



Combinatoriek

Volgorde van belang?	herhaling toegestaan?
✓	✗
✓	✗
✗	✗
✓	✓
✗	✓

Variatie: $V_n^p = \frac{n!}{(n-p)!}$

↳ ketting met parels

Permutatie: $P_n = V_n^n = n!$

↳ n elementen ordenen

Combinatie: $C_n^p = \binom{n}{p} = \frac{n!}{p!(n-p)!}$

↳ "p kiezen uit n"

Herhalingsvariatie: $\bar{V}_n^p = n^p$

↳ nummerplaten

Herhalingscombinatie: $\bar{C}_n^p = \binom{n+p-1}{p}$

↳ "stipjes & streepjes"

Herhalingspermutatie: $\bar{P}_n^{p,q,\dots,r} = \binom{n}{p, q, \dots, r} = \frac{n!}{p!q!\dots r!}$

↳ anagrammen

Het aantal deelverzamelingen van een verzameling met n elementen = $2^n = \#D(A)$

Binomium van Newton

$$(a+b)^n = \sum_{i=0}^n \binom{n}{i} a^{n-i} \cdot b^i$$

Multinomiale ontwikkeling

$$(x_1 + x_2 + \dots + x_k)^n = \sum (n_1, n_2, \dots, n_k) x_1^{n_1} x_2^{n_2} \dots x_k^{n_k} \quad (n_1 + n_2 + \dots + n_k = n)$$

①

Kansruimten

Basisbegrippen

- **universum** Ω : alle mogelijke uitkomsten
- **gebeurtenissen**: deelverz. v.h. universum
- **sigma-algebra's van deelverzamelingen**: een klasse \mathcal{A} van deelverz. van Ω heet een σ -algebra als het aan de volgende axioma's voldoet:
 - ① $\Omega \in \mathcal{A}$
 - ② $A \in \mathcal{A} \Rightarrow A^c \in \mathcal{A}$
 - ③ $\forall n \in \mathbb{N}: A_n \in \mathcal{A} \Rightarrow \bigcup_{n \in \mathbb{N}} A_n \in \mathcal{A}$
 ↳ (Ω, \mathcal{A}) is een **meetbare ruimte**
 - ↳ $\forall n \in \mathbb{N}: A_n \in \mathcal{A} \Rightarrow \bigcap_{n \in \mathbb{N}} A_n \in \mathcal{A}$
 - ↳ $\forall A, B \in \mathcal{A}: A \Delta B = (A \cup B) \setminus (A \cap B) \in \mathcal{A}$
- **discrete σ -algebra**: machtsverz. = $\mathcal{D}(\Omega) = \mathcal{P}(\Omega) = 2^\Omega$
- $P: \mathcal{A} \rightarrow \mathbb{R}$ is een **kansmaat** indien:
 - ① $P(\Omega) = 1$
 - ② $\forall A \in \mathcal{A}: P(A) \geq 0$
 - ③ $\forall n \in \mathbb{N}: A_n \in \mathcal{A} \vee \forall i \neq j: A_i \cap A_j = \emptyset$ (paarsgewijs disjunct)
 - $\Rightarrow P(\bigcup_{n \in \mathbb{N}} A_n) = \sum_{n \in \mathbb{N}} P(A_n)$
 - ↳ axioma van σ -additiviteit
- **kansruimte** (Ω, \mathcal{A}, P) : Ω = universum, \mathcal{A} = σ -algebra, P = kansmaat over Ω
- **stijgende/dalende rij van verzamelingen**:

$$\lim_{n \rightarrow \infty} A_n = \begin{cases} \bigcup_{n=1}^{\infty} A_n & \text{als } A_n \uparrow \\ \bigcap_{n=1}^{\infty} A_n & \text{als } A_n \downarrow \end{cases}$$

st. Zij (Ω, \mathcal{A}, P) een kansruimte:

- ① **Eindige additiviteit**: Als $\{A_n | n \in \{1, \dots, N\}\}$ paarsgewijs disjunct en $\forall n \in \{1, \dots, N\}, A_n \in \mathcal{A}$, dan

$$P(\bigcup_{n=1}^N A_n) = \sum_{n=1}^N P(A_n)$$
- ② $\forall A \in \mathcal{A}: P(A^c) = 1 - P(A)$
- ③ Als $\forall n \in \mathbb{N}_0: A_n \in \mathcal{A}$ en (A_n) monotoon, dan

$$P(\lim_{n \rightarrow \infty} A_n) = \lim_{n \rightarrow \infty} P(A_n)$$
 - ↳ $\forall A, B \in \mathcal{A}$:
 - $P(B) = P(B \cap A) + P(B \cap A^c)$
 - $A \subset B \Rightarrow P(A) \leq P(B)$
 - $P(A) \leq 1$
 - $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

niet-aftelbaar universum

- ↳ gebeurtenissen in vorm $]-\infty, a]$
- ↳ σ -algebra voortgebracht door een collectie verzamelingen $\mathcal{C} \subset 2^\Omega$ = de kleinste σ -algebra die \mathcal{C} bevat:
 - $\sigma(\mathcal{C})$ is een σ -algebra
 - $\mathcal{C} \subset \sigma(\mathcal{C})$
 - $\forall \sigma$ -algebra \mathcal{A} met $\mathcal{C} \subset \mathcal{A}: \sigma(\mathcal{C}) \subset \mathcal{A}$

eig. Voor elke $\mathcal{C} \subset 2^\Omega$ bestaat $\sigma(\mathcal{C})$
bewijs p. 13

uniekheid van $\sigma(\mathcal{C})$

Voorwaardelijke kans en onafhankelijkheid

def. kans op A , gegeven B met $P(B) \neq 0$:

$$P(A|B) = P(A \cap B) / P(B)$$

$$\hookrightarrow P(A|A) = 1, P(A|\Omega) = P(A)$$

st. Zij (Ω, \mathcal{A}, P) en $B \in \mathcal{A}$ met $P(B) > 0$, dan definieert

$P_B: \mathcal{A} \rightarrow \mathbb{R}: A \mapsto P(A|B)$ een kansmaat op \mathcal{A} .

st. • **kettingregel**: Zij (Ω, \mathcal{A}, P) , $(A_n)_{n=1}^k \in \mathcal{A}$ met $2 \leq k \in \mathbb{N}$, dan:

$$P(A_1 \cap A_2 \cap \dots \cap A_n) = P(A_1) P(A_2|A_1) P(A_3|A_1 \cap A_2) \dots P(A_k|A_1 \cap A_2 \cap \dots \cap A_{k-1})$$

bewijs via inductie p. 18

st. • **wet van de totale kans**: Zij (Ω, \mathcal{A}, P) en $\{A_n\}_{n \in \mathbb{N}}$ een partitie van Ω met alle $P(A_n) \neq 0$, dan geldt:

$$\forall B \in \mathcal{A}: P(B) = \sum_{n \in \mathbb{N}} P(A_n) P(B|A_n)$$

\hookrightarrow een partitie is een opdeling van een verzameling in niet-lege, paarsgewijs disjuncte deelverzamelingen

st. • **stelling van Bayes**: zie voorwaarden vorige stelling

$$P(A_n|B) = P(A_n) P(B|A_n) / \sum_{k \in \mathbb{N}} P(A_k) P(B|A_k)$$

Onafhankelijkheid

def. $A, B \in \mathcal{A}$ zijn onafhankelijk $\Leftrightarrow P(A \cap B) = P(A) \cdot P(B)$

def. paarsgewijs onafhankelijk $P(A_i \cap A_j) = P(A_i) \cdot P(A_j) \forall i \neq j$

def. onderling onafhankelijk: $\forall m \in \{2, \dots, n\}, i_1 \neq i_2 \neq \dots \neq i_m$

$$P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_m}) = P(A_{i_1}) P(A_{i_2}) \dots P(A_{i_m})$$

\hookrightarrow aantal voorwaarden: paarsgewijze: $\binom{n}{2}$

onderling: $2^n - n - 1$

Systeem betrouwbaarheid

in serie: $P(S \text{ werkt}) = P(A \text{ werkt}) \cdot P(B \text{ werkt})$

in parallel: $P(S \text{ faalt}) = P(A \text{ faalt}) \cdot P(B \text{ faalt})$

②

Stochastische verand.

Inleiding & definitie

def. **stochastische veranderlijke**: een reële functie die aan elke uitkomst $\omega \in \Omega$ van een kansruimte (Ω, \mathcal{A}, P) een reël getal $X(\omega)$ toekent met:
 $\forall B \in \mathcal{B}(\mathbb{R}) : X^{-1}(B) = \{\omega \mid X(\omega) \in B\} \in \mathcal{A}$

↳ **A-meetbare afbeelding**

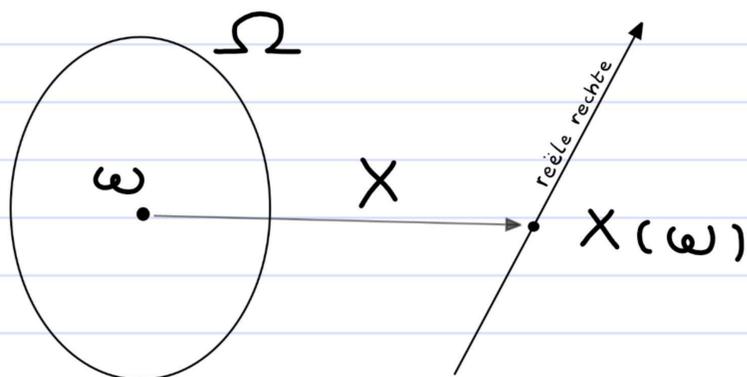
$\mathcal{B}(\mathbb{R})$ is de σ -algebra voortgebracht door $\mathcal{C} = \{]-\infty, a] \mid -\infty < a < +\infty\}$

st. Zij $X: (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$, X is een s.v. a.s.a.:

X onafh. van P

$\forall a \in \mathbb{R} : X^{-1}(]-\infty, a]) = \{\omega \mid X(\omega) \leq a\} \in \mathcal{A}$ (= def.)

↳ als $\Omega = \mathbb{R}$ en $\mathcal{A} = \mathcal{B}(\mathbb{R})$, dan is deze **Borel-meetbaar**



st. X induceert een **kansmaat** $P_x(B)$ of $P(X \in B)$ op $\mathcal{B}(\mathbb{R})$:

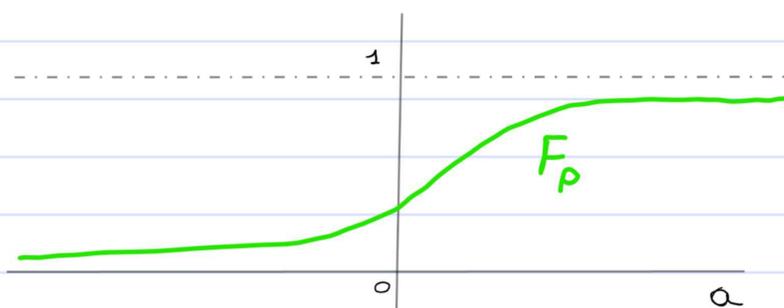
$$P_x(B) = P(X^{-1}(B)) = P(\{\omega \in \Omega \mid X(\omega) \in B\})$$

Noteer: $(\Omega, \mathcal{A}, P) \xrightarrow{X} (\mathbb{R}, \mathcal{B}(\mathbb{R}), P_x)$ kansruimte

zie bewijs p. 29

def. **verdelingsfunctie** F_x : $F_x(a) = P_x(]-\infty, a]) = P(\{\omega \mid X(\omega) \leq a\})$

$$\Rightarrow F_x(a) = P(X \leq a) \quad a \in \mathbb{R}$$



↳ monotoon stijgend:

$$\forall a \leq b : F_x(a) \leq F_x(b)$$

$$\leftarrow \lim_{a \rightarrow +\infty} F_x(a) = 1, \quad \lim_{a \rightarrow -\infty} F_x(a) = 0$$

↳ rechtscontinu:

$$\forall a \in \mathbb{R} : \lim_{h \searrow 0} F_x(a+h) = F_x(a)$$

def. **kwantiel functie** Q_x : de inverse van de verdelingsfunctie F_x . $Q_x(p)$ is de kleinste a met $F_x(a) \geq p$ met $0 < p \leq 1$. Als F_x inverteerbaar: $Q_x = F_x^{-1}$

25%, 50%, 75%
→ 1^o, 2^e en 3^e kwantiel

Types stochastische veranderlijken

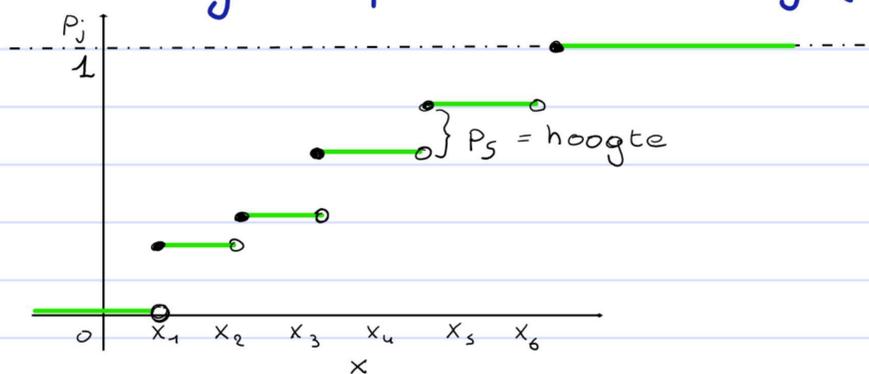
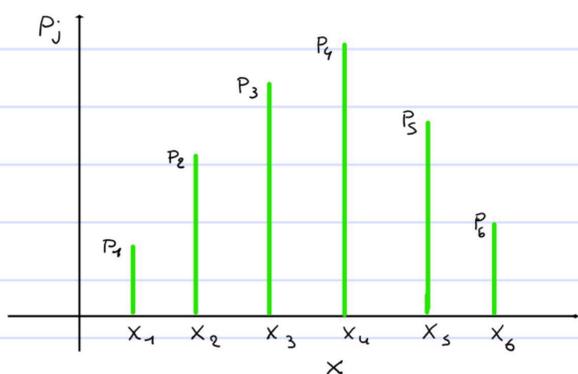
Discrete stochastische veranderlijken

