDELING</b>	$V(T) = \frac{k}{k-2}, \text{ voor } k > 2$	Formularium

**VANAF HIER: STEEKPROEVENVERDELING**

**Steekproefgemiddelde**

<p><b>STEEKPROEFGEMIDDELDE <math>\bar{X}</math></b></p> <p><b>= GEMIDDELDE VAN VERSCHILLENDE STEEKPROEVEN</b></p>	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ <p>!!</p> <p><math>\bar{X}</math> = steekproefgemiddelde voor een bepaalde steekproef, voor een steekproef in algemeen</p> <p><math>\bar{x}</math> = steekproefgemiddelde o.b.v. 1 specifieke steekproef</p>	<p><b>Kennen</b></p>
<p><b>VERWACHTINGSWAARDE VAN HET STEEKPROEFGEMIDDELDE <math>E(\bar{X})</math></b></p>	$E(\bar{X}) = \mu_x$	<p>Formularium</p>
<p><b>VARIANTIE VAN HET STEEKPROEFGEMIDDELDE <math>V(\bar{X})</math></b></p>	$V(\bar{X}) = \frac{\sigma_x^2}{n}$	<p>Formularium</p>
<p><b>VERDELING VAN HET STEEKPROEFGEMIDDELDE</b></p> <p><b>= STEEKPROEFVERDELING VAN HET GEMIDDELDE</b></p>	<p>Bij onafhankelijke, lukrake trekkingen uit populatie dat normaal verdeeld is:</p> $\bar{X} \sim N\left(\mu_x, \frac{\sigma_x^2}{n}\right)$	<p>Formularium</p>
<p><b>STANDAARDISEREN VAN HET STEEKPROEFGEMIDDELDE</b></p>	$Z \leq \frac{x - \mu}{\sqrt{\sigma_x^2 / n}}$ <p>of <math>Z \leq \frac{x - \mu}{\sigma_x / \sqrt{n}}</math></p> <p>Voorwaarde:</p> <ul style="list-style-type: none"> <li>- X komt uit normale verdeling</li> </ul> <p>Hierbij: n maakt niet uit</p> <p>OF</p> <ul style="list-style-type: none"> <li>- <math>n \geq 30</math></li> </ul>	<p><b>Kennen</b></p>
<p><b>Steekproefvariantie</b></p>		
<p><b>STEEKPROEFVARIANTIE <math>SN_x^2</math> of <math>S_x^2</math></b></p>	$SN_x^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ $S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$	<p><b>Kennen</b></p>

VERWACHTINGSWAARDE VOOR  
STEEKPROEFVARIANTIE

$$E(SN_x^2) = \frac{n-1}{n} \sigma_x^2$$

Formularium

$$E(S_x^2) = \sigma_x^2$$

!! daarom: in praktijk meer  $S_x^2$  dan  $SN_x^2$

VERDELING VAN DE  
STEEKPROEFVARIANTIE

Bij onafhankelijke, lukrake trekkingen uit  
populatie dat normaal verdeeld is:

Formularium

$$\frac{(n-1)S_x^2}{\sigma_x^2} \sim \chi_{n-1}^2$$

### SCHATTERS

EEN GOEDE SCHATTER  $\hat{\theta}$  VOOR EEN  
POPULATIEPARAMETER  $\theta$

$\hat{\theta}$  is een goede schatter voor  $\theta$  indien:

Weten

1. De schatter zuiver is: verwachtingswaarde  
schatter = populatieparameter

$$E(\hat{\theta}) = \theta$$

2. De variantie vd schatter  $V(\hat{\theta})$  kleiner  
wordt naarmate de steekproefgrootte  $n$   
 $\nearrow$

= naarmate  $n \nearrow$  wordt de schatter  
nauwkeuriger

STANDAARDDEVIATIE VAN DE SCHATTER =  
STANDAARDFOUT

$$\sqrt{V(\hat{\theta})}$$

Weten

$$\text{Of } \frac{\sigma}{\sqrt{n}} \text{ want } V(\bar{X}) = \frac{\sigma_x^2}{n}$$

