

# HOOFDSTUK 1

## TYPISCHE FOUTEN BIJ STATISTIEK

- ° Foute gegevens
- ° Fouten in berekening kans
- ° Foute interpretatie resultaten

**Statistiek : de wetenschap van het leren uit data & van het meten, controleren en communiceren van onzekerheid**

## 1. Eigenschappen van variabelen

### 1.1 Verschillende schaalfamilies

<b>Nominaal</b>	- Namen → Geen hoeveelheid, gewoon identificatie - Waarden kunnen ook getallen zijn	<i>Variabele = geslacht, waarde = man of vrouw</i> <i>Variabele = nummer tram, waarde = 24, 21..</i> <i>Variabele = land, waarde = België, Frankrijk..</i>
<b>Ordinaal</b>	- Geen hoeveelheid - Hiërarchie : ordening! → volgorde belangrijk, waarde zelf niet	<i>Variabele = uitslag wedstrijd, waarde = goud, zilver, brons</i> <i>Variabele = mate van instemming, waarde = volledig oneens, neutraal, volledig eens...</i>
<b>Interval</b>	- Hiërarchie - Waarde zelf ook belangrijk - Geen absoluut nulpunt - Geen onderlinge verhoudingen - Recht evenredig : onderlinge verschillen blijven even groot	<i>Variabele = temperatuur, waarde = 10°C, 5°C</i> <i>→ 10°C is niet het dubbele van 5°C (want in F is het niet zo op een grafiek)</i> <i>Recht evenredig: even groot verschil tss 10° en 20° &amp; 50° en 60°</i> <i>0° is geen absoluut nulpunt: 32°F</i>
<b>Ratio</b>	- Absoluut nulpunt - Verhoudingen	<i>€10 is dubbel zoveel dan €5</i> <i>Variabele = lengte in cm, geldbedrag in euro..</i>

Dia 49 : vragen 1 = ratio , 2 = interval , 3 = nominaal , 4 = ratio , 5 = ordinaal & ratio , 6 = ordinaal

### 1.2 Discrete & continue variabelen

<b>Continue variabelen</b>	- Tussenwaarden → tussen elke 2 waarden ligt een 3 <sup>e</sup> → oneindig veel waarden	<i>Lengte in cm</i> <i>Temperatuur in °C</i> <i>Tijd in seconden</i>
<b>Discrete variabelen</b>	- Geen tussenwaarden → eindig aantal waarden	<i>Aantal kinderen, aantal volgers op Twitter, aantal GSMs die je al hebt</i>

# Deel I : beschrijvende statistiek

## HOOFDSTUK 2 : VISUALISEREN VAN DATA

<b>Populatie</b>	Volledige verzameling van objecten of personen waarover men info wil
<b>Steekproef</b>	Deelverzameling van de populatie, die ook echt onderzocht wordt → moet representatief zijn : aselect!

<b>CIRKELDIAGRAM</b>	- Nominaal - Relatieve oppervlaktes cirkel = relatieve frequenties
----------------------	---

<b>STAAFDIAGRAM</b>	<ul style="list-style-type: none"> <li>- Nadeel : geen goed overzicht onderlinge verhoudingen</li> <li>- Nominaal &amp; ordinaal</li> <li>- Horiz : waarden variabele</li> <li style="padding-left: 20px;">Vertic : AF of RF</li> <li>- Rechthoeken los van elkaar, breedte en afstand even groot</li> <li>- <b>Voordeel : snel overzicht onderlinge verhoudingen</b></li> </ul>
<b>HISTOGRAM</b>	<ul style="list-style-type: none"> <li>- Interval &amp; ratio</li> <li>- Voordeel : snel overzicht onderlinge verhoudingen</li> <li>- Horiz : waarden variabele</li> <li>- Rechthoeken tegen elkaar <ul style="list-style-type: none"> <li>→ Breedte = klassenbreedte</li> <li>→ Oppervlakte = RF</li> <li>→ Hoogte = RF / klassenbreedte</li> </ul> </li> <li>- Klassenindeling <ul style="list-style-type: none"> <li>→ Beslist onderzoeker zelf</li> <li>→ Klassenbreedte niet altijd even groot, dus rechthoeken niet altijd even breed</li> <li>→ Als klassen toch even groot zijn, kan.. <ul style="list-style-type: none"> <li>.. hoogte rechthoek = AF</li> <li>.. hoogte rechthoek = RF</li> <li>.. oppervlakte rechthoek = RF</li> </ul> </li> <li>→ Uiterste klassen heel kleine AF = samenvoegen</li> <li>→ vuistregel aantal klassen : <math>\sqrt{n}</math></li> </ul> </li> <li>- Verdeling <ul style="list-style-type: none"> <li>→ Symmetrische verdeling</li> <li>→ Scheef naar rechts : staart naar rechts</li> <li>→ Scheef naar links : staart naar links</li> </ul> </li> </ul>

## 1. Algemene begrippen & notatie

<b>ABSOLUTE FREQUENTIES (AF)</b>	Het aantal (per waarde)
<b>ABSOLUTE FREQUENTIEVERDELING</b>	Tabel met absolute frequenties
<b>VARIABELE –NOTATIE-</b>	Met een hoofdletter (vaak X) → Waarden die variabele aanneemt : kleine letters, cijfers als subscriptie
<b>STEEKPROEFGROOTTE (N)</b>	
<b>RELATIEVE FREQUENTIES (RF)</b>	
<b>VERDELING VE VARIABELE</b>	Het geheel van mogelijke waarden, samen met de absolute en/of relatieve frequenties
<b>KLASSEN –NOTATIE- GEGROEPEERDE FREQUENTIEVERDELING</b>	$]a, b]$ tabel, 2 kolommen : klassen & overeenkomstige frequenties

## 2. Cumulatieve frequentiecurve

### 2.1 Ongegroepeerde data

<b># INFO</b>	Meer dan bij gegroepeerde data (indelen in klassen leidt tot informatieverlies)
<b>CUMULATIEVE ABSOLUTE FREQUENTIE</b>	absolute frequenties optellen <ul style="list-style-type: none"> <li>- geeft dan het aantal gegevens weer die gelijk aan of kleiner dan de bijhorende waarde zijn</li> <li>- grootste waarde = steekproefgrootte (want alles is kleiner of gelijk aan)</li> </ul>

**CUMULATIEVE RELATIEVE  
FREQUENTIE**

zelfde systeem als bij rel.abs.fr.

**CUMULATIEVE  
FREQUENTIECURVE**

Horizontaal : gegevens

Verticaal : cumulatieve frequenties

Tekenen

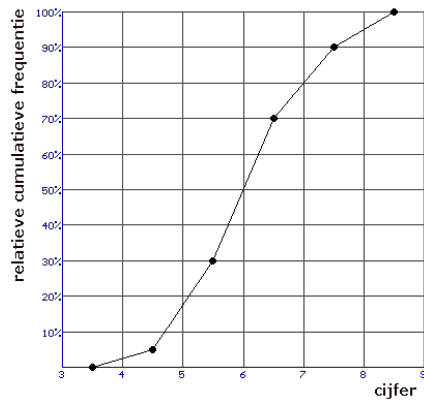
- Alle waarden aanduiden

- Onderling trapsgewijs verbinden

- Laagste waarde & hoogste waarde : horizontale lijn!

## 2.2 Gegroepede data

<b># INFO</b>	Minder dan bij ongegroepeerde data → Je hebt info over een klasse, niet over een specifieke waarde! → NADEEL
<b>CUMULATIEVE FREQUENTIE</b>	Frequenties optellen
<b>CUMULATIEVE FREQUENTIECURVE</b>	Telkens het klassenmidden gebruiken als punt voor op de grafiek



## 2.3 Illustratie methoden

Zie cursus pagina 59 – 63

# HOOFDSTUK 3 : SAMENVATTEN VAN DATA

(Centrummaten & spreidingsmaten, tot aan variantie en standaarddeviatie: zie geschreven samenvatting)

## 2.4 De interkwartielafstand

### 2.4.1 Percentielen $p_k$

Voor een geheel getal  $k$  tussen 0 en 100, is het  $k$ -de percentiel (symbool  $p_k$ ) het getal  $p_k$  waarvoor geldt dat:

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

WAT BETEKENT DIT?

- ° Het  $k$ -de percentiel is de waarde waarvan  $k\%$  van de hele verzameling kleiner of gelijk aan die waarde is
- ° Komt zowat overeen met de cumulatieve relatieve frequentie

VOORBEELD?