↳ Ω wordt afgebeeld op een eindig of aftelbaar aantal verschillende punten x_j

$\Rightarrow p_j = P(\{\omega \in \Omega \mid X(\omega) = x_j\}) = P(X = x_j)$ met $p_j \geq 0$ en $\sum p_j = 1$
de rij $\{p_j \mid j \in \mathbb{N}\}$ is de **discrete verdeling** van de **discrete s.v. X**

↳ x_j zijn de waarden horende bij p_j

\Rightarrow een discrete verdeling bepaald een rij $\{(x_j, p_j) \mid j \in \mathbb{N}\}$



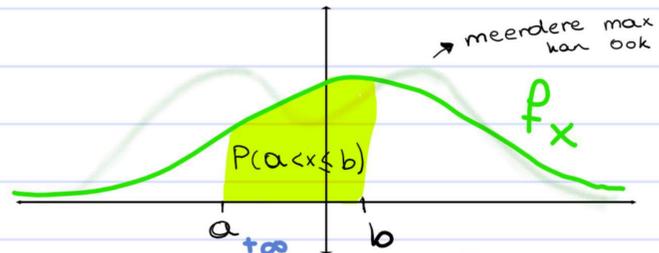
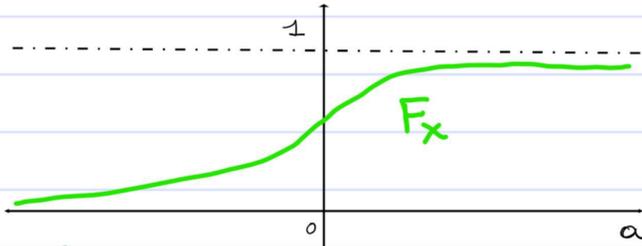
Continue stochastische veranderlijken

↳ $P_x(\{x\}) = 0$, geen sprongen in grafiek

Als F_x continu afleidbaar:

dichtheidsfunctie: $f_x(x) = \frac{dF_x(x)}{dx} = F'_x(x)$

↳ $f_x(x) \geq 0$, want F_x zwak stijgend



$$\int_{-\infty}^a f_x(x) dx = F_x(a) = P(X \leq a) \rightarrow \int_{-\infty}^{\infty} f_x(x) dx = 1$$

$$\int_a^b f_x(x) dx = F_x(b) - F_x(a) = P(a < X \leq b) = P(a \leq X < b) = P(a < X < b)$$

→ F_x heeft geen eenheden en f_x die van $\frac{1}{x}$

Opm.

Er zijn s.v. die niet discreet en niet continu zijn

Momenten van een s.v.

Verwachtingswaarde & variantie

• **verwachtingswaarde $E[X]$** (bestaat enkel als $E[X] < \infty$ volgens formule)

$$E[X] = \int_{-\infty}^{+\infty} x dF_x(x) = \begin{cases} \sum_j x_j P_j & \text{(discreet)} \\ \int_{-\infty}^{+\infty} x \cdot f_x(x) dx & \text{(continu)} \end{cases}$$

→ eigenschappen:

- $E[aX] = a E[X]$ ($a \neq 0$)
- $E[X+b] = E[X] + b$
- $E[b] = b$
- $|E[X]| = E[|X|]$

$E[X]$ is het eerste moment van X en $E[|X|]$ het eerste absolute moment.

• **variantie $Var[X]$**

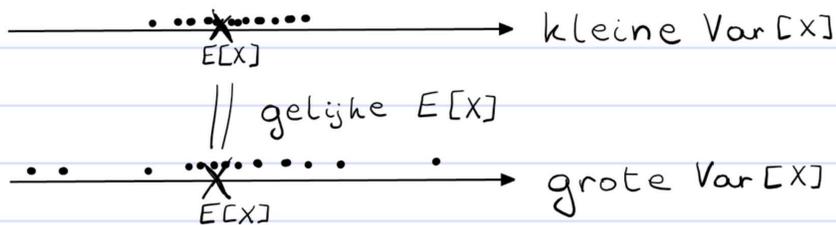
$$Var[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

$$Var[X] = \begin{cases} \sum_j (x_j - E[X])^2 P_j & \text{(discreet)} \\ \int_{-\infty}^{+\infty} (x - E[X])^2 f(x) dx & \text{(continu)} \end{cases}$$

→ eigenschappen:

- $Var[aX] = a^2 Var[X]$
- ↳ $Var[-X] = Var[X]$
- $Var[X+b] = Var[X]$
- $Var[b] = 0$

$Var[X]$ in eenheid van $X^2 \Rightarrow$
Standaardafwijking σ
 $\sigma = \sqrt{Var[X]}$



Som en product van s.v.

• $(X+Y)(\omega) = X(\omega) + Y(\omega)$

↳ $E[X+Y] = E[X] + E[Y]$

• $(XY)(\omega) = X(\omega) Y(\omega)$

→ X en Y zijn **onafhankelijk** a.s.a

$$P((X \in A) \cap (Y \in B)) = P(X \in A) \cdot P(Y \in B),$$

$$P((X \leq a) \cap (Y \leq b)) = P(X \leq a) \cdot P(Y \leq b)$$

↳ $E[XY] = E[X] E[Y]$

↳ $Var[X-Y] =$

$$Var[X+Y] = Var[X] + Var[Y]$$

$$Var[X] + Var[Y]$$

$$Var[XY] = Var[X] Var[Y] + E[X]^2 Var[Y] + Var[X] E[Y]^2 \geq Var[X] Var[Y]$$

bewijzen
 vanaf p. 44

Ongelijkheden & grenzen

• $E[X] < \infty$ als $E[|X|] < \infty$ als $\sum_{n=1}^{\infty} P(|X| \geq n)$ convergeert

bew p. 46

p. 48

• indien $E[|X|^n] < \infty$, dan $E[|X|^k] < \infty$ voor $0 \leq k \leq n$

\Rightarrow Als X positieve s.v. enkel gehele waarden kan aannemen:

$$E[X] = \sum_{n=1}^{\infty} P(X \geq n)$$

p. 49

• **ongelijkheid van Chebyshev**: X s.v., $\phi: \mathbb{R} \rightarrow \mathbb{R}$, $\phi(X)$ s.v., $E[\phi(X)] < \infty$:

① $\phi \geq 0$, even en niet-dalend voor $x \geq 0$, dan:

$$\forall a > 0: P(|X| > a) \leq \frac{1}{\phi(a)} E[\phi(X)]$$

② $\phi \geq 0$ en niet-dalend voor $-\infty < x < +\infty$, dan:

$$\forall a \in \mathbb{R}: P(X \geq a) \leq \frac{1}{\phi(a)} E[\phi(X)]$$

\Rightarrow X s.v. en $E[|X|^n] < \infty$, $n > 0$, dan:

$$\forall a > 0: P(|X| \geq a) \leq a^{-n} E[|X|^n]$$

\Rightarrow X s.v. en $E[X] = \mu$, $\text{Var}[X] = \sigma^2 < \infty$, dan:

$$\forall a > 0: P(|X - \mu| > a) \leq \frac{\sigma^2}{a^2}$$

Hogere momenten en momentgenererende functie

• **ruw moment** van orde $k > 1$: $\alpha_k(X) = E[X^k]$

• **centraal moment** van orde $k > 1$: $\mu_k(X) = E[(X - E[X])^k]$

$$\mu_1(X) = 0 \quad \mu_2(X) = \text{Var}[X]$$

• **momentgenererende functie** van X : $M_X: \mathbb{R} \rightarrow \mathbb{R}^+$, $M_X(t) > 0$

$$M_X(t) = E[e^{tX}] = \begin{cases} \sum_j e^{tx_j} p_j & \text{(discreet)} \\ \int_{-\infty}^{+\infty} e^{tx} f(x) dx & \text{(continu)} \end{cases}$$

$$e^a = 1 + a + \frac{a^2}{2!} + \frac{a^3}{3!} + \dots$$

$$e^{tX} = 1 + (tX) + \frac{(tX)^2}{2!} + \frac{(tX)^3}{3!} + \dots$$

$$E[e^{tX}] = 1 + E[X]t + E[X^2] \frac{t^2}{2!} + E[X^3] \frac{t^3}{3!} + \dots$$

$$\Rightarrow M_X(t) = 1 + \alpha_1 t + \alpha_2 \frac{t^2}{2!} + \alpha_3 \frac{t^3}{3!} + \dots$$

$$\Rightarrow \alpha_k = \frac{d^k}{dt^k} [M_X(t)]_{t=0} \quad (k \geq 1)$$

Stel $Y = aX + b$, dan $M_Y(t) = e^{bt} M_X(at)$

bew p. 52

Als X en Y onafhankelijk zijn, dan:

$$M_{X+Y}(t) = M_X(t) M_Y(t)$$

Kentallen

van Locatie

• **gemiddelde**: $\mu_X = E[X]$

• **mediaan**: $\text{Med}(X) = Q_X(0,5)$

• **modus**: meest voorkomende waarde \rightarrow piek in $f(x)$

\hookrightarrow 1 modus: unimodiaal, 2 modi: bimodaal, multimodaal

van schaal

• **variantie (en standaardafwijking)**: $\sigma^2 = \text{Var}[X] = E[X^2] - E[X]^2$

• **interkwartielafstand**: $IQR = Q_X(\frac{3}{4}) - Q_X(\frac{1}{4})$

• **median absolute deviation**: $MAD = \text{Med}|X - \text{Med}(X)|$

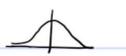
\hookrightarrow 50% van massa tussen $\text{Med}(X) - MAD(X)$ en $\text{Med}(X) + MAD(X)$

• **variatiebreedte**: $Q_X(1) - Q_X(0)$

• **variatiecoëfficiënt**: ($X \geq 0$): $v(X) = \frac{\sigma_X}{\mu_X} \rightarrow$ dimensieloos

van scheefheid

• **centrale derde moment**: $\mu_3(X) = E[(X - E[X])^3]$

\hookrightarrow 0  , > 0  , < 0 

• **scheefheidscoëfficiënt**: $\gamma_1 = \frac{\mu_3}{\sigma^3} \rightarrow$ dimensieloos

$$\frac{3(\mu - \text{Med}(X))}{\sigma} \cdot \frac{(Q(\frac{3}{4}) - Q(\frac{1}{4})) - (Q(\frac{1}{2}) - Q(\frac{1}{4}))}{Q(\frac{3}{4}) - Q(\frac{1}{4})}$$

Belangrijke verdelingen

Discrete verdelingen

- **discrete uniforme verdeling**: elke uitkomst even waarschijnlijk \rightarrow gedefinieerd op eindig universum
 $\Omega = \{x_1, \dots, x_n\}$, $P_j = P(\{x_j\}) = \frac{1}{n}$
- **Bernoulli verdeling**: slechts 2 uitkomsten, $X \rightarrow \mathcal{B}(1, p)$
 $\Omega = \{0, 1\}$ $\begin{cases} P(X=1) = p & \text{(kans op succes)} & (0 < p < 1) \\ P(X=0) = q = 1-p & \text{(kans op mislukking)} \end{cases}$
- **binomiaalverdeling**: aantal successen bij n keer herhalen van het experiment, $Y \sim \mathcal{B}(n, p)$
 $\Omega = \{0, 1, \dots, n\}$, $P(Y=k) = \binom{n}{k} p^k q^{n-k} = \frac{n!}{k!(n-k)!} p^k q^{n-k}$
 \forall als $X_i \sim \mathcal{B}(1, p)$, $i=1, \dots, n$ onderling onafhankelijk ($Y = X_1 + \dots + X_n$)
- **geometrische verdeling**: aantal mislukkingen tot eerste succes bij Bernoulli experiment. $X \sim \text{Geom}(p)$
 $\Omega = \mathbb{N}$, $P(X=j) = q^j p = q^j (1-q)$
- **negatief binomiaalverdeling**: aantal mislukkingen tot r -de succes bij Bernoulli experiment. $X \sim \text{NB}(r, q)$
 $\Omega = \mathbb{N}$, $P(X=j) = \binom{j+r-1}{j} (1-q)^r q^j = \binom{-r}{j} p^r (-q)^j$
- **hypergeometrische verdeling**: s goed uit N opties, n trekken zonder terugleggen \rightarrow # goede. $Y \sim \mathcal{H}(N, s, n)$
 $\Omega = \{0, 1, \dots, n\}$, $\max(0, n-(N-s)) \leq j \leq \min(s, n)$, $P(Y=j) = \frac{\binom{j}{s} \binom{N-j}{n-s}}{\binom{N}{n}}$
 \hookrightarrow met teruglegging: $Y \sim \mathcal{B}(n, \frac{s}{N})$
- **poissonverdeling**: aantal gebeurtenissen met rate α , $0 < \alpha < \infty$, $\mathcal{P}(\alpha)$, $\Omega = \mathbb{N}$, $P(X=j) = \frac{\alpha^j}{j!} \cdot e^{-\alpha}$
 \hookrightarrow aantal gebeurtenissen in disjuncte intervallen is onafhankelijk van elkaar
 $\hookrightarrow X_1 \sim \mathcal{P}(\alpha_1)$ en $X_2 \sim \mathcal{P}(\alpha_2)$ onafhankelijk, dan:
 $X_1 + X_2 \sim \mathcal{P}(\alpha_1 + \alpha_2)$