Schatter met kleinste standaardfout = het  
efficiëntst

**STEEKPROEFGEMIDDELDE: EEN GOEDE  
SCHATTER VOOR  
POPULATIEGEMIDDELDE?**

Steekproefgemiddelde = goede schatter want

**Weten**

1. 
$$E(\bar{X}) = \mu_x$$

Dus gem. van alle steekproefgem.  $\approx$   
populatie gem.

2. Variantie steekproefgemiddelde =

$$V(\bar{X}) = \frac{\sigma_x^2}{n}$$

Dus naarmate  $n \nearrow \rightarrow$  nauwkeuriger

$\Rightarrow$  Ja!! goede schatter

**STEEKPROEFVARIANTIE: EEN GOEDE  
SCHATTER VOOR POPULATIEVARIANTIE?**

2 formules voor steekproefvariantie  $\rightarrow$  bekijken  
beide

**Weten**

Formule 1: 
$$sn_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Deze is geen goede schatter, want

1. 
$$E(SN_x^2) = \frac{n-1}{n} \sigma_x^2$$

Dus  $\sigma_x^2$  (populatievariantie) zal steeds  
te klein worden geschat (door  $(n-1)/n$ )

2. Naarmate  $n \nearrow$  zal het wel  
nauwkeuriger worden

$\Rightarrow$  slechts 1 vd 2 voorwaarden is voldaan

$\Rightarrow$  geen goede schatter

Formule 2: 
$$s_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

1. 
$$E(S_x^2) = \sigma_x^2$$

Verwachting vd steekproefvariantie =  
populatievariantie  $\rightarrow$  goed!!

2. Naarmate  $n \nearrow$  zal het wel  
nauwkeuriger worden

$\Rightarrow$  beide voorwaarden zijn voldaan

$\Rightarrow$  deze formule = voorkeur

Via betrouwbaarheidsinterval (enkel tweezijdig kunnen)

<b>BETROUWBAARHEIDSINTERVAL WANNEER VARIANTIE GEKEND</b>	$\left[ \bar{X} - z_{\frac{\alpha}{2}} \sigma / \sqrt{n}, \bar{X} + z_{\frac{\alpha}{2}} \sigma / \sqrt{n} \right]$ <p><math>Z_{\alpha}</math> = waarde vd standaardnormale verdeling zodat de opp. vd curve <u>rechts</u> vd waarde = <math>\alpha</math></p>	<b>Kennen</b>
<b>BETROUWBAARHEIDSINTERVAL WANNEER VARIANTIE ONGEKEND IS</b>	$\left[ \bar{X} - t_{n-1; \alpha/2} S_x / \sqrt{n}, \bar{X} + t_{n-1; \alpha/2} S_x / \sqrt{n} \right]$ <p><math>T_{n-1; \alpha/2}</math> = waarde vd t-verdeling zodat de opp. vd curve <u>rechts</u> vd waarde = <math>\alpha</math></p>	<b>Formularium</b>

Via statistisch toetsen (eenzijdig en tweezijdig kunnen)

(Hierbij gaan we er altijd vanuit dat x normaal verdeeld is of  $n \geq 30$  is)

<b>NULHYPOTHESE</b>	Bij tweezijdig, linkszijdig en rechtszijdig  $H_0: \mu = \mu_0$ <p><math>\mu_0</math> = gegeven waarde</p>	<b>Formularium</b>
<b>ALTERNATIEVE HYPOTHESE</b>	Tweezijdig:  $H_a: \mu \neq \mu_0$ Linkszijdig:  $H_a: \mu < \mu_0$ Rechtszijdig:  $H_a: \mu > \mu_0$	<b>Formularium</b>

$$G = \frac{\bar{X} - \mu_0}{S_X / \sqrt{n}}$$

Wanneer nulhypothese waar is ( $\mu = \mu_0$ ):