*P10: tiende percentiel: de waarden van een variabele, waarvoor 10% van de waarden hetzelfde of kleiner zijn*

SPECIAAL PERCENTIEL

mediaan = percentiel P50

### 2.4.2 Kwartielen

<b>1<sup>E</sup> KWARTIEL</b>	P25	→ 25% van alle waarden
<b>2<sup>E</sup> KWARTIEL</b>	P50	→ Mediaan : 50% van alle waarden
<b>3<sup>E</sup> KWARTIEL</b>	P75	→ 75% van alle waarden
<b>INTERKWARTIELAFSTAND</b>	Q	° 3 <sup>e</sup> kwartiel – 1 <sup>e</sup> kwartiel ( $P_{75} - P_{25}$ ) ° Interval & ratio
<b>INTERKWARTIELINTERVAL</b>	$[P_{25}, P_{75}]$	° Overspant 50% van alle waarden ° Ordinaal, interval & ratio

## 2.5 De spreidingsmaat $d$

<b>VARIABELEN</b>	Allemaal Vooral nominaal	
<b>UITKOMST</b>	Uitkomst ligt tussen 0 en 1 , geeft op die manier spreiding weer ° 0 = geen spreiding ° 1 = maximale spreiding	
<b>FORMULE</b>	$d = \frac{1 - \frac{f_{mo}}{n}}{1 - \frac{1}{p}}$	KOMT IN FORMULARIUM
	$p$	aantal unieke waarden die een variabele kan aannemen
	$f_{mo}$	frequentie van de modus (kan een waarde of een klasse zijn)
	$n$	steekproefgrootte
	WAT ALS...	
	.. $f_{mo} = n$	Geen spreiding, want alle waarden zijn gelijk aan de modus
	.. $f_{mo} = n/p$	$d = 1$

## 2.6 Gevoeligheid aan outliers

<b>HOE BEREKEN JE OF EEN SPREIDINGSMAAT GEVOELIG IS OF NIET?</b>	1. Bereken spreidingsmaat met alle waarden (inclusief outliers) 2. Bereken spreidingsmaat zonder outliers 3. Als er een groot verschil is tussen deze waarden : de spreidingsmaat is gevoelig aan outliers!	
<b>GEVOELIGHEID?</b>	Variatiebreedte	Ja
	Gemiddelde absolute afwijking	Ja
	Variantie	Ja
	Standaarddeviatie	Ja
	Interkwartielafstand	Neen
	d	Neen

### 3. Boxplot

<b>DATA GROEPEREN?</b>	Hoeft niet - Niet gebruikersafhankelijk - Verschil met histogram!		
<b>VOORDEEL</b>	- Handig om overzicht te krijgen over verdeling van data (mediaan, interkwartielafstand..) - Makkelijk weten wat de outliers zijn		
<b>OUTLIERS VASTSTELLEN – REKENREGEL</b>	Laagste outliers	$P_{25} - 1.5 * Q$	Alles lager dan deze waarde = outlier
	Hoogste outliers	$P_{75} + 1.5 * Q$	Alles hoger dan deze waarde= outlier
<b>HOE TEKENEN</b>	1. Alle waarden op grafiek tekenen 2. Outliers bepalen (via rekenregel) & aanduiden op grafiek 3. Horizontale lijn bij laagste & hoogste waarde die geen outlier is 4. Horizontale lijn voor $P_{25}$ en $P_{75}$ 5. Lijnen voor kwartielen met elkaar verbinden, als een rechthoek 6. Alle waarden van grafiek wissen, behalve outliers 7. Verticale stippellijn tekenen tussen overblijvende horizontale lijn & grens rechthoek → stippellijnen noemen we <b>snorharen</b> of <b>whiskers</b> 8. Mediaan aanduiden met horizontale lijn <b>ILLUSTRATIE: ZIE CURSUS P. 100!</b>		

## HOOFDSTUK 4: SAMENHANG TUSSEN 2 VARIABELEN

Hoofdstuk 2 & 3 : één variabele per keer bekijken → univariate statistiek

Hoofdstuk 4 : twee variabelen tezamen bekijken → bivariate statistiek

# 1. Bivariate frequentieverdeling

WAT?	Frequentietabel, maar dan voor 2 variabelen (in plaats van voor 1)
------	--

VOORDEEL	<ul style="list-style-type: none"> <li>° We kunnen vanuit de bivariate tabel de univariate gegevens afleiden</li> <li>- OPGELET! Het werkt niet langs de andere kant! Vanuit univariate gegevens kunnen we geen bivariate gegevens afleiden</li> <li>- <u>Marginale verdeling</u> = andere naam voor univariate verdeling die je kent via de bivariate verdeling</li> </ul>
NADEEL	<ul style="list-style-type: none"> <li>° Conclusies kunnen anders zijn naarmate de data anders gegroepeerd is → subjectief</li> <li>- Oplossing : spreidingsdiagram &amp; correlatiecoëfficiënten</li> </ul>

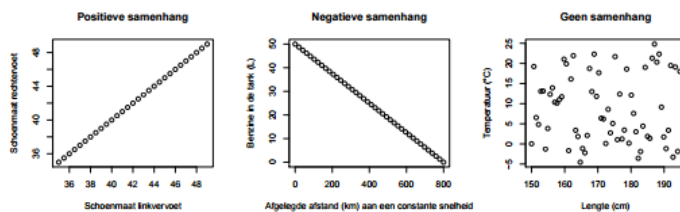
## 2. Spreidingsdiagram

WAT?            ° Geeft samenhang tussen 2 variabelen weer  
 ° Alle waarden als bollen op grafiek

SOORTEN?     *(Veel te extreem weergegeven, in realiteit bijna nooit zo)*

- ° Positieve samenhang
- ° Negatieve samenhang
- ° Geen samenhang

	Klassen Hersengrootte		
IQ Groep	[790, 886]	[886, 982]	[982, 1080]
gemiddeld IQ	9	9	2
hoog IQ	7	9	4



## 3. Maten van samenhang

### 3.1 De covariantie $cov_{XY}$

NOTATIE	$cov_{XY}$
FORMULE	$cov_{XY} = \frac{1}{n-1} * \sum_{i=1}^n (x_i - \bar{x}) * (y_i - \bar{y})$
MEETNIVEAU SAMENHANG?	Beide variabelen : interval, ratio Uitkomst positief                      positieve samenhang Uitkomst negatief                      negatieve samenhang Uitkomst ong 0                          geen samenhang
NADEEL	<ul style="list-style-type: none"> <li>° grootte van de covariantie hangt af van sterkte van samenhang &amp; meeteenheid</li> <li>- Je kan niet echt zeker weten of je samenhang nu écht groot is of niet</li> <li>- oplossing : correlatiecoëfficiënt</li> </ul>

### 3.2 De correlatiecoëfficiënt $r_{XY}$

NOTATIE	$r_{XY}$
FORMULE	$r_{XY} = \frac{COV_{XY}}{S_X * S_Y}$
EIGENSCHAPPEN	<ul style="list-style-type: none"> <li>° Ligt altijd tussen -1 en 1</li> <li>* 1 = perfect positieve samenhang</li> <li>* -1 = perfect negatieve samenhang</li> <li>° Hetzelfde teken als covariantie</li> <li>° Enkel te gebruiken bij lineaire samenhang</li> </ul>
VOORDEEL	Beter dan covariantie, want de uitkomst is onafhankelijk van de meeteenheid → Het is altijd tussen -1 en 1
PROBLEMEN/VRAGEN BIJ OEFENINGEN	Als er geen lineaire samenhang is, kan je $r_{XY}$ nog altijd berekenen zonder probleem, alleen is dat getal dan niet betrouwbaar

### 3.3 Kendall's Tau $\tau$

FORMULE	$\tau = \frac{2 * (\text{aantal concordante paren} - \text{aantal discordante paren})}{n * (n - 1)}$ <p><i>(Komt in formularium)</i></p>
CONCORDANTE PAREN	Mathematisch $\frac{y_j - y_i}{x_j - x_i} > 0$ Grafisch        Positieve rico → stijgende lijn
DISCORDANTE PAREN	Mathematisch $\frac{y_j - y_i}{x_j - x_i} < 0$ Grafisch        Negatieve rico → dalende lijn
FORMULE # PAREN TOTAAL?	Om te weten hoeveel mogelijke paren er zijn (zowel concordant als discordant) $\frac{n * (n - 1)}{2}$
WERKWIJZE	<ul style="list-style-type: none"> <li>° Spreidingsdiagram tekenen</li> <li>° Alle mogelijke rechten tussen 2 punten trekken</li> <li>° Aantal concordante en discordante paren tellen</li> <li>° Formule toepassen (formularium!)</li> </ul>
EIGENSCHAPPEN	° Altijd tussen -1 en 1
MEETNIVEAU	Ordinaal, interval, ratio

### 3.4 Lineaire en niet-lineaire verbanden

LINEAIRE FUNCTIE	<ul style="list-style-type: none"> <li>° Grafisch kan je een rechte lijn tekenen</li> <li>° <u>CORRELATIECOEFFICIENT</u> gebruiken</li> </ul>
MONOTONE FUNCTIE	<ul style="list-style-type: none"> <li>° Bewaart de orde : eenmaal stijgen/dalen, blijven stijgen/dalen → Maar dus niet noodzakelijk in een rechte lijn</li> <li>° <u>KENDALL'S <math>\tau</math></u> gebruiken</li> </ul>
NIET-MONOTONE FUNCTIE	<ul style="list-style-type: none"> <li>° Bewaart de orde niet</li> <li>° correlatiecoëfficiënt en Kendall's <math>\tau</math> allebei niet goed</li> </ul>
<u>BELANGRIJKE TIP</u>	Data eerst visualiseren, dan pas weet je welke spreidingsmaat goed is



**PROBLEMEN/VRAGEN  
BIJ OEFENINGEN**

Zien aan de grafiek of het een lineaire of monotone samenhang is: echt heel algemeen kijken, niet te gedetailleerd!