Continue verdelingen

- **continu uniform**: verdeling op $[a, b]$, $-\infty < a < b < +\infty$,
 $X \sim \mathcal{U}[a, b]$, $f_{a,b}(x) = \frac{1}{b-a} \mathbb{1}_{[a,b]}(x)$. $\Omega = [a, b]$
- **exponentieel**: Levensduur (continue versie van geometrische)
 $X \sim \mathcal{E}(\alpha)$, $0 < \alpha < \infty$, $\Omega = [0, +\infty[$, $f_\alpha(t) = \begin{cases} \alpha \cdot e^{-\alpha t} & t \geq 0 \\ 0 & t < 0 \end{cases}$
 \hookrightarrow geheugenloos: $P(X > s+t | X > t) = P(X > s)$
- **univariate normale verdeling**:
 - * **standaard normaal**: $Z \sim \mathcal{N}(0, 1)$, $\phi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \forall x \in \mathbb{R}$
met ϕ de dichtheidsfunctie, ϕ is even, verdeling:
 $\Phi(x) = \int_{-\infty}^x \phi(t) dt \leftarrow$ tabellen, $P(-z \leq Z \leq z) = 2\phi(z) - 1$
 - * **algemeen normaal**: $X \sim \mathcal{N}(\mu, \sigma^2)$, $X = \sigma Z + \mu$, $Z \sim \mathcal{N}(0, 1)$
 $f_{\mu, \sigma}(x) = \frac{1}{\sigma} f_Z\left(\frac{x-\mu}{\sigma}\right) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
 $F_X(x) = F_Z\left(\frac{x-\mu}{\sigma}\right) = \Phi\left(\frac{x-\mu}{\sigma}\right)$
 $X_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$ en $X_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$ onafhankelijk:
 $X_1 + X_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$
- χ_n^2 : Z_1, \dots, Z_n , $Z_i \sim \mathcal{N}(0, 1)$. $X = Z_1^2 + \dots + Z_n^2$, chi-kwadraat verdeling met n vrijheidsgraden $X \sim \chi_n^2$
 $f_{\chi_n^2}(x) = \begin{cases} \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} e^{-\frac{x}{2}} x^{\frac{n}{2}-2} & (x > 0) \\ 0 & (x \leq 0) \end{cases}$ $\Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx$
 $\hookrightarrow \forall t > 0: \Gamma(t+1) = t \Gamma(t)$
 $\hookrightarrow \forall n \in \mathbb{N}: \Gamma(n+1) = n!$
 $\rightarrow \Gamma(\frac{1}{2}) = \sqrt{\pi}$

Formule van Stirling: $\Gamma(t) \sim \left(\frac{t-1}{e}\right)^{t-1} \sqrt{2\pi(t-1)}$ *
 $\lim_{n \rightarrow \infty} \frac{\Gamma(n)}{n!} = 1$

$\Rightarrow t = n+1 \rightarrow n! \sim \left(\frac{n}{e}\right)^n \sqrt{2\pi n}$

- **gammaverdeling**: modelleren van wachttijden, $X \sim \Gamma_{\gamma, \beta}$, $\gamma, \beta > 0$, $f_{\gamma, \beta}(x) = \frac{x^{\gamma-1} e^{-x/\beta}}{\beta^\gamma \Gamma(\gamma)} \mathbb{1}_{\{x>0\}}$
- t_n : $Z \sim \mathcal{N}(0, 1)$, $X \sim \chi_n^2$ ($n > 1$) onafhankelijk, $T \sim t_n$ een studentverdeling met n vrijheidsgraden $T = \frac{Z}{\sqrt{X/n}}$
 $\hookrightarrow \text{Med}(T) = 0$, $E[T] = 0$, even, $\text{Var}[T] = \frac{n}{n-2}$ ($n > 2$)
- **Cauchy**: t_n voor $n=1$, als $X \sim \mathcal{N}(0, 1)$ en $Y \sim \mathcal{N}(0, 1)$ onafhankelijk, dan $\frac{X}{Y}$ Cauchy verdeeld
 $f(x) = \frac{1}{\pi} \frac{1}{1+x^2}$ ← zwaardere staart dan ϕ
 $F(x) = \frac{1}{\pi} \text{bgtan}(x) + \frac{1}{2}$
- **F**: als $W \sim \chi_m^2$ en $V \sim \chi_n^2$ onafhankelijk, $X = \frac{W/m}{V/n}$ de F-verdeling met m vrijheidsgraden in de teller en n in de noemer, $X \sim F_{m,n}$, α_k bestaat als $k < \frac{n}{2}$,
 $F_{m,n}(x) = 1 - F_{n,m}(x)$
- **Lognormaal**: $Y > 0$, $\mu \in \mathbb{R}$, $\sigma > 0$, $X = \ln(Y) \sim \mathcal{N}(\mu, \sigma^2)$
 $f_Y(y) = \phi\left(\frac{\ln(y) - \mu}{\sigma}\right)$

Transformatie van s.v.

Monotone transformatie

Zij X een continue s.v. met $f_X(x)$, $x \in S$: $f_X(x) = 0$, neem $h: \mathbb{R} \rightarrow \mathbb{R}$, zodat $U = h(X)$ een s.v. Als h differentieerbaar en strikt stijgend/dalend op S , dan:

bew p. 72

$$f_U(u) = \begin{cases} f_X(h^{-1}(u)) \left| \frac{dh^{-1}(u)}{du} \right| & u \in h(S) \\ 0 & u \notin h(S) \end{cases}$$

Integraaltransformatie

bew p. 74

Zij X een continue s.v. met F_X strikt stijgend op $F_X^{-1}(]0, 1[)$, stel $U = F_X(X)$, dan geldt $U \sim \mathcal{U}[0, 1]$

Genereren van s.v.

bew p. 74

Zij F strikt stijgend en $U \sim \mathcal{U}[0, 1]$, $X = F^{-1}(U)$ met F de verdelingsfunctie van X .

3

Bivariate verd.

Verdeling van een stochastisch koppel

Stochastisch koppel $(X, Y) \leftarrow$ stochastische vector van lengte 2
 \hookrightarrow versch. verdelingen.

Gezamenlijke verdeling

\hookrightarrow X en Y samen beschouwen

def. Zij X en Y s.v. op (Ω, \mathcal{A}, P) , dan: $(X, Y): \Omega \rightarrow \mathbb{R}^2: \omega \mapsto (X(\omega), Y(\omega))$
met (X, Y) de stochastische vector $(\Omega, \mathcal{A}, P) \xrightarrow{(X, Y)} (\mathbb{R}^2, \mathcal{B}(\mathbb{R}^2), P_{X, Y})$

$\mathcal{B}(\mathbb{R}^2) \rightarrow]-\infty, a_1] \times]-\infty, a_2]$ $a_1, a_2 \in \mathbb{R}$

$P_{X, Y}(]-\infty, x] \times]-\infty, y]) = P(X^{-1}(]-\infty, x]) \cap Y^{-1}(]-\infty, y]))$

\Rightarrow gezamenlijke verdelingsfunctie $F_{X, Y}$ van X en Y:

$$F_{X, Y}(x, y) = P(X \leq x, Y \leq y) = P_{X, Y}(]-\infty, x] \times]-\infty, y])$$

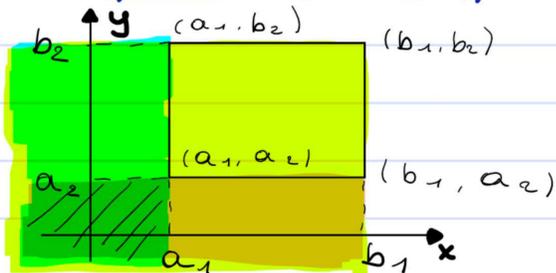
- stijgend in alle argumenten
- rechtscontinu in alle argumenten
- $F_{X, Y}(x, y) \rightarrow 0$ als $x \rightarrow -\infty$ en/of $y \rightarrow -\infty$
- $F_{X, Y}(x, y) \rightarrow 1$ als $x \rightarrow +\infty$ en $y \rightarrow +\infty$

\uparrow elke bivariate F die hieraan voldoet is een verdelingsf.

$F_{X, Y}$ bepaald volledig de verdeling van (X, Y) :

$$P((X, Y) \in]a_1, b_1] \times]a_2, b_2])$$

$$= F_{X, Y}(b_1, b_2) - F_{X, Y}(b_1, a_2) - F_{X, Y}(a_1, b_2) + F_{X, Y}(a_1, a_2)$$



\rightarrow continue bivariate verd.:

$$F_{X, Y}(x, y) = \int_{-\infty}^x \int_{-\infty}^y f_{X, Y}(u, v) du dv$$

met $f_{X, Y}(x, y)$ de gezamenlijke dichtheidsfunctie

- $f_{X, Y}(x, y) \geq 0 \forall x, y \in \mathbb{R}$
 - $\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{X, Y}(u, v) du dv = 1$
- $$f_{X, Y}(x, y) = \frac{\partial^2 F_{X, Y}(x, y)}{\partial x \partial y}$$

Marginale verdeling

\hookrightarrow verdeling van X of verdeling van Y

$$F_X(x) = P(X \leq x) = P(X \leq x, Y \leq +\infty) = F_{X, Y}(x, +\infty) \text{ (idem voor } F_Y(y))$$

\rightarrow marginale dichtheidsfunctie van X: $f_X(x) = \int_{-\infty}^{+\infty} f_{X, Y}(x, y) dy$

Onafhankelijkheid

def. De s.v. X en Y zijn onafhankelijk als voor alle $A, B \in \mathcal{B}(\mathbb{R})$ $\{X \in A\}$ en $\{Y \in B\}$ (onderling) onafhankelijk zijn.

$$\hookrightarrow F_{X, Y}(x, y) = F_X(x) F_Y(y)$$

$\Rightarrow E[XY] = E[X]E[Y], \text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$ + zie s.v.

Voorwaardelijke verdeling

\hookrightarrow verdeling van X, gegeven Y

$$P(X=x | Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)} \text{ als } P(Y=y) > 0$$

$$\text{voorwaardelijke dichtheidsfunctie } f_{X|Y}(x|y) = \frac{f_{X, Y}(x, y)}{f_Y(y)}$$

- $f_{X|Y}(x|y) \geq 0 \forall x$

$$\int_{-\infty}^{+\infty} f_{X|Y}(x|y) dx = \frac{\int_{-\infty}^{+\infty} f_{X, Y}(x, y) dx}{f_Y(y)} = \frac{f_Y(y)}{f_Y(y)} = 1$$

discreet
continu

2 s.v. zijn onafh. a.s.a. $P(X=x | Y=y) = P(X=x) \forall x, y$ + omg
 $f_{X|Y}(x|y) = f_X(x)$

Karakteristieken

Momenten & MGF

def. Zij (X, Y) stochastisch koppel en $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ Borelmeetbaar dan $E[g(X, Y)] = \begin{cases} \sum_x \sum_y g(x, y) P(X=x, Y=y) & \text{als } < \infty \text{ (discreet)} \\ \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x, y) f_{X, Y}(x, y) dx dy & \text{als } < \infty \text{ (continu)} \end{cases}$

- (ruwe) momenten: $g(x, y) = x^{r_1} y^{r_2}$, $E[X^{r_1} Y^{r_2}]$
- centrale momenten: $g(x, y) = (x - \mu_x)^{r_1} (y - \mu_y)^{r_2}$ met $\mu_i = E[i]$
 $E[(X - \mu_x)^{r_1} (Y - \mu_y)^{r_2}]$
- gezamenlijke momentgenererende functie

Covariantie en correlatie

def. $\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$
 $\hookrightarrow \text{Cov}(X, X) = E[(X - E[X])^2] = \text{Var}[X]$
 $\text{Cov}(X, Y) = \text{Cov}(Y, X)$

zie p. 86 $\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$

def. **correlatiecoëfficiënt** $\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}[X]\text{Var}[Y]}}$

\rightarrow beschrijven lineair verb. tussen X en Y

cov: mate van verband (pos / neg \)

corr: maat voor sterkte verband $\in [-1, 1] \rightarrow \pm 1 \rightarrow$ perfect

Eigenschappen

st. Zij X en Y onafhankelijk en $g: \mathbb{R}^2 \rightarrow \mathbb{R}$ Borelmeetbaar, zodat

bew. p. 87 $\forall x, y \in \mathbb{R}: g(x, y) = g_1(x) g_2(y)$ en $E[g_1(x)], E[g_2(y)] < \infty$, dan:
 $E[g(x, y)] = E[g_1(x)] E[g_2(y)]$