- G volgt  $t_{n-1}$ -verdeling
- Waarden van G liggen rond 0
- Waarden van G zijn + en -

**Weten**

Wanneer nulhypothese niet waar is en  $\mu > \mu_0$ :

- Grotere waarden  $\leftrightarrow \mu = \mu_0$
- G heeft enkel positieve waarden

Wanneer nulhypothese niet waar is en  $\mu < \mu_0$ :

- Kleinere waarden  $\leftrightarrow \mu = \mu_0$
- G heeft enkel negatieve waarden

g = waarde van G die we bekomen o.b.v. één steekproef

Maar: wanneer wat doen met g-waarde?

- Als g rond 0 ligt  $\rightarrow H_0$  niet verwerpen
- Als g sterk van 0 verschilt  $\rightarrow H_0$  wel verwerpen

Wat is "rond" en wat is "sterk verschilt"?  $\rightarrow$  **verschillende beslissingsregels**

# 1.

## BESLISSINGSREGELS O.B.V. AANVAARDINGSGEBIED/KRITIEKE WAARDEN

Tweezijdig:

Indien  $g$  tussen het volgende interval ligt:

$$-t_{n-1;\alpha/2} \leq g \leq t_{n-1;\alpha/2}$$

$H_0$  niet verwerpen

Ligt het er niet in  $\rightarrow H_0$  wél verwerpen

M.a.w.:

$$|g| \leq t_{n-1;\alpha/2} \rightarrow H_0 \text{ niet verwerpen} \quad \text{Formularium}$$

$$|g| > t_{n-1;\alpha/2} \rightarrow H_0 \text{ verwerpen, } H_a \text{ aanvaarden}$$

Linkszijdig:

$$g > -t_{n-1;\alpha} \rightarrow H_0 \text{ niet verwerpen}$$

$$g < -t_{n-1;\alpha} \rightarrow H_0 \text{ verwerpen, } H_a \text{ aanvaarden}$$

Rechtszijdig:

$$g \leq t_{n-1;\alpha} \rightarrow H_0 \text{ niet verwerpen}$$

$$g > t_{n-1;\alpha} \rightarrow H_0 \text{ verwerpen, } H_a \text{ aanvaarden}$$

$\alpha$  = significantie niveau

Komen overeen met interval vd tweezijdige toets

**Weten**

Formularium

**Weten**

Formularium

**TYPE I FOUT**

$H_0$  = correct maar we verwerpen ze  $\rightarrow$  **type I fout**

**Weten**

**KANS OP EEN TYPE I FOUT**

$$P(\text{verwerp } H_0 \mid \mu = \mu_0) = \alpha$$

**Weten**

$\alpha$  = het significantieniveau en is altijd gegeven

**KANS OP CORRECT BESLUIT INDIEN**  
 $\mu = \mu_0$

$P(\text{verwerp } H_0 \text{ niet} \mid \mu = \mu_0) = 1 - \alpha$

**Weten**

OF

Kijken naar ...% betrouwbaarheidsinterval

Bv. Bij 95% betrouwbaarheidsinterval:

1.  $\alpha = 0,05$

$1 - 0,05 = 0,95$

Dit: kans op correct besluit

OF

2. 95% betrouwbaarheidsinterval dus 95% kans op correct besluit

**TYPE II FOUT**

$H_0 \neq$  correct maar we verwerpen ze niet  $\rightarrow$  **type II fout** **Weten**