### 3.5 Gevoeligheid aan outliers

GEVOELIG AAN OUTLIERS	Covariantie Correlatiecoëfficiënt
NIET GEVOELIG AAN OUTLIERS	Kendall's $\tau$

## 4. De regressielijn

WAT?	Regressielijn zorgt ervoor dat we de $r_{XY}$ kunnen visualiseren op een spreidingsdiagram → onze vroegere 'functie' dus
MEETNIVEAU	Interval, ratio
FORMULE (VOOR DE RECHTE)	$Y = b_0 + b_1 * X$ <p><math>b_1</math> regressiecoëfficiënt – helling van de rechte (- de rico)</p> $b_1 = \frac{y_j - y_i}{x_j - x_i}$ <p><math>b_0</math> Snijpunt met de verticale as - intercept</p> $b_0 = y_i - b_1 * x_i$ <p>!!! Deze 2 formules te gebruiken bij perfecte samenhang, anders de 2 hieronder!!!</p>
KLEINSTE KWADRATENMETHODE	<ul style="list-style-type: none"> <li>° Erg vaak is de samenhang niet perfect : onmogelijk om een regressielijn te tekenen die door alle punten gaat</li> <li>° We willen uiteindelijk een rechte die toch zo goed mogelijk door de punten gaat</li> <li>° Methode om dit te bereiken : kleinste kwadratenmethode</li> </ul> <p>→ <math>\sum_{i=1}^n (y_i - (b_0 + b_1 * x_i))^2</math></p> <p>→ Logica erachter :</p> <ul style="list-style-type: none"> <li>* Je hebt je regressielijn &amp; je eigenlijke punten</li> <li>* De afstand tussen elk punt en de lijn is je fout (want je zit ernaast)</li> <li>* Het kwadraat van deze afstand wil je zo klein mogelijk, op die manier is je fout zo klein mogelijk</li> </ul> <p>→ Kwadraat van het bolletje (<math>y_i</math>) en de regressielijn (<math>(b_0 + b_1 * x_i)</math>)</p> <p>→ Via deze 2 formules zijn je <math>b_0</math> en <math>b_1</math> het meest geschikt volgens de methode</p> <ul style="list-style-type: none"> <li>* deze waarden kan je dan integreren in de 1<sup>e</sup> formule</li> <li>* <math>b_1 = r_{XY} * \frac{s_Y}{s_X}</math></li> <li>* <math>b_0 = \bar{y} - b_1 * \bar{x}</math></li> </ul>

## 5. Samenhang en causaliteit

Samenhang betekent GEEN causaliteit!

Er kan een derde variabele zijn, die niet bestudeerd is

# Deel II : Kansrekening

WAAR GAAT ELK DEEL OVER?

Deel I steekproef

Deel II populatie

Deel III inductieproces

We willen altijd iets weten over de populatie, maar dat is te veel om allemaal te onderzoeken. Daarom doen we aan steekproeftrekking. Maar uiteindelijk willen we deze resultaten veralgemenen naar de gehele populatie.

## HOOFDSTUK 5: DE POPULATIE EN VERDELINGSFUNCTIES

VERDELINGSFUNCTIE	<p><i>Frequentieverdeling, maar dan voor een populatie (geen steekproef)</i></p> <p>Frequentieverdeling – steekproef Verdelingsfunctie – populatie</p> <p>Hoe deze eruit ziet hangt af van het soort variabele: discreet of continu</p>
-------------------	---

### 1. Verdelingsfunctie discrete variabelen

DISCRETE VARIABELEN	<ul style="list-style-type: none"> <li>◦ kan geen tussenwaarden aannemen</li> <li>◦ Eindig aantal waarden             <ul style="list-style-type: none"> <li>→ Aantal mogelijke waarden: p</li> <li>→ OPGELET : het is de variabele die een eindig aantal waarden heeft, niet de populatie (want erg vaak heeft populatie zodanig veel waarden, dat het wiskundig makkelijker is om aan te nemen dat er oneindig veel zijn)</li> </ul> </li> </ul>
NOTATIE REL.FR. VAN DE POPULATIE	$P(X = x_i) = \lim_{n \rightarrow \infty} \frac{f_i}{n}$ <ul style="list-style-type: none"> <li>◦ De kans dat de variabele X de waarde <math>x_i</math> aanneemt</li> <li>◦ <math>\frac{f_i}{n}</math> = de relatieve frequentie</li> </ul>
SOORTEN?	<ul style="list-style-type: none"> <li>◦ kansverdeling</li> <li>◦ Cumulatieve verdelingsfunctie</li> </ul>

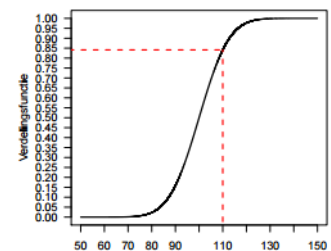
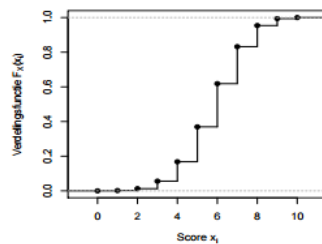
## 1.1 Kansverdeling

WAT ?	<ul style="list-style-type: none"> <li>° <i>relatieve frequentieverdeling van de populatie (geen steekproef)</i></li> <li>° tabel met 2 kolommen : de waarden van <math>x_i</math> &amp; de overeenkomstige kansen</li> </ul>
-------	---

## 1.2 Cumulatieve verdelingsfunctie / verdelingsfunctie $F_X(x)$

WAT?	<ul style="list-style-type: none"> <li>° <i>Cumulatieve relatieve frequentie van de populatie (geen steekproef)</i></li> <li>° Geeft de kans weer dat een waarde van X kleiner of gelijk aan x is</li> </ul>
FORMULE?	$F_X(x) = P(X \leq x)$
GRAFIEK OF TABEL?	<p>Kan beide zijn! Het gaat gewoon om het feit dat er aan elke X-waarde een bijhorende Y-waarde wordt gekoppeld, maakt niet uit op welke manier dat weergegeven wordt</p> <p>→ OPGELET! Beide zijn alleen mogelijk bij discrete variabelen! Bij continue variabelen is een tabel niet mogelijk, want er zijn oneindig veel punten</p>
GRAFIEK	Trapsgewijs

Score $x_i$	$F_X(x_i) = P(X \leq x_i)$
0	0.00013
1	0.00171
2	0.01265
3	0.05593
4	0.16798
5	0.36938
6	0.61860
7	0.83201
8	0.95385
9	0.99388
10	1.00000



## 2. Verdelingsfunctie continue variabelen

Score $x_i$	$P(X = x_i)$
0	0.00013
1	0.00158
2	0.01094
3	0.04328
4	0.11205
5	0.20140
6	0.24922
7	0.21341
8	0.12184
9	0.04003
10	0.00612

CONTINUE VARIABELEN	<ul style="list-style-type: none"> <li>° Kan oneindig veel tussenwaarden aannemen</li> <li>° Kan oneindig veel waarden aannemen</li> </ul>
---------------------	--

Probleem : doordat er oneindig veel waarden zijn, is de kans dat 1 specifieke waarde voorkomt quasi 0

→  $P(X = x_i) = 0$

→ We gaan een andere manier moeten vinden om kansen te berekenen (zie volgende titels)

### 2.1 Cumulatieve verdelingsfunctie $F_X(x)$

WAT?	De kans dat een waarde van X kleiner of gelijk aan x is
FORMULE	$F_X(x) = P(X \leq x)$ of $F_X(x) = P(X < x)$

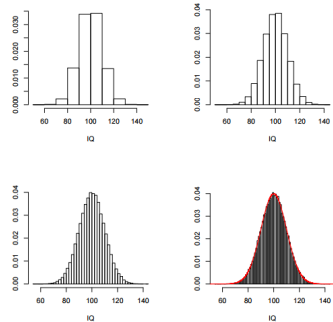
GRAFIEK

≤ of < maakt niet uit, want de kans is uiteindelijk toch 0  
 continu (niet trapsgewijs)

2.2 De dichtheidsfunctie of de kansdichtheid  $f_X(x)$

WAT?

- Formele uitleg: de afgeleide van de verdelingsfunctie
- Duidelijkere uitleg:
  - Histogram tekenen, waarbij oppervlakte rechthoek gelijk is aan relatieve frequentie
  - Er zijn een oneindig aantal waarden, dus we kunnen histogram opstellen met oneindig aantal klassen
  - Hoe meer klassen, hoe meer het lijkt op de dichtheidsfunctie



FORMULE

AFGELEIDEN NIET ZELF BEREKENEN, GEWOON OM WAT TE KUNNEN VATTEN WAT DICHTHEIDSFUNCTIE IS

$$f_X(x) = \lim_{b \rightarrow 0} \frac{F_X(x+b) - F_X(x)}{b}$$

- De kans dat X valt binnen het interval  $]x, x + b]$  gedeeld door b
- b = de breedte van het interval ; gaat richting 0

KANSEN BEREKENEN?

- Welke soort kansen?
  - Van de vorm  $P(x_1 \leq X \leq x_2)$
- Via integralen (niet zelf uitvoeren, gewoon begrijpen)
  - Grafisch : de oppervlakte tussen de 2 grenswaarden is wat je zoekt
  - Algemene formules (niet zelf kennen of gebruiken, gwn begrijpen)
 
$$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} f_X(x) dx$$

$$P(X \leq x) = \int_{-\infty}^x f_X(x) dx$$

$$P(X > x) = \int_x^{+\infty} f_X(x) dx$$

EIGENSCHAPPEN

- $P(x_1 \leq X \leq x_2) = F_X(x_2) - F_X(x_1)$
- Altijd positief, nooit negatieve waarden
  - functie is gebaseerd op een kans, een kans kan nooit negatief zijn
- Oppervlakte onder dichtheidsfunctie = 1
  - Komt altijd overeen met de kans ; de volledige kans is 100% of 1
- $P(X > x) = 1 - P(X < x)$

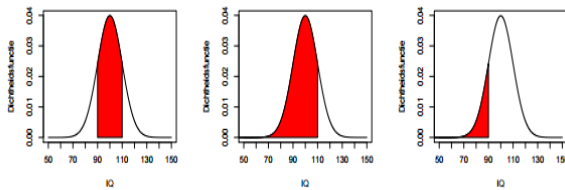
Voorbeeld: in de linker figuur zie je de kans dat X tussen 110 en 90 ligt. Dit is gelijk aan de kans dat X onder 110 ligt (middelste figuur), min de kans de X onder 90 ligt (rechter figuur)