$\Rightarrow M_{X, Y}(t_1, t_2) = M_X(t_1) M_Y(t_2)$

st. Zij X, Y s.v., dan:

bew. p. 87 $\bullet | \text{Cov}(X, Y) | \leq \sqrt{\text{Var}[X] \text{Var}[Y]}$

\bullet Zij X, Y onafh., dan $E[XY] = E[X]E[Y]$ en $\text{Cov}(X, Y) = 0$
 \hookrightarrow omgekeerd niet altijd waar

\Rightarrow Voor X, Y s.v. geldt:

bew. p. 88 $\bullet -1 \leq \text{Corr}(X, Y) \leq 1$

\bullet Zij $X_2 = aX_1 + b$ met $a, b \in \mathbb{R}$, dan $\text{Corr}(X, Y) = \frac{a}{|a|} = \text{sgn}(a)$

st. Zij X en Y s.v. en $a_1, a_2 \in \mathbb{R}$, dan:

bew. p. 89 $\bullet E[a_1 X + a_2 Y] = a_1 E[X] + a_2 E[Y]$

$\bullet \text{Var}[a_1 X + a_2 Y] = a_1^2 \text{Var}[X] + a_2^2 \text{Var}[Y] + 2a_1 a_2 \text{Cov}(X, Y)$

\Rightarrow zij X en Y onafh. en $a_1, a_2 \in \mathbb{R}$, dan:

$\text{Var}[a_1 X + a_2 Y] = a_1^2 \text{Var}[X] + a_2^2 \text{Var}[Y]$

Bivariate normale verdeling

def. (X, Y) heeft een **bivariate normale verdeling** als de gezamenlijke verdelingsfunctie wordt gegeven door:

$$f_{X, Y}(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x-\mu_1}{\sigma_1}\right)^2 - 2\rho \left(\frac{x-\mu_1}{\sigma_1}\right) \left(\frac{y-\mu_2}{\sigma_2}\right) + \left(\frac{y-\mu_2}{\sigma_2}\right)^2 \right]\right)$$

$\forall (x, y) \in \mathbb{R}^2$ met $\mu_1, \mu_2 \in \mathbb{R}$, $\sigma_1, \sigma_2 \in \mathbb{R}_0^+$ en $\rho \in]-1, 1[$.

Hierin is: $X \sim \mathcal{N}(\mu_1, \sigma_1^2)$ met $\mu_1 = E[X]$ en $\sigma_1^2 = \text{Var}[X]$

en $Y \sim \mathcal{N}(\mu_2, \sigma_2^2)$ met " en $\rho = \text{Cor}(X, Y)$

$\Rightarrow \text{Cov}(X, Y) = \rho \sigma_1 \sigma_2$

\rightarrow **standaard bivariate normale verdeling**

$\hookrightarrow \mu_1 = \mu_2 = 0, \sigma_1 = \sigma_2 = 1, \rho = 0$

$\Rightarrow f_{X, Y}(x, y) = \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x^2 + y^2)\right)$

4

Benaderingen

Limietstellingen

def. Zij X_1, X_2, \dots s.v. met verdelingsfuncties $F_1(x), F_2(x), \dots$
 Dan convergeert X_n in verdeling naar X als er een $F(x)$ bestaat met $\lim_{n \rightarrow \infty} F_n(x) = F(x) \quad \forall x$ waar F continu is
 ↳ notatie: $F_n \xrightarrow{D} F$ of $X_n \xrightarrow{D} X$
 † puntsgewijze conv. van F_n naar F in continuïteitspunten van F
 • steekproefgemiddelde: X_1, \dots, X_n i.i.d., dan $\bar{X} = \frac{X_1 + \dots + X_n}{n}$
 Dan is: • $E[\bar{X}] = E[X_1]$
 • $\text{Var}[\bar{X}] = \frac{1}{n} \text{Var}[X_1]$
 $\Rightarrow X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ en onafh. $\Rightarrow \bar{X} \sim (\mu, \frac{\sigma^2}{n})$

Centrale Limietstelling

st. Zij (X_n) een rij i.i.d. s.v. met $E[X_1] = \mu$ en $\text{Var}[X_n] = \sigma^2 < \infty$.
 Stel $\bar{X} = (X_1 + \dots + X_n) / n$, dan geldt:
 $\forall x \in \mathbb{R}: \lim_{n \rightarrow \infty} P\left(\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \leq x\right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}u^2} du$ ofwel
 $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \xrightarrow{D} Z$ met $Z \sim \mathcal{N}(0, 1)$

- $E[\bar{X}] = \frac{1}{n} E[X_1 + \dots + X_n] = \frac{1}{n} \cdot n \cdot E[X_1] = E[X_1]$
- $\text{Var}[\bar{X}] = \text{Var}[X_1 + \dots + X_n] / n^2 = \frac{n}{n^2} \text{Var}[X_1] = \text{Var}[X_1] / n$

Herschrijving CLS:

$$S_n = n\bar{X} = X_1 + \dots + X_n \Rightarrow E[S_n] = n\mu, \text{Var}[S_n] = n\sigma^2$$

$$\text{Voor } n \rightarrow \infty: P\left[\left(\frac{S_n - n\mu}{\sqrt{n}\sigma}\right) \leq x\right] \rightarrow \Phi(x)$$

→ in praktijk: Zij (X_n) i.i.d. en n voldoende groot ($\mu, \sigma < \infty$):

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \approx \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

$$S_n = \sum_{i=1}^n X_i \approx \mathcal{N}(n\mu, n\sigma^2)$$

Opm.:

- benadering
- verd. v. X_i moet niet continu zijn
- X_i moet niet symmetrisch zijn (hoe schever, hoe groter n)

Stelling van De Moivre-Laplace

st. Bernoulli verd.: $S_n = X_1 + \dots + X_n$ met $X_i \sim \mathcal{B}(1, p)$, dan $S_n \sim \mathcal{B}(n, p)$

$$\hookrightarrow E[X_1] = \mu = p, \text{Var}[X_1] = \sigma^2 = p(1-p)$$

$$E[S_n] = np, \text{Var}[S_n] = np(1-p)$$

$$\frac{S_n - np}{\sqrt{npq}} \xrightarrow{D} Z \text{ met } Z \sim \mathcal{N}(0, 1) \Rightarrow \mathcal{B}(n, p) \approx \mathcal{N}(np, np(1-p))$$

$$\hat{=} n \geq 30, np > 5 \text{ en } n(1-p) > 5$$

Continuïteitscorrectie

$$\left. \begin{aligned} P(Y \geq a) &\approx 1 - \Phi\left(\frac{a - 0,5 - np}{\sqrt{npq}}\right) \\ P(Y \leq b) &\approx \Phi\left(\frac{b + 0,5 - np}{\sqrt{npq}}\right) \end{aligned} \right\} P(a \leq Y \leq b) \approx \Phi\left(\frac{b + 0,5 - np}{\sqrt{npq}}\right) - \Phi\left(\frac{a - 0,5 - np}{\sqrt{npq}}\right)$$

Limietstelling van Poisson

st. Zij (X_n) een rij s.v. met $X_n \sim \mathcal{B}(n, p_n)$ met als $n \rightarrow \infty$, dan $p_n \rightarrow 0$,

zodat $np_n \rightarrow \alpha$ met $0 < \alpha < \infty$, dan $X_n \xrightarrow{D} X$ met $X \sim \mathcal{P}(\alpha)$

$$\Rightarrow \mathcal{B}(n, p) \approx \mathcal{P}(np)$$

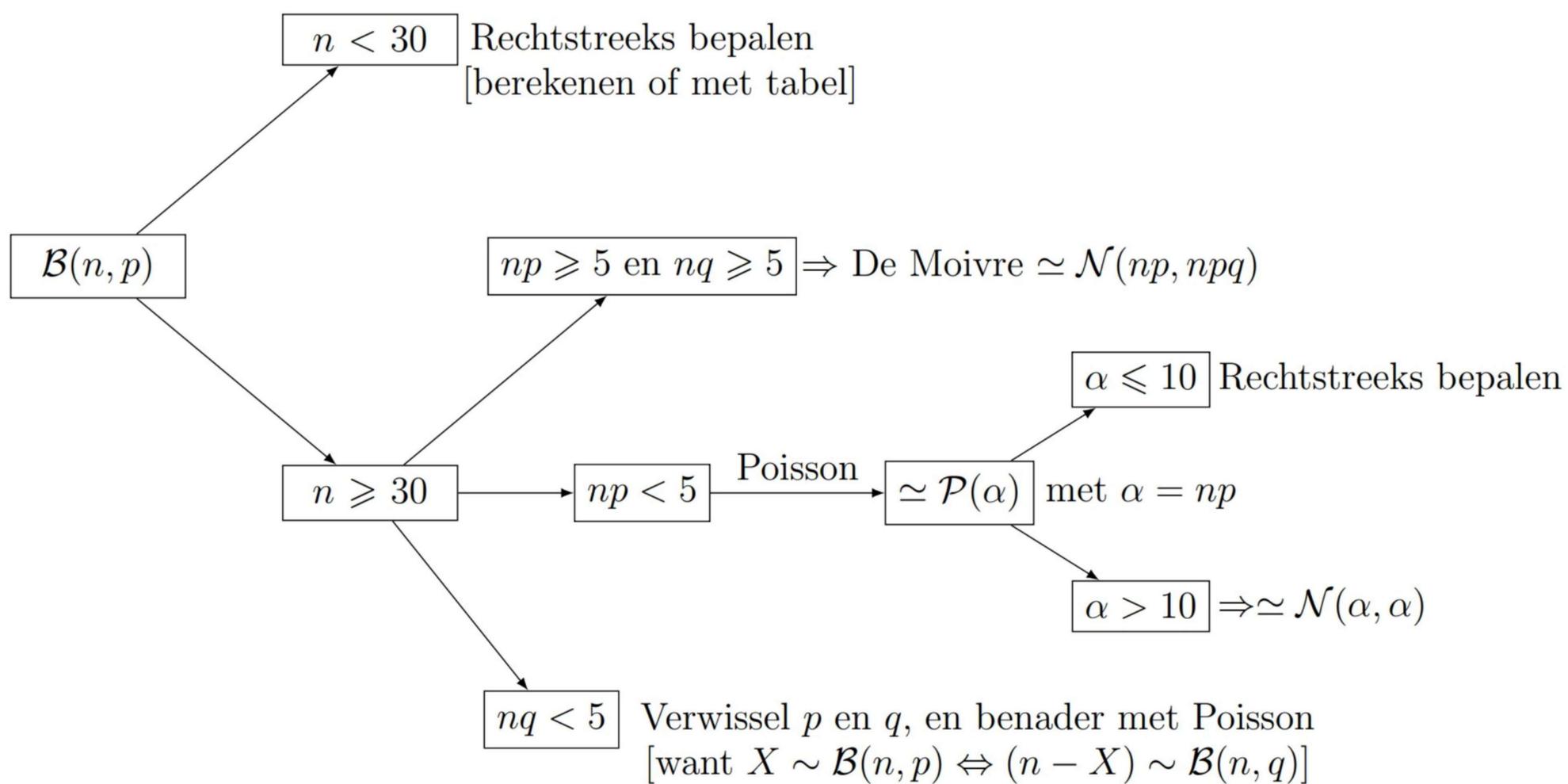
↳ $n \geq 30$ en $np < 5$ of $nq < 5$ → p en q van plaats wisselen

→ Normale benadering voor Poissonverdeling:

st. Zij $X \sim \mathcal{P}(\alpha)$, dan geldt voor $\alpha \rightarrow \infty$ dat $\frac{X - \alpha}{\sqrt{\alpha}} \xrightarrow{D} \mathcal{N}(0, 1)$

bew slides $\Rightarrow \mathcal{P}(\alpha) \approx \mathcal{N}(\alpha, \alpha)$

$$\hat{=} \alpha > 10$$



①

Beschrijvende statistiek

Beschrijven van frequenties

3 stappen bij statistische analyse:

- 1) Verzamelen van gegevens → steekproef van n observaties
- 2) beschrijvende statistiek → waar dit hoofdstuk over gaat
- 3) statistische inferentie → beslissingen nemen

Wat willen we weten?

↳ hypothesetoetsen

- grote verschillen tussen geg. of \pm gelijk?
- hoe groot zijn de verschillen?
- uitschieters?
- trends?

⇒ kansmodel bepalen

Soorten variabelen

- **kwantitatieve var.**: geven categorie aan
 - ↳ **nominale var.**: geen rangorde (nationaliteit)
 - ↳ **ordinale var.**: kunnen geordend \bar{w} (weinig → veel)
- **kwantitatieve var.**: numerieke waarden (lengte)
 - ↳ continu of discreet
 - ↳ bewerkingen zinvol (meet eenheid)

Frequentietabellen en grafieken

steekproef: (x_1, \dots, x_n) , uitkomstenverz.: $\Omega = \{m_1, \dots, m_n\}$

- **absolute frequentie**: $n_j = \#$ keer dat m_j voorkomt
 - **relatieve frequentie**: $h_j = n_j/n$
- $\sum_j n_j = n$ $\sum_j h_j = 1$

→ in frequentietabel (relatief) staafjesdiagram of taartdiagram zetten

- **histogram**: voor oneindig universum v. kwantitatieve
 - ↳ geg. **discretiseren** / groeperen in **klassen**
 - ↳ **klassenmiddens** c_j , **klassebreedte** Δ
 - ↳ continue tegenhanger v.h. staafjesdiagram
 - ↳ frequenties \propto opp.