**KANS OP EEN TYPE II FOUT**

$P(\text{verwerp } H_0 \text{ niet} \mid \mu \neq \mu_0) = \beta$

**Weten**

$\beta = \text{bèta}$

**KANS OP CORRECT BESLUIT INDIEN**  
 $\mu \neq \mu_0$

$P(\text{verwerp } H_0 \mid \mu \neq \mu_0) = 1 - \beta$

**Weten**

= **ONDSCHIEDINGSKANS OF POWER**

Invloeden

- Significantieniveau

Wanneer  $\alpha \nearrow \rightarrow \beta \searrow$

- Steekproefgrootte

Wanneer  $n \nearrow \rightarrow \beta \searrow$

En dus kans op correct besluit  $(1 - \beta) \nearrow$

**2.**

Indien  $\mu_0$  tussen het volgende interval ligt:

Formularium

**BESLISSINGSREGELS O.B.V.**  
**BETROUWBAARHEIDSINTERVAL**

$$\left[ \bar{X} - t_{n-1; \alpha/2} S_x / \sqrt{n}, \bar{X} + t_{n-1; \alpha/2} S_x / \sqrt{n} \right]$$

$H_0$  niet verwerpen

Ligt het er niet in  $\rightarrow H_0$  wél verwerpen

### 3.

#### BESLISSINGSREGELS O.B.V. P-WAARDE OF OVERSCHRIJDINGSKANS

p-waarde

Formularium

- Wordt berekend in veronderstelling dat  $H_0$  waar is
- Hangt af vd  $H_a$

Basisregel:

$p \geq \alpha \rightarrow H_0$  niet verwerpen

$p < \alpha \rightarrow H_0$  verwerpen

Linkszijdig:  $P(G < g | \mu = \mu_0)$

1.  $P(T < g)$  berekenen
2. Deze  $g$ : in R output "pt()" steken
3. Bekomen waarde vergelijken met  $\alpha$  + basisregel toepassen

Rechtszijdig:  $P(G > g | \mu = \mu_0)$

1.  $P(T > g)$  berekenen
2.  $1 - \text{pt}(g, (n-1))$  doen
3. Bekomen waarde vergelijken met  $\alpha$  + basisregel toepassen

Tweezijdig:

1.  $g$  waarde berekenen
2.  $g$  waarde vergelijken met 0
  - Als  $g > 0$ :  
 $p = 2 \cdot P(T > g)$
  - Als  $g \leq 0$ :  
 $p = 2 \cdot P(T < g)$

Deze kansen  $\rightarrow$  aflezen in R

## VERDUIDELIJK POPULATIEPARAMETERS

	Bij discrete variabelen	Bij continue variabelen
<b>Populatie gemiddelde / verwachtingswaarde</b> $E(X)$ , $\mu_x$ of $\mu$	$E(X) = \sum_{i=1}^p P(X = x_i)x_i$	$P(X = x_i) = 0 \rightarrow$ andere definitie nodig: $E(X) = \int_{-\infty}^{+\infty} f_x(x)dx$
<b>Populatie variantie</b> $V(X)$ , $\sigma_x^2$ of $\sigma^2$	$V(X) = \sum_{i=1}^p P(X = x_i)(x_i - E(x))^2$	$V(X) = \int_{-\infty}^{+\infty} f_x(x)(x - E(X))^2 dx$
<b>Populatie covariantie</b>	$\text{COV}(X, Y) = \sum_{i=1}^p \sum_{j=1}^q P(X = x_i \text{ en } Y = y_j)(y_j - E(Y))$	$\text{COV}(X, Y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{x,y}(x,y)(x - E(X))(y - E(Y))dx dy$
<b>Populatie correlatiecoëfficiënt</b>		$\rho_{XY} = \frac{\text{COV}(X, Y)}{\sigma_X \sigma_Y}$
	Formularium (behalve correlatie)	Moeten we niet kunnen uitrekenen

## STAPPENPLAN BETROUWBAARHEIDSINTERVAL

Wanneer? Populatie gem. is ongekend maar we willen er toch uitspraak over doen

Vraag 1: X normaal verdeeld?		
JA		NEE
Vraag 2: Populatievariantie gekend?		
JA	NEE	NEE