### 3. Populatieparameters

#### 3.1 Populatiegemiddelde of verwachtingswaarde $E(X)$ , $\mu_x$ of $\mu$

Weetje : de E komt van 'expectation', vandaar verwachtingswaarde

##### 3.1.1. Discrete variabelen



FORMULE	$E(X) = \sum_{i=1}^p P(X = x_i) * x_i$		
WAT IS VERANDERD?	relatieve frequenties	→	kansen $P(X = x_i)$
	$x_i^u$	→	$x_i$

##### 3.1.2. Continue variabelen

We kunnen de vorige definitie niet gebruiken, want  $P(X = x_i) = 0$   
DUS gaan we gebruik maken van integralen (weeral niet zelf uitrekenen)

#### 3.2 Populatievariantie $V(X)$ , $\sigma_X^2$ of $\sigma^2$

##### 3.2.1. Discrete variabelen

FORMULE	$V(X) = \sum_{i=1}^p P(X = x_i) * (x_i - E(X))^2$		
WAT IS VERANDERD?	relatieve steekproeffrequenties	→	kansen $P(X = x_i)$
	Steekproefgemiddelde	→	$E(X)$
	$x_i^u$	→	$x_i$
STANDAARDDEVIATIE $\sigma$	$\sigma = \sqrt{V(X)}$		

##### 3.2.2. Continue variabelen

gebruik maken van integralen: niet zelf kunnen, wel begrijpen

## 4. Bivariate kansverdelingen

- Twee variabelen samen bekijken (in tabel), daar dan kansen over uitspreken
- Vergelijkbaar met hoofdstuk 4 (alleen was H4 over steekproefverdeling)

### 4.1 Discrete variabelen

<p>MARGINALE VERDELINGEN AFLEIDEN</p>	<p>° Zelfde als normaal : de aparte waarden optellen</p> <p>° <math>P(X = x_i) = \sum_{j=1}^q P(X = x_i \text{ en } Y = y_j)</math></p> <p>→ X wordt vastgehouden bij de waarde <math>x_i</math> (verandert dus niet), dan gaan we alle mogelijke y-waarden optellen die passen bij die waarde <math>x_i</math></p> <p>→ q = aantal mogelijke waarden van y</p> <p>° <math>P(Y = y_j) = \sum_{i=1}^p P(X = x_i \text{ en } Y = y_j)</math></p> <p>→ Y wordt vastgehouden bij de waarde <math>y_j</math> (verandert dus niet), dan gaan we alle mogelijke x-waarden optellen die passen bij die waarde <math>y_j</math></p>
<p>STATISTISCHE ONAFHANKELIJKHEID</p>	<p>° 2 discrete variabelen X en Y zijn onafhankelijk als de gelijkheid <math>P(X = x_i \text{ en } Y = y_j) = P(X = x_i) * P(Y = y_j)</math> geldt voor alle mogelijke combinaties van i en j</p> <p>→ De kans dat zowel X als Y een specifieke waarde aannemen moet altijd gelijk zijn aan de kans dat X haar specifieke waarde aanneemt, vermenigvuldigd met de kans dat Y zijn specifieke waarde aanneemt</p> <p>→ Als dit inderdaad zo is, zijn de variabelen onafhankelijk. Anders zijn ze afhankelijk.</p>
<p>COVARIANTIE COV(X,Y)</p>	$COV(X, Y) = \sum_{i=1}^p \sum_{j=1}^q P(X = x_i \text{ en } Y = y_j) * (x_i - E(X)) * (y_j - E(Y))$
<p>CORRELATIECOËFFICIËNT <math>\rho_{XY}</math></p>	$\rho_{XY} = \frac{COV(X, Y)}{\sigma_X * \sigma_Y}$

Illustratie statistische onafhankelijkheid: uit de tabel lees je af dat

$$P(X=4 \text{ en } Y=10) = 0.09119$$

$$\text{maar... } P(X=4) \times P(Y=10) = 0.35461 \times 0.35286 = 0.1251277$$

Deze twee uitkomsten zijn niet hetzelfde, dus zijn X en Y niet onafhankelijk.

## 4.2 Continue variabelen

- Erg vaak kennis nodig van integralen : deze gebruiken we niet zelf, dus continue variabelen wordt maar zeer beknopt besproken
- Lijkt behoorlijk hard op univariate continue variabelen
  - o  $P(X = x_i \text{ en } Y = y_j) = 0$

CUMULATIEVE BIVARIATE VERDELINGSFUNCTIE $F_{X,Y}(x, y)$	$F_{X,Y}(x, y) = P(X \leq x \text{ en } Y \leq y)$
BIVARIATE DICHTHEIDSFUNCTIE $f_{X,Y}(x, y)$	° De afgeleide van de cumulatieve bivariate verdelingsfunctie
STATISTISCHE ONAFHANKELIJKHEID	Continue variabelen zijn onafhankelijk als voor alle mogelijke waarden van x en y geldt dat: $P(X \leq x \text{ en } Y \leq y) = P(X \leq x) * P(Y \leq y)$
COVARIANTIE COV(X,Y)	° Maakt gebruik van integralen : zullen we nooit zelf moeten berekenen ° $COV(X, Y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{X,Y}(x, y) * (x - E(X)) * (y - E(Y)) dx dy$
CORRELATIECOËFFICIËNT $\rho_{XY}$	$\rho_{XY} = \frac{COV(X, Y)}{\sigma_X * \sigma_Y}$

Score X	Leeftijd Y	
	10	11
0	0.00341	0.00021
1	0.02730	0.00404
2	0.08275	0.03291
3	0.12110	0.13337
4	0.09119	0.26342
5	0.02711	0.21319

## 5. Nuttige stellingen

1. **Als X en Y onafhankelijk zijn dan COV(X,Y) = 0**
  - a. OPGELET: omgekeerd niet altijd waar! Niet altijd dat als de covariantie nul is, dat ze onafhankelijk zijn
    - i. Kan ook bijvoorbeeld door niet-lineaire samenhang komen
    - b. Alleen bij populatie, niet bij steekproef
2. **Als Y = X+a dan E(Y) = E(X)+a**
  - a. Hierbij is a een constante
    - i. Het gemiddelde van a = a
  - b. Bij populatie en steekproef
3. **Als Y = a\*X dan E(Y) = a\*E(X)**
  - a. Zelfde logica als stelling 2
  - b. Bij populatie en steekproef
4. **E(X+Y) = E(X) + E(Y)**  
**E(X-Y) = E(X) - E(Y)**
  - a. Zowel bij onafhankelijke als afhankelijke variabelen
  - b. Bij populatie en steekproef

5. Als X en Y onafhankelijk zijn  
dan  $E(X*Y) = E(X)*E(Y)$ 
  - a. Alleen populatie, niet steekproef
6. Als  $Y = X+a$   
dan  $V(Y) = V(X)$ 
  - a. Hierbij is a een constante
7. Als  $Y = a*X$   
dan  $V(Y) = a^2*V(X)$ 
  - a. Hierbij is a een constante
8.  $V(X+Y) = V(X) + V(Y) + 2*COV(X,Y)$ 
  - a. Als X en Y onafhankelijk zijn  
dan  $V(X+Y) = V(X) + V(Y)$ 
    - i. Komt voort uit stelling 1 en 8
  - b. Populatie en steekproef
9.  $V(X-Y) = V(X) + V(Y) - 2*COV(X,Y)$ 
  - a. Als X en Y onafhankelijk zijn  
dan  $V(X-Y) = V(X) + V(Y)$ 
    - i. Komt voort uit stelling 1 en 9
  - b. Bij populatie en steekproef

## 6. Bijzondere (kans)verdelingen

### 6.1. De binomiale verdeling

WAT?	<ul style="list-style-type: none"> <li>° Geeft de kans weer om k successen te halen bij N mogelijkheden</li> <li>→ Bv. Een meerkeuze-examen met N vragen, hoe groot is je kans om k correcte antwoorden te hebben?</li> </ul>
ILLUSTRATIE NODIG?	<p>(Cursus p.164-166!!)</p> <ul style="list-style-type: none"> <li>° Situatie : meerkeuze-examen, elke vraag heeft 4 antwoordmogelijkheden, proefpersonen moeten ad random antwoorden invullen</li> <li>° Stel dat er maar 1 vraag is: 25% van de proefpersonen zullen A aangeduid hebben, 25% B enzovoort</li> <li>→ Stel dat A het juiste antwoord is ; kans op succes = <math>p = 25\%</math></li> <li>° Er komt een 2<sup>e</sup> vraag bij: van de mensen die A geantwoord hebben bij vraag 1, gaat weer een kwart A bij de 2<sup>e</sup> vraag antwoorden, een kwart B enzovoort</li> <li>→ Antwoord A is weer het juiste antwoord</li> <li>° Nu zijn er 3 mogelijkheden</li> <li>→ Persoon heeft beide antwoorden fout <ul style="list-style-type: none"> <li>- Zowel op vraag 1 als 2 gekozen voor antwoord B, C of D</li> <li>- 9 van de 16 groepen</li> <li>- <math>P(X=0) = 9/16</math></li> </ul> </li> <li>→ Persoon heeft 1 van de 2 antwoorden juist <ul style="list-style-type: none"> <li>- OF op vraag 1 antwoord A en op vraag 2 geen A</li> <li>OF op vraag 1 geen A en op vraag 2 wel antwoord A</li> <li>- 6 van de 16 groepen</li> <li>- <math>P(X=1) = 6/16</math></li> </ul> </li> <li>→ Persoon heeft beide antwoorden juist <ul style="list-style-type: none"> <li>- op beide vragen antwoord A</li> </ul> </li> </ul>