→ **klassiek histogram** bij vaste Δ

→ klassen met \neq breedtes Δ_j \neq vertekend beeld

↳ schalen: hoogte = $\frac{h_j}{\Delta_j}$ bij **dichtheidshistogram**

totale opp: $\sum_j \Delta_j \frac{h_j}{\Delta_j} = 1$

Kernel dichtheid schatter

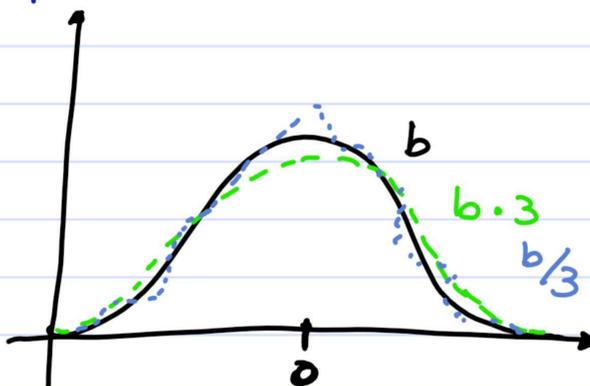
↳ (dichtheids) histogram is ruwe schatter v. dichth. en hangt af v. aantal (en locatie) v. balken

↳ Kernel nauwkeuriger: **kernel functie** K

↳ dichtheidsfunctie symm. rond 0

$\hat{f}(x) = \frac{1}{nb} \sum_{i=1}^n K\left(\frac{x_i - x}{b}\right)$ met b de bandbreedte

→ optimale bandbreedte kiezen (niet stabiel voor kleine b)



→ weerspiegeling v. verdeling van beschikbare steekproefgeg.

→ correcte vorm v. verd. identificeren

ecdf

Empirische verdelingsfunctie en kwantiel functie

$$\hat{F}_n(x) = \frac{1}{n} (\# x_i \leq x; i=1, \dots, n) = \frac{\sum_{i=1}^n \mathbb{1}_{x_i \leq x}}{n} \quad \mathbb{1}_{x_i \leq x} = \begin{cases} 1 & \text{als } x_i \leq x \\ 0 & \text{als } x_i > x \end{cases}$$

- trapfunctie
- geen samenvallende waarden $\Rightarrow \hat{F}_n(x_i) = \frac{i}{n}$
- enkel gegroepede geg. \Rightarrow **cumulatieve freq.** voor klassengrenzen met n_j en h_j
- **empirische kwantiel functie** \hat{Q}_n : \hat{F}_n trap, dus \hat{Q}_n (inverse van \hat{F}_n) als volgt bepalen:
 - geen x met $\hat{F}_n(x) = p$, dan \hat{Q}_n kleinste x met $\hat{F}_n(x) \geq p$
 - interval met $\hat{F}_n(x) = p$, $\hat{Q}_n(p) =$ kleinste x in interval $\Rightarrow \hat{Q}_n(p)$ is kleinste x met $\hat{F}_n(x) \geq p$
 - $\hat{Q}_n(i/n) = x_i$
 - $\hat{Q}_n(p) = x_i$ als $\frac{i-1}{n} < p \leq \frac{i}{n}$
 - $\hat{F}(x_i) = \frac{i}{n+1}$ of $\frac{i-0,5}{n}$ (zie slides)

Kenmerken voor Locatie en schaal

gegroepeerde gegevens

- ### Centrumkenmerken
- **steekproefgemiddelde**: $\bar{x} = \frac{1}{n} (x_1 + \dots + x_n) = \frac{1}{n} \sum x_i = \frac{1}{n} \sum n_j c_j$
 - **steekproefmediaan**: $\text{med}(x) = \begin{cases} x_{(n+1)/2} & (n \text{ oneven}) \\ (x_{n/2} + x_{n/2+1}) / 2 & (n \text{ even}) \end{cases}$
 ↳ benadering: $\frac{y-y_0}{x-x_0} = \frac{y_1-y_0}{x_1-x_0}$ $x=50\%$ $y=\text{med}$
 - **modus**: meest voorkomende waarde
 ↳ **modale klasse**: klasse met grootste verh. $\frac{\text{frequentie}}{\text{klassebreedte}}$

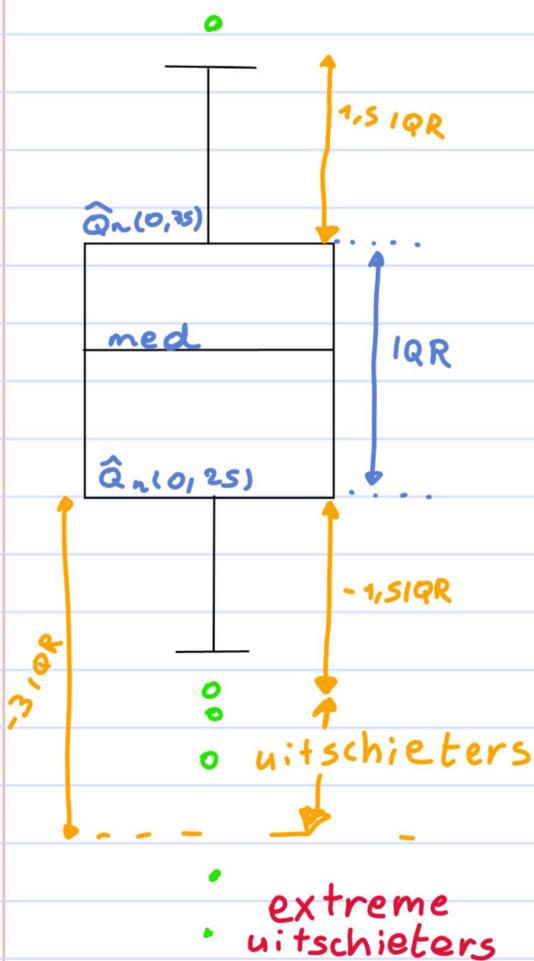
Spreadingskenmerken

- **empirische var.** s^2 en standaardafwijking: $s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ en $s = \sqrt{s^2}$
- **empirische interkwartielafstand**: $IQR_n = \hat{Q}_n(3/4) - \hat{Q}_n(1/4)$ → als n groot genoeg: $\frac{IQR_n}{s} \approx 1,34$
- **median absolute deviation**: $\text{Mad}(x) = \text{med}(|x - \text{med}(x)|)$
 → als n voldoende groot: $\frac{\text{Mad}(x)}{s} \approx 0,67$
- **bereik (range)**: $R = x_n - x_1$

Boxplot

zichtbare informatie over verdeling:

- centrum: mediaan (soms ook gem)
- IQR = lengte v.d. doos
- scheefheid: symmetrie & whiskers
- zwaarte v.d. staarten



Covariantie en correlatie

↳ Lineair verband tussen twee (metrische) var. X en Y?

→ grafisch: via scatterplot

→ analytisch: via steekproef covariantie en -correlatie

• **steekproef covariantie:**

$$s_{xy} = \text{Cov}(x, y) = \frac{1}{n-1} \sum_i (x_i - \bar{x})(y_i - \bar{y}) \quad , \quad s_{xx} = s_x^2$$

→ Lineariteitseigenschap: als $u_i = a_1 x_i + b_1$ en $v_i = a_2 y_i + b_2$, dan

$$s_{uv} = a_1 a_2 s_{xy}$$

↳ sterk afhankelijk van meeteenheid

• (Pearson) **correlatiecoëfficiënt:**

$$r = r_{xy} = \text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{s_x s_y}$$

$$\text{als } s_x > 0 \text{ en } s_y > 0: r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}} = \frac{1}{n-1} \sum_i \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

$$\rightarrow -1 \leq r \leq 1$$

→ Lineariteitseigenschap: $|r_{uv}| = |r_{xy}|$

→ maat voor lineair verb. $|r| \uparrow$ sterker verb. $r > 0$ / $r < 0$ \

Kruistabellen

↳ geeft verband tussen twee kwalitatieve variabelen

→ verwachte waarde voor elk onderdeel berekenen

en vergelijken. ong. gelijk \Rightarrow onafhankelijk

$$\text{verwachte waarde} = \frac{\text{rijtotaal} \cdot \text{kolomtotaal}}{n}$$

QQ-plot (kwantielplot)

↳ normaliteit (of andere continue verdeling) nagaan

→ **Normale QQ-plot:**

- zet empirische kwantielen $\hat{Q}_n(p)$ uit tegenover theoretische kwantielen $Q(p) = \Phi^{-1}(p)$ v.e. standaardnormale verd voor versch. $0 \leq p \leq 1$

↳ Lineair verb. bij normaal verd. geg.

$$Q(p) = \mu + \sigma \Phi^{-1}(p) \quad \forall 0 < p < 1 \rightarrow (\Phi^{-1}(p), Q(p)) \text{ vormen rechte}$$

→ bij lin verb.: $(\Phi^{-1}(p), \hat{Q}_n(p)) \approx$ rechte

welke p gebr.? $\rightarrow \hat{Q}_n(i/n) = x_i \Rightarrow p = 1/n, \dots, (n-1)/n, 1$

→ continuïteitscorrectie: $(\Phi^{-1}(\frac{i-0.5}{n}), x_i)$

⚠ let op welke as waar staat ⚠

→ exponentiële verdeling

$$\bullet 1 - F_\alpha(x) = e^{-\alpha x} \quad (x > 0)$$

$$\bullet \text{standaard verd.: } 1 - F_1(x) = e^{-x}$$

$$\bullet Q_\alpha(p) = -\frac{1}{\alpha} \log(1-p) \quad (0 < p < 1) \quad \bullet Q_1(p) = -\log(1-p)$$

$$\bullet \text{plot } (-\log(1-p), \hat{Q}_n(p)) \rightarrow (-\log(1 - \frac{i-0.5}{n}), x_i)$$

• rechtlijnig patroon met rico $\frac{1}{\alpha}$

algemeen:

$$\bullet \text{lognormale verd.f.: } F_x(x) = \Phi\left(\frac{\ln(x) - \mu}{\sigma}\right) \Rightarrow \ln(x_i) = \mu + \sigma \Phi^{-1}\left(\frac{i}{n}\right)$$

②

Schatten v. parameters & bi

Puntschatters

\bar{X} : schatter (s.v.), \bar{x} : schatting (van 1 steekproef)

Ξ : uitkomstenverzameling / waarnemingsruimte

\mathcal{B} : bijbehorende σ -algebra

θ : parameter, Θ : parameter ruimte

$\{P_\theta; \theta \in \Theta\}$ verz. kansverdelingen op meetbare ruimte (Ξ, \mathcal{B})

\rightarrow 1 kansverdeling per $\theta \in \Theta$

def. $(\Xi, \mathcal{B}, \{P_\theta; \theta \in \Theta\})$ is een **statistisch model** a.s.a.

$\forall \theta \in \Theta: (\Xi, \mathcal{B}, P_\theta)$ is een kansruimte

def. Een reële (of univariate) **statistiek** T is een meetbare afbeelding $T: (\Xi^n, \mathcal{B}^{\otimes n}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$

$\rightarrow T$ is een functie v.d. steekproef en niet v.d. onbekende parameter θ

$\Rightarrow T$ is een s.v. en heeft dus een verdeling

def. Beschouw een statistisch model $(\Xi, \mathcal{B}, \{P_\theta; \theta \in \Theta\})$ en een willekeurige afbeelding $g: \Theta \rightarrow \mathbb{R}$ (estimand).

Een **(punt)schatter** (estimator) van $g(\theta)$ is een reële statistiek $T: (\Xi^n, \mathcal{B}^{\otimes n}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ (onafh. van θ).

Een realisatie $T(x)$ voor een $x \in \Xi^n$ heet een

(punt)schatting (estimate) van $g(\theta)$.

zie p. 28-29 voor voorbeeld

Onvertekende schatters

def. **Bias** (vertekening) van een schatter is

$$b_\theta(\hat{\theta}) = E_\theta[\text{error}] = E_\theta[\hat{\theta} - \theta] = E_\theta[\hat{\theta}] - \theta$$

Een schatter $\hat{\theta}$ voor θ is **onvertekend** of **zuiver** als

$$E_\theta[\hat{\theta}] = \theta \quad (b_\theta(\hat{\theta}) = 0)$$

\rightarrow **mean squared error** van een schatter $\hat{\theta}$:

bew. p. 29

$$MSE_\theta(\hat{\theta}) = E_\theta[(\hat{\theta} - \theta)^2] = b_\theta(\hat{\theta})^2 + \text{Var}_\theta[\hat{\theta}] \quad \text{Var}_\theta[\hat{\theta}] = E_\theta[(\hat{\theta} - E_\theta[\hat{\theta}])^2]$$

\rightarrow zo klein mogelijk voor nauwkeurige $\hat{\theta}$

zie slide 34 voor rekenvoorbeeld (in. comb.)