$\left[ \bar{X} - z_{\frac{\alpha}{2}} \sigma / \sqrt{n}, \bar{X} + z_{\frac{\alpha}{2}} \sigma / \sqrt{n} \right]$ <p>Ligt het tss dit interval <math>\rightarrow</math> interval bevat pop.gem.</p> <p>Bv. Bij 95% betrouwbaarheidsinterval <math>\rightarrow</math> in 95% vd gevallen ligt pop. gem. erin</p> <p><u>Invloeden</u></p> <ol style="list-style-type: none"> <li>1. Steekproefgrootte Naarmate <math>n \nearrow \rightarrow</math> interval smaller</li> <li>2. Als <math>\alpha \nearrow \rightarrow</math> interval smaller Want <math>1 - \alpha</math> zal <math>\searrow</math> en hiertss = interval</li> </ol> <p>Indien je de kans wil <math>\nearrow</math> dat pop.gem. erin ligt <math>\rightarrow</math> interval moet breed zijn <u>maar</u> indien breed: niet zo informatief meer</p> <p><math>\Rightarrow</math> compromis tss beide vinden (vaak 95%)</p>	$\left[ \bar{X} - t_{n-1; \alpha/2} S_x / \sqrt{n}, \bar{X} + t_{n-1; \alpha/2} S_x / \sqrt{n} \right]$ <p>Nu: t-verdeling ipv normale verdeling</p> <p>Deze: lijken op elkaar</p> <p>Toch verschillen:</p> <ul style="list-style-type: none"> <li>- <math>T_{n-1}</math>-verdeling heeft een grotere variantie</li> <li>- <math>T_{n-1; \alpha/2}</math>-waarde is groter dan <math>z_{\frac{\alpha}{2}}</math>-waarde</li> </ul> <p>Maar ook hierbij: naarmate <math>n \nearrow \rightarrow</math> steeds betere benadering standaardnormale verdeling</p>	<p>Indien grote steekproef: centrale limietstelling</p> <p>Indien kleine steekproef: GEEN LEERSTOF VAN STATISTIEK I</p>
--	--	---

### STAPPENPLAN STATISTISCH TOETSEN

1.  $H_0$  en  $H_a$  opstellen (uit gekregen opgave)
2. Significantieniveau vaststellen
3. Gem. en standaarddeviatie vd specifieke steekproef berekenen (of uit opgave halen)
4. Toetsingsgrootte  $g$  berekenen

$$G = \frac{\bar{X} - \mu_0}{S_X / \sqrt{n}}$$

5. Beslissingsregels toepassen
  - M.b.v. kritieke waarde  
!! bij kritische waarde:  $1 - P(T < t\text{-waarde})$  doen  
Dit =  $P(T > t\text{-waarde})$  wat altijd zo is (want t-waarde = waarde rechts vd grafiek)
  - M.b.v. betrouwbaarheidsinterval
  - M.b.v. p-waarde

!! andere beslissingsregels voor eenzijdig, linkszijdig en rechtszijdig
6. Conclusie formuleren:  $H_0$  verwerpen of niet?

- Niet verwerpen =  $H_0$  aanvaarden
- Verwerpen =  $H_a$  aanvaarden

### **Opletten !!**

Voor zowel de  $Z_\alpha$  - waarde, de  $T_{n-1; \alpha/2}$ -waarde als de kritische waarde: gaat het om de waarde rechts vd grafiek

Maar: R-output geeft standaard wat links onder grafiek ligt

### Oplissing?

1. Eigenschap normale/t-verdeling gebruiken

1 - ... doen

OF

2. Bij R-output extra info geven: "lower.tail"

### Voorbeeld

$$\alpha = 0,025 \Rightarrow P(Z > Z_{0,025}) = 0,025$$

Dus  $1 - 0,025 = 0,975$

Dit in R-output: `qnorm(0,975)`

Uitkomst van deze = juiste z-waarde