	<ul style="list-style-type: none"> <li>- 1 van de 16 groepen</li> <li>- <math>P(X=2) = 1/16</math></li> </ul>
FORMULE	$P(X = k) = \frac{N!}{k! * (N-k)!} * p^k * (1 - p)^{N-k}$ <p>(Formularium!)</p> <ul style="list-style-type: none"> <li>° ! = faculteit : vermenigvuldiging, afrollend tot 1</li> <li>→ <math>N! = N \times (N-1) \times (N-2) \times (N-3) \times \dots \times 2 \times 1</math></li> <li>→ Bv. <math>4! = 4 \times 3 \times 2 \times 1 = 24</math></li> <li>→ <math>0! = 1</math></li> <li>° p = kans op succes</li> <li>° k = aantal successen</li> <li>° N = maximaal aantal successen</li> </ul>
NOTATIE & FORMULES	<ul style="list-style-type: none"> <li>° Variabele met binomiale verdeling</li> <li><math>X \sim \text{Binom}(N, p)</math></li> <li>° Verwachtingswaarde</li> <li><math>E(X) = N * p</math></li> <li>° Variantie</li> <li><math>V(X) = N * p * (1-p)</math></li> </ul>
VOORWAARDEN	<ul style="list-style-type: none"> <li>° N is vast</li> <li>° de kans op succes (p) blijft ongewijzigd</li> </ul>
VARIABELE?	Altijd bij discrete variabelen
PROBLEMEN/VRAGEN BIJ OEFENINGEN	<ul style="list-style-type: none"> <li>- Er is geen typische vorm qua grafiek</li> <li>- k en N bepalen: <ul style="list-style-type: none"> <li>° N is vaak je steekproefgrootte: hoe vaak de test gedaan wordt</li> <li>° k is hoeveel je daarvan nodig hebt</li> </ul> </li> </ul>

## 6.2. De normale verdeling

HEEL BELANGRIJK STUK!!

NORMAAL VERDEELDE VARIABELE	<ul style="list-style-type: none"> <li>° Notatie:</li> <li><math>X \sim N(\mu, \sigma^2)</math></li> <li>° Dichtheidsfunctie, formule</li> </ul>
DICHTHEIDSFUNCTIE	<ul style="list-style-type: none"> <li>° Formule</li> <li> <math display="block">f_X(x) = \frac{1}{\sigma * \sqrt{2\pi}} * e^{-\frac{(x-\mu)^2}{2\sigma^2}}</math> </li> <li><math>e \approx 2.71 ; \pi \approx 3.14</math></li> <li>° Symmetrisch</li> <li>° Hoogste punt : ter hoogte van <math>\mu</math> (op de x-as) – <u>mediaan = gemiddelde!!</u></li> <li>° Hoe groter <math>\sigma^2</math>, hoe breder en minder hoog de functie</li> <li>° De functie wordt nergens 0</li> <li>° Volledige oppervlakte is nog steeds gelijk aan 1</li> </ul>
VARIABELE?	Altijd bij continue variabelen

### 6.2.1. De standaardnormale verdeling

- De normale verdeling met  $\mu = 0$  en  $\sigma^2 = 1$ 
  - o Symmetrisch rond 0 (want dat is het gemiddelde)

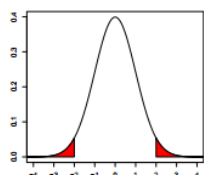
ENKELE EIGENSCHAPPEN..

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

Voorbeeld: de kans dat  $X$  groter is dan 2 (het rechter rode deel), is even groot als de kans dat  $X$  kleiner of gelijk aan -2 is (het linker rode deel)

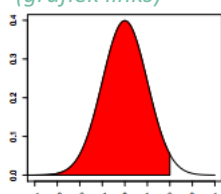
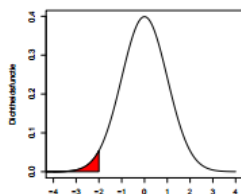
2)  $P(X \leq -x) = 1 - P(X \leq x)$

Voorbeeld: de kans dat  $X$  kleiner of gelijk aan -2 is (rode deel linker grafiek), is gelijk aan 1 – de kans dat  $X$  kleiner of gelijk aan 2 is (het witte deel van de rechter grafiek)



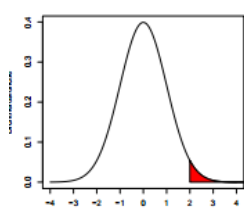
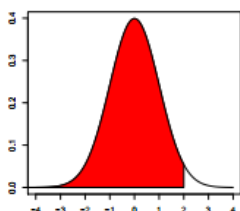
3)  $P(X > x) = 1 - P(X \leq x)$

Voorbeeld: de kans dat  $X$  groter is dan 2 (grafiek rechts) is gelijk aan 1 – de kans dat  $X$  kleiner of gelijk aan 2 is (grafiek links)



### 6.2.2. Kansen berekenen

- Er is een tabel die geldt voor de standaardnormale verdeling om kansen te berekenen



- o Zie formularium
- o OPGELET! Geldt enkel voor de standaardnormale verdeling
- o Het is de cumulatieve verdelingsfunctie van de standaardnormale variabele  $X$ 
  - de kansen die je uitkomt zijn telkens de kansen dat de variabele waarden aanneemt die kleiner of gelijk aan je specifieke  $x$ -waarde zijn :  $P(X \leq x)$

- Hoe de tabel gebruiken?

- o Je wil de kans berekenen dat  $X$  kleiner of gelijk aan  $x$  is :  $P(X \leq x)$
- o Je  $x$ -waarde ga je aflezen in de kolommen bovenaan en links
  - Linkerkant: de eerste twee cijfers van  $x$  (tot 1 getal na de komma)
  - Bovenaan: het laatste cijfer van  $x$  (het 2<sup>e</sup> getal na de komma)
- o Zoek het kruispunt van deze 2 kolommen in de tabel zelf
- o De waarde die je daar treft, is de kans dat  $X$  kleiner of gelijk aan je  $x$ -waarde is

- Wat als het geen standaardnormale verdeling is, maar een normale verdeling van een andere vorm?

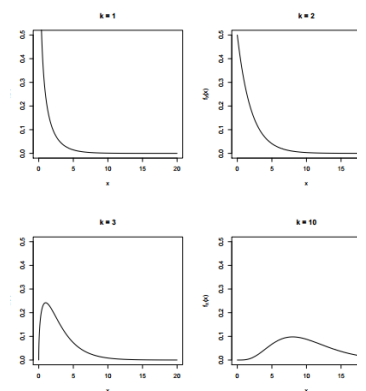
- o Dan moet je de variabele standaardiseren
- o *Eigenschap*: als  $X$  een niet-standaardnormale verdeling is, dan heeft de variabele  $Z$  wel een standaardnormale verdeling via deze formule:  $Z = \frac{X-\mu}{\sigma}$
- o Als  $X \sim N(\mu, \sigma^2)$ , dan geldt dat :  $P(X \leq x) = P(Z \leq \frac{x-\mu}{\sigma})$ 
  - De waarde die je uitkomt na het uitvoeren van de bewerking  $\frac{x-\mu}{\sigma}$  (die dus rechts van het ongelijkheidsteken bij  $Z$  staat) is de  $x$ -waarde die je wil zoeken in de zijkanten van je tabel in het formularium
  - Van daaruit is het weer hetzelfde: het getal in de tabel dat overeenkomt met die  $x$ -waarde, is de kans dat  $P(X \leq x)$

o *Illustratie*:

- $P(X \leq 3) = P\left(Z \leq \frac{3-1}{\sqrt{4}}\right) = P(Z \leq 1)$
- Je zoekt 1 dan als x-waarde op in de tabel
- De bijhorende kans is 0.8413
- $P(X \leq 3) = 0.8413$

### 6.3. De $\chi^2$ -verdeling

WAT?	<ul style="list-style-type: none"> <li>° Je neemt k aantal onafhankelijke standaardnormale variabelen</li> <li>° de <math>\chi^2</math>-verdeling is de som van die gekwadrateerde variabelen → som van kwadraten: altijd alleen maar positieve waarden</li> </ul>
NOTATIE	$Y \sim \chi_k^2$
FORMULES	<ul style="list-style-type: none"> <li>° Algemene formule: <math display="block">Y = X_1^2 + X_2^2 + X_3^2 + \dots + X_k^2</math></li> <li>° Verwachtingswaarde: <math>E(Y) = k</math></li> <li>° Variantie: <math>V(Y) = 2 \cdot k</math></li> </ul>
K?	k = het aantal vrijheidsgraden → hoeveel er zijn, dus.
VARIABLEN?	Altijd bij continue variabelen
PROBLEMEN/VRAGEN BIJ OEFENINGEN	Er is geen typische vorm qua grafiek

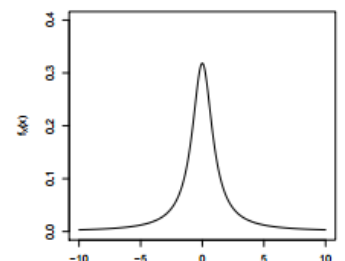


#### 6.3.1. Kansen berekenen

- Zie tabel in formularium
  - o Werkt wel wat anders dan tabel voor normale verdeling!
  - o Linker kolom: het aantal vrijheidswaarden k zoeken
  - o Bovenste kolom: de waarden van de verdelingsfunctie  $F_Y(y)$
  - o Tabel zelf: de waarden y van de variabele → datgene wat tussen haakjes staat als 'y' bij ' $F_Y(y)$ '