Parameters v. normale populatie schatten

$x_i \sim \mathcal{N}(\mu, \sigma^2)$, $\hat{\mu} = \bar{x}$, want onvertekend ($E[\bar{x}] = \mu$) en $MSE[\bar{x}] = \text{Var}[\bar{x}] = \frac{\sigma^2}{n}$

\rightarrow onvertekende schatter met kleinste variantie

$\Rightarrow \bar{x}$ is een UMVU (uniform minimum variantie onvertekend) voor μ

! enkel als alle datapunten uit normale verdeling

• σ^2 schatten met $s^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2$

\rightarrow delen door $n-1$ voor onvertekend

\rightarrow 1 v.d. n vrijheidsgr. gebruikt voor \bar{x}

st. Zij X_1, \dots, X_n steekproeven uit $\mathcal{N}(\mu, \sigma^2)$ -verdeelde populatie, dan $E[S^2] = \sigma^2$

zie slide

43-44

voor bew

Voor $S_*^2 = \frac{1}{n} \sum_i (x_i - \bar{x})^2$ geldt: $E[S_*^2] = E[\frac{n-1}{n} S^2] = \frac{n-1}{n} \sigma^2$

$\rightarrow S_*^2$ geen onvertekende schatter voor σ^2

st. **Stelling v. Helmert**: $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ onafh. met $n \geq 2$, dan

• $\bar{x} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n})$

• \bar{x} en S^2 zijn onafhankelijk

• $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$

st. X normaal verdeeld a.s.a. \bar{x} en S^2 onafh.

zie oef

1.1 a en c

• $\text{Var}[S^2] = \frac{2\sigma^4}{n-1}$

zie slides 55-56

Schatten van proporties of kansen

→ relatieve frequentie / proportie successen in steekproef:

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \bar{x}$$

• $E[\hat{p}] = p$ → onvertkend

• $\text{Var}[\hat{p}] = \frac{p(1-p)}{n}$

• $\frac{\hat{p}-p}{\sqrt{\frac{p(1-p)}{n}}} \stackrel{\text{CLS}}{\approx} \mathcal{N}(0,1)$ indien n voldoende groot

Betrouwbaarheidsintervallen

↳ afhankelijk v.d. steekproef

def. Stel X_1, \dots, X_n steekproeven uit X en dichtheid van X hangt af van onbekende parameter $\theta \in \Theta \subset \mathbb{R}$. Kies $0 < \alpha < 1$.

Een $(1-\alpha)$ -betrouwbaarheidsinterval is een interval

van de vorm $[L(X_1, \dots, X_n), R(X_1, \dots, X_n)]$ met

$$P(L(X_1, \dots, X_n) \leq \theta \leq R(X_1, \dots, X_n)) = 1 - \alpha$$

↳ α is het significantieniveau

$|\frac{R-L}{2}|$ is de foutenmarge bij een tweezijdige test

def. • tweezijdig: $P(\theta < L) = P(\theta > R) = \frac{\alpha}{2}$

• Linkseenzijdig: $]-\infty, R]$

• rechtseenzijdig: $[L, +\infty[$

Strategie om BI te construeren:

1) goede schatter T voor θ

2) verdeling van $h(T, \theta)$ onafh. v. θ

3) a en b zodat $P(a \leq h(T, \theta) \leq b) = 1 - \alpha$

4) herschrijf zodat θ middelste deel is

BI voor μ met bekende σ^2

$$\bar{X}_n \sim \mathcal{N}(\mu, \frac{\sigma^2}{n}), \quad Z = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim \mathcal{N}(0,1)$$

$$\Rightarrow P(-1,96 \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq 1,96) = 0,95$$

$$\Rightarrow P(\bar{X}_n - 1,96 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + 1,96 \frac{\sigma}{\sqrt{n}}) = 95\%$$

$$\Rightarrow \text{BI: } [\bar{X}_n - 1,96 \frac{\sigma}{\sqrt{n}}, \bar{X}_n + 1,96 \frac{\sigma}{\sqrt{n}}]$$

→ deze intervallen zullen in 95% v.d. gevallen μ bevatten

$1-\alpha$ = het betrouwbaarheidsniveau

$$\text{algemeen: } [\bar{X}_n - z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}]$$

• $\alpha \downarrow \Rightarrow \text{BI} \leftarrow \leftarrow$ • $\sigma \uparrow \Rightarrow \text{BI} \leftarrow \leftarrow$ • $n \uparrow \Rightarrow \text{BI} \rightarrow \leftarrow$

⚠ meestal is σ niet op voorhand gekend

Variantie van normale verdeling

$$X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2), \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2, \quad E[S^2] = \sigma^2$$

→ st. v. Helmert → $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$

$$\text{BI: } \left[\frac{(n-1)S^2}{\chi_{n-1, 1-\alpha/2}^2}, \frac{(n-1)S^2}{\chi_{n-1, \alpha/2}^2} \right]$$

⚠ niet symmetrisch rond s^2

→ wortels nemen voor σ

BI voor μ , σ ongekend

$$X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2), \quad \bar{X}_n \sim (\mu, \frac{\sigma^2}{n})$$

bewp. 41 st. Als $X \sim \mathcal{N}(\mu, \sigma^2)$, dan $T = \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \sim t_{n-1}$

↳ $n-1$ vrijheidsgraden

$$\text{BI: } \left[\bar{X}_n - t_{n-1, 1-\alpha/2} \frac{S}{\sqrt{n}}, \bar{X}_n + t_{n-1, 1-\alpha/2} \frac{S}{\sqrt{n}} \right]$$

$$t_{n-1, 1-\alpha/2} = -t_{n-1, \alpha/2}$$

→ exact resultaat

$n \rightarrow \infty \Rightarrow S \rightarrow \sigma \Rightarrow t_{n-1, 1-\alpha/2} \rightarrow z_{1-\alpha/2}$ (1^e resultaat)

Proporctie

$$\hat{p} = \frac{\# \text{ waargenomen successen}}{n} = \frac{1}{n} \sum_{i=1}^n x_i = \bar{X}_n$$

$$E[\hat{p}] = \frac{1}{n} \sum_{i=1}^n E[x_i] = \frac{np}{n} = p$$

$$n \hat{p} = \sum_{i=1}^n x_i \sim B(n, p)$$

als $n \geq 30$, $np \geq 5$ en $nq \geq 5$, dan:

$$B(n, p) \approx N(np, npq)$$

$$\Rightarrow n \hat{p} \approx N(np, npq)$$

$$\Rightarrow \hat{p} \approx N\left(p, \frac{pq}{n}\right)$$

n groot genoeg $\Rightarrow \sigma = \sqrt{\frac{p(1-p)}{n}}$ \bar{w} benaderd door $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

$$\Rightarrow Z = \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \sim N(0, 1)$$

foutenmarge

$$BI: \left[\hat{p} - z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

2 benaderingen ●

$$n \uparrow \Rightarrow BI \rightarrow \leftarrow$$

$$\alpha \downarrow \Rightarrow BI \leftarrow \rightarrow$$

③

hypothesetesten

Algemeen

- **hypothese**: bewering over verdeling of parameter van een verdeling.
- **testen van hypothesen**: op basis van steekproef uitspraak doen over geldigheid v. hypothese.

Stappenplan van 6 stappen:

1) Hypothesen opstellen

- ↳ H_0 en H_1 , twee complementaire testen opstellen
 - ↳ ofwel H_0 verwerpen (H_1 waar)
 - ofwel H_0 niet verwerpen (geen besluit)

→ H_0 moet gekozen \bar{w} dat het gedrag v.d. toetsingsgrootheid te bep. is in de veronderstelling dat H_0 waar is (gelijkheid moet mee in H_0 staan).

→ H_1 is als enige te besluiten of ze waar is.

=> wat we willen bewijzen in H_1

$$\begin{array}{l}
 H_0 \\
 H_1
 \end{array}
 \left\{ \begin{array}{l}
 \theta = \theta_0 \\
 \theta \neq \theta_0
 \end{array} \right.
 \begin{array}{l}
 \text{tweezijdig} \\
 \end{array}
 \left\{ \begin{array}{l}
 \theta \leq \theta_0 \\
 \theta > \theta_0
 \end{array} \right.
 \begin{array}{l}
 \text{rechtseenzijdig} \\
 \end{array}
 \left\{ \begin{array}{l}
 \theta \geq \theta_0 \\
 \theta < \theta_0
 \end{array} \right.
 \begin{array}{l}
 \text{linkseenzijdig} \\
 \end{array}$$

2) Voorwaarden noteren en nagaan

- steekproef willekeurig genomen?
- specifieke voorwaarden

3) Bepalen v.d. teststatistiek

- ↳ statistiek met gekende verdeling
- afhankelijk van de situatie

4) Berekenen v.d. testwaarde

- significantieniveau α

5) Aanvaardingsgebied / p-waarde

→ AG tweezijdige test:

↳ we kunnen steeds a_1 en a_2 vinden zodat

$$P(a_1 \leq T \leq a_2) = 1 - \alpha \mid H_0$$

=> we kunnen steeds een interval opstellen zodat de kans dat, onder H_0 , de waarde van de teststatistiek in het interval valt steeds gelijk is aan $1 - \alpha$.

=> Beslissingsregel:

- $T \in [a_1, a_2] \rightarrow H_0$ niet verwerpen: **zwak besluit**
- $T \notin [a_1, a_2] \rightarrow H_0$ verwerpen: **sterk besluit**

→ complement aanvaardingsgebied = **verwerpingsgebied**

↳ scheidingspunten: **kritische punten**

→ testen kunnen ook uitgevoerd \bar{w} a.d.h.v. BI

→ AG eenzijdige test

	rechtseenzijdig	linkseenzijdig
Onder H_0	$P(T \leq a_3) = 1 - \alpha$	$P(T \geq a_4) = 1 - \alpha$
Aanvaardingsgebied	$] -\infty, a_3]$	$[a_4, \infty[$
Verwerpingsgebied	$] a_3, \infty[$	$] -\infty, a_4[$

→ p-waarde

def. De p-waarde of overschrijdingskans is de kans, onder H_0 , dat de teststatistiek T in een even grote steekproef een waarde aanneemt die minstens zo extreem ligt in de richting v.h. alternatief als de waargenomen testwaarde.
 $p < \alpha \Leftrightarrow$ verwerp H_0 op significantieniveau α
 $\forall p$ onafh. v. α
 ↳ kleinste niveau waarop H_0 \bar{w} verworpen

6) Besluit formuleren

↳ 3 delen:

- testwaarde in aanvaardingsgeb. / $p > \alpha$?
- H_0 verwerpen of niet?
- inhoudelijk besluit

μ normale verd., σ^2 gekend

Tweezijdige z-test

1) $\begin{cases} H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0 \end{cases}$

2) willekeurig?

\bar{X} normaal verdeeld?

3) $Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}} \sim \mathcal{N}(0, 1) \quad | H_0 \rightarrow \text{"onder nulhypothese"}$

4) $z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$

	tweezijdig	$p = 2P(Z \geq z)$
Onder H_0	$P(-z_{1-\alpha/2} \leq Z \leq z_{1-\alpha/2}) = 1 - \alpha$	
AG	$[-z_{1-\alpha/2}, z_{1-\alpha/2}]$	
VG	$] -\infty, -z_{1-\alpha/2}[\cup] z_{1-\alpha/2}, \infty[$	

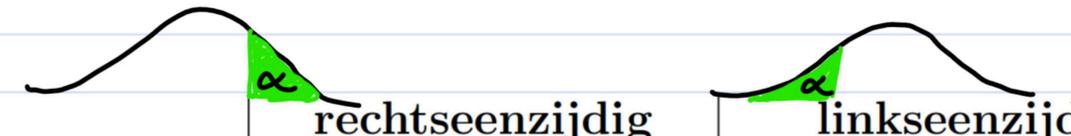
Eenzijdige z-test

1) rechtseenzijdig:

$\begin{cases} H_0: \mu \leq \mu_0 \\ H_1: \mu > \mu_0 \end{cases}$

linkseenzijdig:

$\begin{cases} H_0: \mu \geq \mu_0 \\ H_1: \mu < \mu_0 \end{cases}$



	rechtseenzijdig	linkseenzijdig
Onder H_0	$P(Z \leq z_{1-\alpha}) = 1 - \alpha$	$P(Z \geq -z_{1-\alpha}) = 1 - \alpha$
AG	$] -\infty, z_{1-\alpha}]$	$[-z_{1-\alpha}, +\infty[$
VG	$] z_{1-\alpha}, +\infty[$	$] -\infty, -z_{1-\alpha}[$

$p = P(Z \geq z)$

$p = P(Z \leq z)$

μ norm. verd., σ^2 ongekend

t-test

1) en 2) analoog aan vorige

3) $T = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim_{H_0} t_{n-1}$

4) $t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$

5)

	tweezijdig
Onder H_0	$P(-t_{n-1, 1-\frac{\alpha}{2}} \leq T \leq t_{n-1, 1-\frac{\alpha}{2}}) = 1 - \alpha$
AG	$[-t_{n-1, 1-\frac{\alpha}{2}}, t_{n-1, 1-\frac{\alpha}{2}}]$
VG	$] -\infty, -t_{n-1, 1-\frac{\alpha}{2}} [\cup] t_{n-1, 1-\frac{\alpha}{2}}, \infty [$

	rechtseenzijdig	linkseenzijdig
Onder H_0	$P(T \leq t_{n-1, 1-\alpha}) = 1 - \alpha$	$P(T \geq -t_{n-1, 1-\alpha}) = 1 - \alpha$
AG	$] -\infty, t_{n-1, 1-\alpha} [$	$[-t_{n-1, 1-\alpha}, \infty [$
VG	$] t_{n-1, 1-\alpha}, \infty [$	$] -\infty, -t_{n-1, 1-\alpha} [$

	tweezijdig	rechtseenzijdig	linkseenzijdig
p-waarde	$2P(T \geq t)$	$P(T \geq t)$	$P(T \leq t)$