### 6.4. De t-verdeling

WAT?	<ul style="list-style-type: none"> <li>° Er is een standaardnormale variabele en een variabele volgens de <math>\chi^2</math>-verdeling.</li> <li>° De variabelen zijn onderling onafhankelijk.</li> <li>° De <math>t_k</math>-verdeling is de verdeling van de variabele <math display="block">T = \frac{X}{\sqrt{\frac{1}{k} Y}}</math></li> </ul>
NOTATIE	$T \sim t_k$
FORMULES	<ul style="list-style-type: none"> <li>° Verwachtingswaarde <math>E(T) = 0</math></li> <li>° Variantie <math>V(T) = \frac{k}{k-2}</math> voor <math>k &gt; 2</math></li> </ul>
VARIABLEN?	Altijd bij continue variabelen



### 6.4.1. Kansen berekenen

- Zie tabel in formularium
  - o Zelfde logica als bij chi-kwadraat verdeling
  - o Linker kolom: het aantal vrijheidswaarden  $k$  zoeken
  - o Bovenste kolom: de waarden van de verdelingsfunctie  $F_T(t)$
  - o Tabel zelf: de waarden van de variabele  $\rightarrow$  datgene wat tussen haakjes staat als 't' bij ' $F_T(t)$ '

## HOOFDSTUK 6: DE STEEKPROEVENVERDELING

### ALGEMEEN

Wat bestuderen we?	Eigenschappen van variabelen die we bekommen door op willekeurige wijze een steekproef te trekken uit de populatie
Dus...?	<u>We kijken ook naar een steekproef, alleen is het één van de ontelbaar mogelijke steekproeven uit de populatie. We weten niet exact wat de waarden zijn, het is een stuk abstracter en algemener</u>
Wat belangrijk?	Reproduceerbaarheid: als je zelfde experiment zou uitvoeren met andere steekproef, gelijkaardige resultaten bekommen
→ Probleem	Vaak maar tijd & geld om experiment 1x uit te voeren
→ Oplossing	statistische formules

### 1. Steekproeftrekking

ASELECTE STEEKPROEFTREKKING	<ul style="list-style-type: none"> <li>° Uit de populatie worden op random wijze <math>n</math> elementen geselecteerd  <math>\rightarrow</math> Die <math>n</math> elementen zijn onderling afhankelijk (daar gaan we toch van uit)</li> <li>° Er zijn veel meer soorten steekproeftrekkingen, wij bespreken alleen deze</li> <li>° Alle methoden &amp; formules die we zien, zijn van toepassing op aselecte steekproef</li> </ul>
NOTATIE	<u>Alles met een hoofdletter schrijven</u> <ul style="list-style-type: none"> <li>° De variabele: hoofdletter <math>X</math></li> <li>° De specifieke waarden: <math>X_1, X_2, X_3 \dots</math>  <math>\rightarrow</math> Bij één expliciete steekproef: <math>x_1, x_2, x_3 \dots</math></li> </ul>
TOEVALSVARIABELE	Een variabele die bekomen wordt door op toevallige wijze een steekproef uit de populatie te selecteren
KANSBEREKENING	$P(X = x_i)$ kan 2 betekenissen hebben: <ul style="list-style-type: none"> <li>° De relatieve frequentie in de populatie</li> <li>° Kansberekening  <math>\rightarrow</math> Kans op gebeurtenis = rel.fr. van de gebeurtenis als experiment <math>\infty</math> keer herhaald wordt  <math>\rightarrow</math> In praktijk is dat onmogelijk: we moeten dit proberen te benaderen              * Hoe vaker je een experiment uitvoert, hoe dichter bij <math>\infty</math>, hoe dichter bij de echte kans, hoe betrouwbaarder</li> </ul>

→ Een experiment meerdere keren herhalen = herhaalde steekproeftrekking

## 2. Steekproevenverdeling van het gemiddelde

NOTATIE	$\bar{X}$ (van een steekproef in het algemeen dus, niet van specifieke waarden in 1 steekproef)
VARIABELE	<ul style="list-style-type: none"> <li>° Als je een herhaalde steekproeftrekking doet, dan merk je dat de waarden voor <math>\bar{X}</math> veranderen</li> <li>→ Logisch, want de waarden van de steekproef zijn telkens anders</li> <li>→ <math>\bar{X}</math> verandert, dus is letterlijk variabel</li> </ul>
FORMULE	$\bar{X} = \frac{1}{n} * \sum_{i=1}^n X_i$
STEEKPROEFGROOTHEID OF EEN STATISTIEK	<ul style="list-style-type: none"> <li>° Een bewerking toegepast op de variabelen <math>X_1, \dots, X_n</math></li> <li>° <math>\bar{X}</math> is hier een voorbeeld van</li> <li>° Andere voorbeelden : modus, variantie..</li> </ul>
STEEKPROEVENVERDELING	<ul style="list-style-type: none"> <li>° Uiteindelijk willen we toch weten hoe alle verschillende gemiddelden vd steekproeven vd populatie zich verhouden onder elkaar!</li> <li>° Stappenplan: <ul style="list-style-type: none"> <li>→ Oneindig aantal steekproeven trekken &amp; gemiddelden berekenen</li> <li>→ Histogram opstellen van alle steekproefgemiddelden (is dan een histogram met oneindig veel klassen)</li> <li>→ Wat krijgen we: de dichtheidsfunctie van het gemiddelde of de steekproevenverdeling van het gemiddelde</li> </ul> </li> </ul> <p>OPGELET: verschil frequentieverdeling – steekproevenverdeling</p> <ul style="list-style-type: none"> <li>° frequentieverdeling = verdeling van een variabele</li> <li>° steekproevenverdeling = verdeling van een steekproefgrootheid</li> </ul>

### 2.1 Stellingen

$E(\bar{X}) = \mu_X$	<p>De verwachtingswaarde van het steekproefgemiddelde <math>\bar{X}</math> = het populatiegemiddelde van de variabele X</p> <p>→ Het steekproefgemiddelde is een zuivere schatter voor het populatiegemiddelde!</p>
$V(\bar{X}) = \frac{\sigma_X^2}{n}$	<p>De variantie van het steekproefgemiddelde = populatievariantie van de variabele, gedeeld door de steekproefgrootte</p> <p>→ Logica: hoe groter de steekproef, hoe minder de gemiddeldes zullen variëren, hoe dichter bij het échte steekproefgemiddelde</p>
$\bar{X} \sim N(\mu_X, \frac{\sigma_X^2}{n})$	<p>ALS <math>X_1, \dots, X_n</math> random, onafhankelijk en normaal verdeeld zijn</p> <p>DAN is <math>\bar{X}</math> ook normaal verdeeld: <math>\bar{X} \sim N(\mu_X, \frac{\sigma_X^2}{n})</math></p> <p>→ Logica: als populatie normaal verdeeld is, dan steekproefgemiddelde ook</p> <p>→ Gaat op voor elke grootte van n, maar alleen als het normaal verdeeld is</p>
<b>CENTRALE LIMIETSTELLING</b>	<p>Het maakt niet uit hoe de steekproef verdeeld is</p> <p>ALS de steekproefgrootte groot genoeg is (vuistregel: &gt; 30)</p> <p>DAN is het steekproefgemiddelde (+/-) normaal verdeeld</p> <p>→ We kunnen deze dan, net als gewone normale verdelingen, standaardiseren &amp; op die manier kansen berekenen!</p> <p>→ Gaat op voor elke soort verdeling, maar alleen als n groot genoeg is</p>

## 2.2 OPGELET!!!

Bij het standaardiseren van  $\bar{X}$ : je moet de waarden van  $\bar{X}$  gebruiken, niet van  $X$ !!

→ Het gemiddelde blijft hetzelfde (zie formule)

→ de standaarddeviatie blijft NIET hetzelfde!! Je moet deze berekenen via de formule!

ALGEMENE FORMULE:  $Z = \frac{\bar{X} - \mu_X}{\frac{\sigma_X}{\sqrt{n}}}$

## 3. Steekproevenverdeling van de variantie

ALGEMEEN	° Ook een voorbeeld van een steekproefgrootte of een statistiek ° Ook een variabele
NOTATIE	$SD_X^2$ of $S_X^2$
FORMULES	$SD_X^2 = \frac{1}{n} * \sum_{i=1}^n (X_i - \bar{X})^2$ $S_X^2 = \frac{1}{n-1} * \sum_{i=1}^n (X_i - \bar{X})^2$ $E(SD_X^2) = \frac{n-1}{n} * \sigma_X^2$ $E(S_X^2) = \sigma_X^2$

### 3.1 Stelling

$\frac{(n-1) * S_X^2}{\sigma_X^2} \sim \chi_{n-1}^2$	ALS $X_1, \dots, X_n$ random, onafhankelijk en normaal verdeeld zijn DAN geldt $\frac{(n-1) * S_X^2}{\sigma_X^2} \sim \chi_{n-1}^2$
--	--

## HOOFDSTUK 7: BETROUWBAARHEIDSINTERVALLEN EN STATISTISCHE TOETSEN VOOR HET POPULATIEGEMIDDELDE

### ALGEMEEN – Waar gaat het hoofdstuk over?

- We gaan proberen om op basis van een steekproef een uitspraak te formuleren over de populatie.
  - Eerst moeten we populatieparameter schatten op basis van resultaten uit steekproef
  - Dan Betrouwbaarheidsinterval opstellen & statistische toets gebruiken
  - Je kan nooit 100% zeker zijn dat je uitspraak over de populatie correct is
- Hier gaan we dat enkel doen bij het populatiegemiddelde. (Bij andere populatieparameters kan dat ook, maar dat is voor latere cursussen)

### 1. Schatters

NOTATIE	- Een populatieparameter: $\theta$ → een populatieparameter is een heel algemeen woord voor iets wat iets zegt over de
---------	---

	toestand van uw populatie zelf. Bijvoorbeeld: het gemiddelde, de variantie..
WAT MAAKT EEN GOEIE SCHATTER?	<ul style="list-style-type: none"> <li>- Een schatter van een populatieparameter: <math>\hat{\theta}</math></li> <li>- Ze moet <u>zuiver</u> zijn <ul style="list-style-type: none"> <li>◦ De verwachtingswaarde van de schatter = de populatieparameter</li> <li>◦ <math>E(\hat{\theta}) = \theta</math></li> <li>◦ Dat geeft aan dat de populatieparameter niet systematisch te groot of te klein wordt geschat</li> </ul> </li> <li>- De variantie van de schatter wordt kleiner naarmate de steekproefgrootte groter wordt <ul style="list-style-type: none"> <li>◦ <math>V(\hat{\theta}) \downarrow</math> als <math>n \uparrow</math></li> <li>◦ Nauwkeuriger, want meer info</li> </ul> </li> </ul>
$\sqrt{V(\hat{\theta})}$	<ul style="list-style-type: none"> <li>◦ Standaarddeviatie van de schatter</li> <li>◦ <u>Standaardfout</u> van de schatter</li> <li>→ De schatter met de kleinste standaardfout is de beste: het efficiëntste</li> </ul>
SCHATTER VS SCHATTING?	<p>Schatter: algemeen, veranderlijk</p> <p>Schatting: van 1 specifieke steekproef, vaste waarde</p>

### 1.1 Het gemiddelde

HOE GAAN WE POPULATIEGEMIDDELDE SCHATTEN?	Een logische optie: het steekproefgemiddelde gebruiken
IS HET STEEKPROEFGEMIDDELDE WEL EEN GOEDE SCHATTER?	<p>Voldoet het aan de 2 voorwaarden van een goede schatter?</p> <ul style="list-style-type: none"> <li>◦ Is ze zuiver? <ul style="list-style-type: none"> <li>→ <math>E(\bar{X}) = \mu_X</math></li> <li>→ Dat was een stelling uit het vorige hoofdstuk</li> </ul> </li> <li>◦ Daalt de variantie als de steekproefgrootte stijgt? <ul style="list-style-type: none"> <li>→ Stelling uit vorig hoofdstuk: <math>V(\bar{X}) = \frac{\sigma_X^2}{n}</math></li> <li>→ Als n groter wordt, is <math>\frac{\sigma_X^2}{n}</math> kleiner, dus minder variantie</li> </ul> </li> </ul>
CONCLUSIE	<p>Ja, het <u>steekproefgemiddelde is een goede schatter!</u></p> <p>→ Je wilt populatiegemiddelde schatten: gebruik steekproefgemiddelde</p>

### 1.2 De variantie

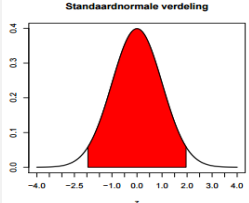
HOE GAAN WE POPULATIEVARIANTIE SCHATTEN?	Een logische optie: de steekproefvariantie gebruiken
PROBLEEM – WELKE FORMULE GEBRUIKEN?	<p>We hebben voor de steekproefvariantie 2 formules: <math>S_X^2</math> en <math>SD_X^2</math></p> <p>→ Stelling uit vorig hoofdstuk: <math>E(S_X^2) = \sigma_X^2</math></p>
CONCLUSIE	<p><u><math>S_X^2</math> is een goede schatter voor de populatievariantie</u></p> <p><math>SD_X^2</math> is GEEN goede schatter voor de populatievariantie</p>

## 2. Betrouwbaarheidsintervallen

- De steekproevenverdeling laat ons toe betrouwbaarheidsintervallen te construeren
- Via een betrouwbaarheidsinterval kunnen we met een bepaalde zekerheid een uitspraak doen over het populatiegemiddelde
- Er zijn verschillende werkwijzen, afhankelijk van de verdeling en de kennis over de populatievariantie



## 2.1 X normaal verdeeld en gekende populatievariantie

<p><math>z_\alpha</math></p>	<ul style="list-style-type: none"> <li>◦ De waarde van de standaardnormale verdeling, waarvoor de oppervlakte onder de curve rechts ervan gelijk is aan <math>\alpha</math></li> <li>◦ <math>P(Z &gt; z_\alpha) = \alpha</math> met <math>Z \sim N(0,1)</math></li> <li>◦ OPGELET: als je in de tabel waarden wilt aflezen, gaat het altijd over 'inclusief alle waarden rechts ervan'             <ul style="list-style-type: none"> <li>→ je moet gebruik maken van <math>1 - P(Z \leq z_\alpha)</math></li> </ul> </li> <li>◦ De standaardnormale verdeling is symmetrisch rond 0:             <ul style="list-style-type: none"> <li>→ <math>P(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) = 1 - \alpha</math></li> <li>→ of <math>P(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}) = 1 - \alpha</math></li> <li>* We gaan standaardiseren, zodat het van toepassing is op een standaardnormale verdeling (is vereist, zie bovenste zin)</li> <li>→ Waarom <math>\alpha/2</math>?</li> </ul> </li> </ul> <div style="display: flex; align-items: center;">  <div style="border: 1px solid black; padding: 5px; margin-left: 10px;"> <p>Het is <math>z_{\alpha/2}</math> want hier is <math>\alpha</math> de som van de 2 witte gebieden, zowel links als rechts! Dus uw boven- of onderwaarde apart moet je delen door 2</p> </div> </div>
<p>BETROUWBAARHEIDS-INTERVAL BI</p>	<p><math>P(\bar{X} - z_{\alpha/2} * \sigma/\sqrt{n} \leq \mu \leq \bar{X} + z_{\alpha/2} * \sigma/\sqrt{n}) = 1 - \alpha</math></p> <p>→ Logica van de formule:</p> <ul style="list-style-type: none"> <li>◦ Zie cursus p.211</li> <li>◦ Je hebt alle elementen van de Z-score aan alle kanten van de ongelijkheidstekens toegevoegd, zodat je uiteindelijk in het midden alleen <math>\mu</math> uitkomt, want dat is uiteindelijk wat je zoekt</li> </ul> <p>→ Uiteindelijk wil je dat je getal dat voor <math>\mu</math> staat tussen die 2 grenzen ligt</p> <p>Grenzen van het betrouwbaarheidsinterval: <math>[\bar{X} - z_{\alpha/2} * \sigma/\sqrt{n}, \bar{X} + z_{\alpha/2} * \sigma/\sqrt{n}]</math></p> <p>→ Deze formule ga ik vooral gebruiken bij oefeningen</p> <p>Dat is dan het <math>(1-\alpha)*100\%</math> betrouwbaarheidsinterval</p> <p>→ Bv. <math>\alpha = 0.05</math> : het <math>(1-0.05)*100\%</math> BI : het 95%-BI</p> <p>→ <u>Wat betekent het? Dat is 95% van de gevallen van de steekproeftrekkingen het reële populatiegemiddelde er zich ook echt in bevindt</u></p> <p>Het steekproefgemiddelde gaat altijd exact in het midden van het BI liggen</p>
<p>INTERPRETATIE</p>	<p>Wat gebeurt er als het experiment herhaalt op basis van een nieuwe steekproef?</p> <p>→ <math>\sigma</math> en <math>z_{\alpha/2}</math> zijn vaste waarden</p> <p>→ Maar nieuwe steekproef = nieuwe gegevens = nieuw gemiddelde</p> <p>Resultaat: gezien het gemiddelde variabel is, zal het betrouwbaarheidsinterval ook variabel zijn</p> <p>→ de grenzen van het betrouwbaarheidsinterval zullen verschillen per steekproef</p> <p>→ een 95%-betrouwbaarheidsinterval garandeert dat 95% van al die variabele intervallen het reële populatiegemiddelde zullen bevatten als we het experiment een oneindig aantal keer zullen herhalen</p>
<p>EIGENSCHAPPEN</p>	<p>Breedte van het interval <math>[a, b] = 2 * z_{\alpha/2} * \sigma/\sqrt{n}</math></p> <p>→ Hoe smaller het interval, hoe nauwkeuriger</p> <p>→ Het kleinste betrouwbaarheidsinterval mogelijk: <math>\alpha = 0.01</math></p> <ul style="list-style-type: none"> <li>◦ interval is dan 99% betrouwbaar</li> </ul>