Fouten

↳ foute conclusie

→ 2 mogelijke fouten: type I en type II fout

	verwerp H_0	verwerp H_0 niet	Type I Error (false-positive)	Type II Error (false-negative)
H_0 waar	Type I fout α	geen fout $1-\alpha$		
H_0 niet waar	geen fout $1-\beta$	Type II fout β		

zie p. 66 bovenaan



H_0 : niet zwanger
 H_1 : zwanger

Type I fout

$$\begin{aligned}
 P(\text{type-I fout}) &= P(H_0 \text{ verwerpen} | H_0 \text{ waar}) \\
 &= P(T \text{ in verwerpingsgeb.} | \mu = \mu_0) \\
 &= P(T \in]-\infty, t_{n-1, 1-\frac{\alpha}{2}} [\cup] t_{n-1, 1-\frac{\alpha}{2}}, \infty [| \mu = \mu_0) \\
 &= \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha \quad (T \sim t_{n-1} \text{ als } \mu = \mu_0)
 \end{aligned}$$

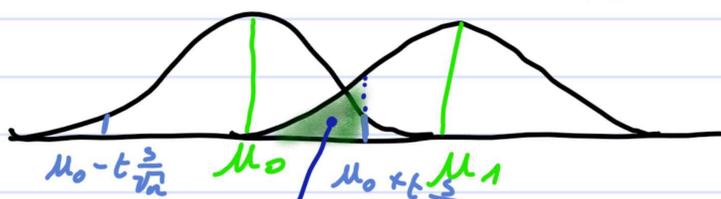
Type II-fout

↳ $\beta = P(\text{type II fout})$, $1 - \beta$ is de power / onderscheidingsvermogen

↳ kan niet "gekozen" \bar{w}

→ enkel bepalen voor $H_1: \mu = \mu_1 \neq \mu_0$

$$\begin{aligned}
 P(\text{type II-fout} | \mu = \mu_1) &= P(H_0 \text{ niet verwerpen} | H_1: \mu = \mu_1 \text{ waar}) \\
 &= P(-t_{n-1, 1-\frac{\alpha}{2}} \leq \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \leq t_{n-1, 1-\frac{\alpha}{2}} | \mu = \mu_1) \\
 \text{voor } \mu = \mu_1: T &= \frac{\bar{x} - \mu_1}{s/\sqrt{n}} \sim t_{n-1} \\
 \Rightarrow \beta &= P(-t_{n-1, 1-\frac{\alpha}{2}} - \frac{\mu_1 - \mu_0}{s/\sqrt{n}} \leq \frac{\bar{x} - \mu_0}{s/\sqrt{n}} - \frac{\mu_1 - \mu_0}{s/\sqrt{n}} \leq t_{n-1, 1-\frac{\alpha}{2}} - \frac{\mu_1 - \mu_0}{s/\sqrt{n}} | \mu = \mu_1) \\
 &= P(-t_{n-1, 1-\frac{\alpha}{2}} - \frac{\mu_1 - \mu_0}{s/\sqrt{n}} \leq T \leq t_{n-1, 1-\frac{\alpha}{2}} - \frac{\mu_1 - \mu_0}{s/\sqrt{n}})
 \end{aligned}$$



↳ kans op type II-fout → verkleinen door:

- $s \downarrow$
- $n \uparrow$
- $\alpha \uparrow$ ∴ type I \uparrow

→ α beter begr. dan $\beta \Rightarrow H_0$ & H_1 zo kiezen zodat fouten v.d. 1^e soort "belangrijker" zijn.

→ onterecht verwerpen is erger

Normaliteitsassumptie

$\begin{cases} H_0: \text{geg. zijn normaal verdeeld} \\ H_1: \text{geg. zijn niet normaal verdeeld} \end{cases}$

→ verwerp H_0 als r_Q (correlatiecoëff. v. normale kwantielplot) $\ll 1$

→ X niet normaal verdeeld?

→ Logaritme wel? → $\Rightarrow \text{Med}(X) = e^{\mu_0}$

$\begin{cases} H_0: \mu_{\log(X)} = \log(\mu_0) \\ H_1: \mu_{\log(X)} \neq \log(\mu_0) \end{cases} \Leftrightarrow \begin{cases} H_0: \mu_X = e^{\mu_0} \\ H_1: \mu_X \neq e^{\mu_0} \end{cases}$

Variantie van normale verd.

$X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ met μ en σ onbekend

1) & 2) zoals voorheen

3) $X^2 = \frac{(n-1)S^2}{\sigma_0^2} \sim_{H_0} \chi^2_{n-1}$

4) $X^2 = \frac{(n-1)S^2}{\sigma_0^2}$

5)

	tweezijdig
Onder H_0	$P(\chi^2_{n-1, \frac{\alpha}{2}} \leq X^2 \leq \chi^2_{n-1, 1-\frac{\alpha}{2}}) = 1 - \alpha$
AG	$[\chi^2_{n-1, \frac{\alpha}{2}}, \chi^2_{n-1, 1-\frac{\alpha}{2}}]$
VG	$[0, \chi^2_{n-1, \frac{\alpha}{2}}[\cup]\chi^2_{n-1, 1-\frac{\alpha}{2}}, +\infty[$

$p = 2 \min \{P(X^2 \geq x^2), P(X^2 \leq x^2)\}$

	rechtseenzijdig	linkseenzijdig
Onder H_0	$P(X^2 \leq \chi^2_{n-1, 1-\alpha}) = 1 - \alpha$	$P(X^2 \geq \chi^2_{n-1, \alpha}) = 1 - \alpha$
AG	$[0, \chi^2_{n-1, 1-\alpha}]$	$[\chi^2_{n-1, \alpha}, +\infty[$
VG	$] \chi^2_{n-1, 1-\alpha}, +\infty[$	$[0, \chi^2_{n-1, \alpha}[$

$p = P(X^2 \geq x^2)$

$p = P(X^2 \leq x^2)$

Proportie

$X_1, \dots, X_n \sim \mathcal{B}(1, p)$

→ $\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$ of als n voldoende groot: $\hat{p} \approx \mathcal{N}(p, \frac{p(1-p)}{n})$

⇒ 2 mogelijke testen

Exacte test

3) $X \sim_{H_0} \mathcal{B}(n, p_0)$

4) $x = n\hat{p} = \# \text{successen in de steekproef}$

→ enkel p -waarde (AG ingewikkeld door discrete verd.)

$p = P(|X - np_0| \geq |x - np_0|)$

$= P(X - np_0 \geq |x - np_0|) + P(X - np_0 \leq -|x - np_0|)$

2 mogelijkheden:

① $x \geq np_0$

⇒ $P(X - np_0 \geq x - np_0) + P(X - np_0 \leq -x + np_0) = p$
 $= P(X \geq x) + P(X \leq 2np_0 - x)$

② $x \leq np_0$

⇒ $p = P(X \leq x) + P(X \geq 2np_0 - x)$

5)

	tweezijdig
$x \geq np_0$	$P(X \geq x) + P(X \leq 2np_0 - x)$
$x \leq np_0$	$P(X \leq x) + P(X \geq 2np_0 - x)$

rechtseenzijdig	linkseenzijdig
$P(X \geq x)$	$P(X \leq x)$

Benaderende test

2) $\forall n > 30, np > 5$ en $n(1-p) > 5 \rightarrow p$ benaderen door \hat{p}

3) $Z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \underset{H_0}{\approx} \mathcal{N}(0,1) \quad E[\hat{p}] = p_0, \text{Var}[\hat{p}] = \frac{p_0(1-p_0)}{n}$

4) $z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$

5)

	tweezijdig
Onder H_0	$P(-z_{1-\alpha/2} \leq Z \leq z_{1-\alpha/2}) = 1 - \alpha$
AG	$[-z_{1-\alpha/2}, z_{1-\alpha/2}]$
VG	$] -\infty, -z_{1-\alpha/2}[\cup] z_{1-\alpha/2}, \infty[$

$p = 2P(Z \geq |z|)$

	rechtseenzijdig	linkseenzijdig
Onder H_0	$P(Z \leq z_{1-\alpha}) = 1 - \alpha$	$P(Z \geq -z_{1-\alpha}) = 1 - \alpha$
AG	$] -\infty, z_{1-\alpha}]$	$[-z_{1-\alpha}, \infty[$
VG	$] z_{1-\alpha}, \infty[$	$] -\infty, -z_{1-\alpha}[$

$p = P(Z \geq z) \quad p = P(Z \leq z)$

Vergelijken van varianties

$X_1, \dots, X_n \sim \mathcal{N}(\mu_1, \sigma_1^2)$, μ_1 en σ_1 onbekend

$Y_1, \dots, Y_m \sim \mathcal{N}(\mu_2, \sigma_2^2)$, μ_2 en σ_2 onbekend

$\rightarrow \sigma_1^2$ en σ_2^2 schatten:

$S_1^2 = \frac{1}{n-1} \sum_i (X_i - \bar{X}_n)^2$

$S_2^2 = \frac{1}{m-1} \sum_j (Y_j - \bar{Y}_m)^2$

st. als σ_1 en σ_2 onbekend, dan: $\frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \sim F_{n-1, m-1}$

$\Rightarrow P(F_{n-1, m-1, \frac{\alpha}{2}} \leq \frac{S_1^2/\sigma_1^2}{S_2^2/\sigma_2^2} \leq F_{n-1, m-1, 1-\frac{\alpha}{2}}) = 1 - \alpha$

met $F_{n, m, \frac{\alpha}{2}} = \frac{1}{F_{m, n, 1-\frac{\alpha}{2}}}$

\Rightarrow BI voor σ_1^2/σ_2^2 :

$[\frac{S_1^2}{S_2^2} \frac{1}{F_{n-1, m-1, 1-\frac{\alpha}{2}}}, \frac{S_1^2}{S_2^2} F_{n-1, m-1, \frac{\alpha}{2}}]$

\rightarrow hypothesetest

1) $\sigma_1^2 = \sigma_2^2$ komt overeen met $\frac{\sigma_1^2}{\sigma_2^2} = 1$, anderen kunnen op analoge manier herschreven worden

∇ let op nummering groepen

2) • willekeurig?

• normaal verd. (voor beide groepen)?

• kwantitatief en groepen onafh.?

• geen uitschieters?

3) omdat $\frac{(n-1)S_1^2}{\sigma_1^2} \sim \chi_{n-1}^2$ en $\frac{(m-1)S_2^2}{\sigma_2^2} \sim \chi_{m-1}^2$ met S_1^2 en S_2^2 onafh.

$F = \frac{S_1^2}{S_2^2} \underset{H_0}{\sim} F_{n-1, m-1}$

4) $f = \frac{S_1^2}{S_2^2}$

5)

	tweezijdig
Onder H_0	$P(F_{n_1-1, n_2-1, \frac{\alpha}{2}} \leq F \leq F_{n_1-1, n_2-1, 1-\frac{\alpha}{2}}) = 1 - \alpha$
AG	$[F_{n_1-1, n_2-1, \frac{\alpha}{2}}, F_{n_1-1, n_2-1, 1-\frac{\alpha}{2}}]$
VG	$[0, F_{n_1-1, n_2-1, \frac{\alpha}{2}}[\cup] F_{n_1-1, n_2-1, 1-\frac{\alpha}{2}}, +\infty[$

	tweezijdig
$f > 1$	$2P(F \geq f)$
$f \leq 1$	$2P(F \leq f)$

	rechtseenzijdig	linkseenzijdig
O. H_0	$P(F \leq F_{n_1-1, n_2-1, 1-\alpha}) = 1 - \alpha$	$P(F \geq F_{n_1-1, n_2-1, 1-\alpha}) = 1 - \alpha$
AG	$[0, F_{n_1-1, n_2-1, 1-\alpha}]$	$[F_{n_1-1, n_2-1, \alpha}, +\infty[$
VG	$] F_{n_1-1, n_2-1, 1-\alpha}, +\infty[$	$] 0, F_{n_1-1, n_2-1, \alpha}[$

rechtseenzijdig $P(F \geq f)$ | linkseenzijdig $P(F \leq f)$

\rightarrow benadering

Vergelijken van gemiddelden

- onderscheid gepaarde en ongepaarde waarnemingen
- geen verband tussen X_1, \dots, X_n en $Y_1, \dots, Y_m \Rightarrow$ ongepaard
- gepaard als verband tussen elke X_i en een Y_j
 $\Rightarrow n=m$
- checken via boxplots

Ongepaarde waarnemingen

1) 2) zie vorige

3) onderscheid gelijke/ongelijke varianties
 - gelijke varianties

als $\sigma = \sigma_1 = \sigma_2$, dan $\bar{X} - \bar{Y} \sim \mathcal{N}(\mu_1 - \mu_2, \frac{\sigma^2}{n} + \frac{\sigma^2}{m})$
 of $\frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\sigma \sqrt{1/n + 1/m}} \sim \mathcal{N}(0, 1)$