## 2.2 X normaal verdeeld en ongekende populatievariantie

FORMULE	<ul style="list-style-type: none"> <li>◦ <math>P(\bar{X} - t_{n-1;\alpha/2} * S_X/\sqrt{n} \leq \mu \leq \bar{X} + t_{n-1;\alpha/2} * S_X/\sqrt{n}) = 1-\alpha</math></li> <li>◦ Grenzen betrouwbaarheidsinterval:  <math>[\bar{X} - t_{n-1;\alpha/2} * S_X/\sqrt{n}, \bar{X} + t_{n-1;\alpha/2} * S_X/\sqrt{n}]</math></li> </ul>
LOGICA ACHTER FORMULE	<ul style="list-style-type: none"> <li>◦ Combinatie van 2 stellingen: <ul style="list-style-type: none"> <li>→ <math>\frac{(n-1)*S_X^2}{\sigma_X^2} \sim \chi_{n-1}^2</math></li> <li>→ <math>\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim N(0,1)</math></li> </ul> </li> <li>◦ Een <math>\chi^2</math>-verdeling en een standaardnormale verdeling vormen samen een t-verdeling</li> <li>◦ Het is uiteindelijk een <math>t_{n-1}</math>-verdeling, omdat het in de eerste stelling ook een <math>\chi_{n-1}^2</math>-verdeling is</li> </ul>
EIGENSCHAPPEN $t_{n-1}$ -VERDELING	<ul style="list-style-type: none"> <li>◦ Een grotere variantie dan een standaardnormale verdeling</li> <li>◦ <math>t_{n-1;\alpha/2}</math>-waarde is groter dan <math>z_{\alpha/2}</math>-waarde <ul style="list-style-type: none"> <li>→ OPGELET: ook hier is de <math>t_{n-1;\alpha/2}</math>-waarde de waarde voor alles rechts ervan, niet links! Moet je weer rekening mee houden bij het aflezen van de tabel</li> </ul> </li> <li>◦ Het betrouwbaarheidsinterval is groter, dus extra variabiliteit <ul style="list-style-type: none"> <li>→ Logisch, want je het de populatiestandaarddeviatie moeten schatten, waardoor een grotere kans op fouten</li> </ul> </li> <li>◦ OPGELET: die 'n-1' betekent niet dat je voor het aflezen van de tabel er eentje moet bijtellen bij het aantal vrijheidswaarden! 'n-1'=k!</li> </ul>

## 2.3 X niet normaal verdeeld & populatievariantie niet gekend

GROTE STEEKPROEF	<ul style="list-style-type: none"> <li>° Centrale limietstelling gebruiken: <math>\bar{X}</math> is normaal verdeeld</li> <li>° Zelfde werkwijze vorige titel: → interval BI: <math>[\bar{X} - t_{n-1;\alpha/2} * S_X/\sqrt{n}, \bar{X} + t_{n-1;\alpha/2} * S_X/\sqrt{n}]</math></li> </ul>
KLEINE STEEKPROEF	Wordt niks over gezegd, moeten we niet kunnen..

## 2.4 Algemeen overzicht werkwijzen

<b>X NORMAAL VERDEELD</b> $\sigma_X^2$ GEKEND	Grenzen betrouwbaarheidsinterval: $[\bar{X} - z_{\alpha/2} * \sigma/\sqrt{n}, \bar{X} + z_{\alpha/2} * \sigma/\sqrt{n}]$
<b>X NORMAAL VERDEELD</b> $\sigma_X^2$ ONGEKEND	Grenzen BI : $[\bar{X} - t_{n-1;\alpha/2} * S_X/\sqrt{n}, \bar{X} + t_{n-1;\alpha/2} * S_X/\sqrt{n}]$
<b>X NIET-NORMAAL VERDEELD</b> $\sigma_X^2$ ONGEKEND	Grote steekproef: centrale limietstelling Grenzen BI: $[\bar{X} - t_{n-1;\alpha/2} * S_X/\sqrt{n}, \bar{X} + t_{n-1;\alpha/2} * S_X/\sqrt{n}]$ (Kleine steekproef: niet kunnen)

## 3. Statistische toetsen

ALGEMENE VOORWAARDEN VOOR HET GEBRUIK VAN STATISTISCHE TOETSEN	<ul style="list-style-type: none"> <li>° Het moet over <u>1 grote steekproef</u> gaan</li> <li>° De variabele is <u>normaal verdeeld</u></li> <li>° <math>\sigma_X^2</math> is <u>ongekend</u></li> </ul>
DE 2 HYPOTHESEN	<p>Er zijn telkens 2 stellingen ivm de hypothese, 1 van de 2 is juist:</p> <ul style="list-style-type: none"> <li>° Nulhypothese <math>H_0</math> → <math>H_0: \mu = \mu_0</math> → We bewijzen dat het reële populatiegemiddelde gelijk is aan het geschatte populatiegemiddelde → In praktijk is het erg vaak dat we deze hypothese proberen te verwerpen</li> <li>° Alternatieve hypothese <math>H_a</math> → <math>H_a: \mu \neq \mu_0</math> → We bewijzen dat het reële populatiegemiddelde niet gelijk is aan het geschatte populatiegemiddelde → Ook wel de tweezijdig alternatieve hypothese genoemd: <math>\mu_0</math> kan uiteindelijk langs beide kanten van <math>\mu</math> liggen, erboven of eronder</li> </ul>
EN DAN?	<ul style="list-style-type: none"> <li>° Als <math>\bar{x}</math> ongeveer gelijk is aan <math>\mu_0</math>, dan gaan we <math>H_0</math> aanvaarden &amp; <math>H_a</math> verwerpen</li> <li>° Als <math>\bar{x}</math> ver afwijkt van <math>\mu_0</math>, dan gaan we <math>H_0</math> verwerpen &amp; <math>H_a</math> aanvaarden</li> <li>° MAAR nu is de vraag: wat is 'ongeveer gelijk aan' en 'ver afwijken van'?? → Gaan we bewijzen via een statistische toets: een toetsingsgrootheid → We gaan proberen voldoende bewijs te vinden om <math>H_0</math> te verwerpen</li> </ul>

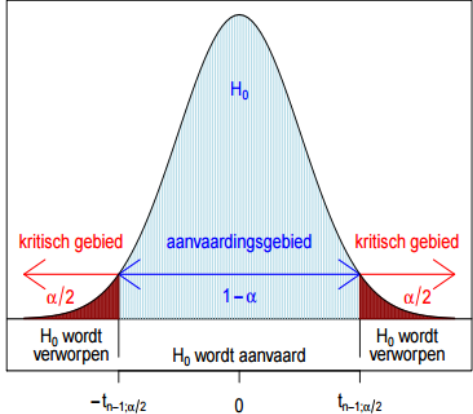
### 3.1 Toetsingsgrootheid

EXTRA VARIABELE	We gaan een extra variabele G invoeren:
-----------------	---

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

→ Als  $H_0$  correct is, dan volgt G een  $t_{n-1}$ -verdeling

### 3.2 Beslissingsregels

WAT?	Regels die je moet volgen om te beslissen of we $H_0$ gaan aanvaarden of verwerpen
HOE GELDEN ZE?	Als $H_0$ waar is, dan ligt $G$ rond 0 Als $H_0$ niet waar is, dan ligt $G$ ver van 0 (zowel positief als negatief kan)
BESLISSINGSINTERVAL & KRITISCHE WAARDEN	<p>Werkt ook langs de andere kant:</p> <ul style="list-style-type: none"> <li>- Als <math>g</math> rond 0 ligt, dan is <math>H_0</math> waar</li> <li>- als <math>g</math> ver van 0 ligt, dan is <math>H_0</math> fout</li> </ul> <p><math>-t_{n-1;\alpha/2} \leq g \leq t_{n-1;\alpha/2}</math>  <math>\rightarrow</math> Als je <math>g</math>-waarde (die je bekomt via formule <math>G = \frac{\bar{x} - \mu_0}{s_x/\sqrt{n}}</math>) tussen deze 2 waarden ligt, dan aanvaarden we <math>H_0</math></p> <p>Kritische waarden: <math>-t_{n-1;\alpha/2}</math> en <math>t_{n-1;\alpha/2}</math></p> 

### 7.3 Type I en type II fout

TYPE I FOUT	<ul style="list-style-type: none"> <li>◦ <math>H_0</math> verwerpen, maar eigenlijk is die juist</li> <li>◦ De kans hiertoe : <math>P(\text{verwerp } H_0 \mid \mu = \mu_0) = \alpha</math>  <math>\rightarrow H_0</math> correct aanvaarden: <math>P(\text{aanvaard } H_0 \mid \mu = \mu_0) = 1 - \alpha</math></li> </ul>
TYPE II FOUT	<ul style="list-style-type: none"> <li>◦ <math>H_0</math> aanvaarden, maar eigenlijk is die fout</li> <li>◦ De kans hiertoe: : <math>P(\text{aanvaard } H_0 \mid \mu \neq \mu_0) = \beta</math>  <math>\rightarrow H_0</math> correct verwerpen: <math>P(\text{verwerp } H_0 \mid \mu \neq \mu_0) = 1 - \beta</math></li> </ul>
ONDERLING VERBAND	<ul style="list-style-type: none"> <li>◦ <math>\alpha</math> daalt = <math>\beta</math> stijgt</li> <li>◦ <math>n</math> stijgt = <math>\beta</math> daalt</li> </ul>

	In werkelijkheid is $H_0$	
	juist	fout
We aanvaarden $H_0$	Juiste beslissing (A) Betrouwbaarheid $(1 - \alpha)$	Foute beslissing (C) Type II fout $\beta$
We verwerpen $H_0$	Foute beslissing (B) Type I fout $\alpha$	Juiste beslissing (D) Onderscheidingsvermogen $(1 - \beta)$

### 3.4 Beslissingsregels op basis van het betrouwbaarheidsinterval

$\mu_0$ ligt binnen betrouwbaarheidsinterval	$H_0$ aanvaarden
$\mu_0$ ligt buiten betrouwbaarheidsinterval	$H_0$ verwerpen

# Extra info: rekenmachine

Je kan redelijk wat zaken ook ingeven in de GRM en daarna info hierover aflezen

!! OPGELET Ik denk niet dat we een grafische rekenmachine mogen gebruiken op het examen !!

- Info ingeven in lijsten
  - o 1-Var-Stat = in verband met 1 variabele
  - o 2-Var-Stat = in verband met 2 variabelen
    - $'a' = b_1$
    - $'b' = b_0$
    - Deze twee kan je gebruiken om het functievoorschrift van de regressielijn te noteren
- $'s_X'$  = de standaarddeviatie die in de formule deelt door  $'n - 1'$
- $'\sigma_X'$  = de standaarddeviatie die in de formule deelt door  $'n'$