→ σ^2 schatten met gewogen/gepoolde/gemengde variantie
 $\bar{S}^2 = \frac{(n-1)S_1^2 + (m-1)S_2^2}{(n-1) + (m-1)} \rightarrow E[\bar{S}^2] = \sigma^2$ (onvertekende sch.)

bew. p. 92 St. $T = \frac{\bar{X} - \bar{Y} - (\mu_1 - \mu_2)}{\bar{S} \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{(n-1)+(m-1)} = t_{n+m-2}$

\Rightarrow BI: $[\bar{x} - \bar{y} \pm t_{n+m-2, 1-\frac{\alpha}{2}} \bar{S} \sqrt{\frac{1}{n} + \frac{1}{m}}]$ ← voor $\mu_1 - \mu_2$
 → hypothesetest

3) $T = \frac{\bar{X} - \bar{Y}}{\bar{S} \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim_{H_0} t_{n+m-2}$ 4) $t = \frac{\bar{x} - \bar{y}}{\bar{S} \sqrt{\frac{1}{n} + \frac{1}{m}}}$

5)

	tweezijdig
Onder H_0	$P(-t_{n_1+n_2-2, 1-\frac{\alpha}{2}} \leq T \leq t_{n_1+n_2-2, 1-\frac{\alpha}{2}}) = 1 - \alpha$
AG	$[-t_{n_1+n_2-2, 1-\frac{\alpha}{2}}, t_{n_1+n_2-2, 1-\frac{\alpha}{2}}]$
VG	$] -\infty, -t_{n_1+n_2-2, 1-\frac{\alpha}{2}} [\cup] t_{n_1+n_2-2, 1-\frac{\alpha}{2}}, \infty [$

$p = 2 P(T \geq |t|)$

	rechtseenzijdig	linkseenzijdig
Onder H_0	$P(T \leq t_{n_1+n_2-2, 1-\alpha}) = 1 - \alpha$	$P(T \geq -t_{n_1+n_2-2, 1-\alpha}) = 1 - \alpha$
AG	$] -\infty, -t_{n_1+n_2-2, 1-\alpha} [$	$[-t_{n_1+n_2-2, 1-\alpha}, \infty [$
VG	$] t_{n_1+n_2-2, 1-\alpha}, \infty [$	$] -\infty, t_{n_1+n_2-2, 1-\alpha} [$

$p = P(T \geq t)$

$p = P(T \leq t)$

→ hypothesetest ongelijke varianties

• n en m groot $\Rightarrow s_1^2 \approx \sigma_1^2, s_2^2 \approx \sigma_2^2$
 $\Rightarrow \bar{X} - \bar{Y} \sim \mathcal{N}(\mu_1 - \mu_2, \frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m})$ of $Z = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{m}}} \sim \mathcal{N}(0, 1)$

• voor kleinere steekproeven:

3) $T = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{m}}} \approx t_r$ met $r = \frac{(\frac{s_1^2}{n} + \frac{s_2^2}{m})^2}{\frac{(s_1^2/n)^2}{n-1} + \frac{(s_2^2/m)^2}{m-1}}$

4) $t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{m}}}$

5)

	tweezijdig
Onder H_0	$P(-t_{r, 1-\frac{\alpha}{2}} \leq T \leq t_{r, 1-\frac{\alpha}{2}}) = 1 - \alpha$
AG	$[-t_{r, 1-\frac{\alpha}{2}}, t_{r, 1-\frac{\alpha}{2}}]$
VG	$] -\infty, -t_{r, 1-\frac{\alpha}{2}} [\cup] t_{r, 1-\frac{\alpha}{2}}, \infty [$

$p = 2 P(T \geq |t|)$

	rechtseenzijdig	linkseenzijdig
Onder H_0	$P(T \leq t_{r, 1-\alpha}) = 1 - \alpha$	$P(T \geq -t_{r, 1-\alpha}) = 1 - \alpha$
AG	$] -\infty, -t_{r, 1-\alpha} [$	$[-t_{r, 1-\alpha}, \infty [$
VG	$] t_{r, 1-\alpha}, \infty [$	$] -\infty, t_{r, 1-\alpha} [$

$p = P(T \geq t)$

$p = P(T \leq t)$

Gepaarde waarnemingen

→ nieuwe variabele invoeren: $D_i = X_i - Y_i$ ($n=m$)

$\sigma^2 \neq \sigma_1^2 + \sigma_2^2$
want afh.

↳ als (X_i, Y_i) bivariaat normaal verdeeld zijn, dan
 $D_1, \dots, D_n \sim \mathcal{N}(\mu, \sigma^2)$ met onbekende $\sigma^2 \neq \sigma_1^2 + \sigma_2^2$ en $\mu = \mu_1 - \mu_2$
→ voer t-test voor μ bij norm. verd. uit (voor D , $\mu \stackrel{H_0}{=} 0$)

Vergelijken van twee proporties

→ onafhankelijke experimenten: $p = p_1 - p_2$, $\sigma^2 = \sigma_1^2 + \sigma_2^2 = \frac{p_1 q_1}{n} + \frac{p_2 q_2}{m}$

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n} + \frac{p_2(1-p_2)}{m}}} \approx \mathcal{N}(0, 1)$$

→ BI: $[(\hat{p}_1 - \hat{p}_2) \pm z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n} + \frac{\hat{p}_2(1-\hat{p}_2)}{m}}]$

→ hypothesetest

↳ 2 manieren

① gegevens van BI gebr. en met normale verdeling werken

② met gemengde proportie werken

omdat getest \bar{w} of $p_1 = p_2$

$$\rightarrow \hat{p}_0 = \frac{n \hat{p}_1 + m \hat{p}_2}{n+m}$$

$$3) Z = \frac{\hat{p}_1 - \hat{p}_2}{\hat{p}_0(1-\hat{p}_0) \sqrt{\frac{1}{n} + \frac{1}{m}}} \approx \mathcal{N}(0, 1)$$

Correlatie tussen twee variabelen

Kwantitatieve variabelen

↳ populatie correlatiecoëfficiënt ρ , als $\rho \neq 0 \Rightarrow X$ en Y afhankelijk

$$1) \begin{cases} H_0: \rho = 0 & \rightarrow \text{als } H_0 \text{ verworpen} \Rightarrow \text{lineaire afh.} \\ H_1: \rho \neq 0 & \forall H_0 \text{ niet verw.} \Rightarrow \text{Lin. onafh.} \end{cases}$$

↳ enkel tweezijdige test

2) • X en Y bivariaat normaal verdeeld?

↳ X en Y normaal verdeeld (grafisch/test)

↳ scatterplot (elliptische wolk)

• willekeurig?

3) R = steekproef correlatiecoëfficiënt

$$T = \frac{R \sqrt{n-2}}{\sqrt{1-R^2}} \sim_{H_0} t_{n-2}$$

$$4) Z = \frac{r \sqrt{n-2}}{\sqrt{1-r^2}}$$

$$5) [-t_{n-2, 1-\frac{\alpha}{2}}, t_{n-2, 1-\frac{\alpha}{2}}]$$

$$p = 2P(T > |Z|)$$

Nominale variabelen

1) $\begin{cases} H_0: X \text{ en } Y \text{ onafhankelijk} \\ H_1: X \text{ en } Y \text{ afhankelijk} \end{cases}$

↳ weer enkel tweezijdig

2) • willekeurig?

• X en Y kwalitatief met minstens 2 mog. uithomsten?

• klassen disjunct?

• som v. alle frequenties > 20 ?

• frequentie per categorie ≥ 5 ? (samenemen)

• alle waarnemingen onafhankelijk?

↳ met bijdrage tot frequentie

Er moet gelden dat $P(X = m_x, Y = m_y) = P(X = m_x)P(Y = m_y)$ voor onafhankelijkheid.

↳ meestal beschikken we niet over de juiste geg.

zie p. 108

⇒ benadering: $\frac{n_{xy}}{n} = \frac{n_x}{n} \frac{n_y}{n}$

→ komt overeen met wat we eerder zagen:

verwachte waarde = $n_{xy} = \frac{\text{rij totaal} \times \text{kolom totaal}}{n} = \frac{n_x n_y}{n} = \frac{n_x}{n} \frac{n_y}{n} n$

• waargenomen frequenties O_i

• verwachte frequenties E_i

⇒ $n = \sum_i O_i = \sum_i E_i$

zie redenering p. 109

3) $\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i} \approx_{H_0} \chi^2_{\substack{(n_1-1)(n_2-2) \\ \# \text{rijen} \quad \# \text{kolommen in kruistabel}}}$

4) $\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i}$

5) $[0, \chi^2_{(m_1-1)(m_2-1), 1-}$

(*)

Regressie

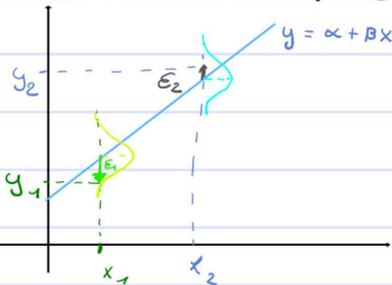
Introductie

- onderzoek of/hoe Y op een lineaire manier voorspeld kan worden door X met lineaire regressie
- Y = responsvariabele / afh. var.
- X = verklarende / onafh. var.
- Hoe ziet dit verband eruit? Welk model gebr?
 - ↳ ideaal: $y_i = \alpha + \beta x_i$
 - ↳ in hoeverre wijkt de realiteit af?
- Welke vwdⁿ voor goed model?
 - ↳ hoe nagaan?
- Hoe betrouwbaar?
 - ↳ welke afwijkingen?

Enkelvoudig lineair regressiemodel

$$Y_i = \alpha + \beta x_i + \epsilon_i$$

↳ fouten ϵ_i zijn onafhankelijk en $\epsilon_i \sim N(0, \sigma_\epsilon^2)$



↳ zelfde voor alle ϵ_i

→ opmerkingen

- $E(Y_i | X_i) = E(\alpha + \beta x_i + \epsilon_i) = \alpha + \beta x_i + \overbrace{E[\epsilon_i]}^{=0} = \alpha + \beta x_i$
 - ↳ β geeft invloed van x_i op gem. v. Y_i
- $\text{Var}[Y_i] = \text{Var}[\epsilon_i] = \sigma_\epsilon^2 \Rightarrow$ homoscedasticiteit
- geen veronderstellingen over x_i en marginale verd. v. Y .

Kleinste kwadratenmethode

→ Hoe $y = \alpha + \beta x$ bepalen?

→ schattingen $\hat{\alpha}$ en $\hat{\beta}$ voor α en β

• definieer $e_i(a, b)$ als vert. afw. van (x_i, y_i) t.o.v. $y = a + bx$

$$e_i(a, b) = y_i - (a + bx_i)$$

→ zoek $\hat{\alpha}$ en $\hat{\beta}$ waarvoor $\sum_i e_i^2(a, b) = \sum_i (y_i - (a + bx_i))^2$ minimaal

$$\Rightarrow \hat{\beta} = \frac{\sum_i (y_i - \bar{y})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} = r \frac{s_y}{s_x}, \quad \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

correlatie

⚠ uit schieters weggooien

⇒ geschatte rechte: $y = \hat{\alpha} + \hat{\beta} x$

• schatting voor elke obs: $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$

• residu's: $e_i = e_i(\hat{\alpha}, \hat{\beta}) = y_i - \hat{y}_i = y_i - (\hat{\alpha} + \hat{\beta} x_i)$

• $\sum_i e_i = 0$ altijd!

• (\bar{x}, \bar{y}) ligt op $y = \hat{\alpha} + \hat{\beta} x$

• als $r = r(x, y) = 0$, dan horizont. rechte ($\hat{\beta} = 0$ en $\hat{\alpha} = \bar{y}$)

⇒ geen lin. verb, dan \bar{y} beste voorsp. voor $Y | x_i$ ongeacht waarde van x_i

→ Schatting voor σ_ϵ^2 !

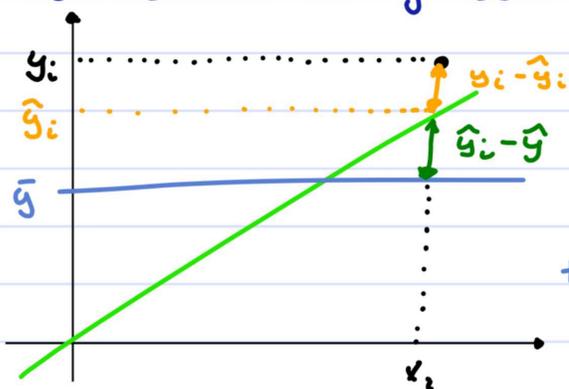
↳ $E(\epsilon_i) = 0 \Rightarrow \sigma_\epsilon^2 = \text{Var}(\epsilon_i) = E(\epsilon_i^2)$

$$s_\epsilon^2 = \hat{\sigma}_\epsilon^2 = \frac{1}{n-2} \sum_i e_i^2(\hat{\alpha}, \hat{\beta}) = \frac{1}{n-2} \sum_i (y_i - \hat{y}_i)^2$$

↳ (

ANOVA-tabel en determinatiecoëfficiënt

Determinatiecoëfficiënt: geeft weer hoe goed het model de variatie in y -waarden verklaart



$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i) \text{ met } \hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$$

$$\Rightarrow \sum_i (y_i - \bar{y})^2 = \sum_i (\hat{y}_i - \bar{y})^2 + \sum_i (y_i - \hat{y}_i)^2$$

$$SST = SSM + SSE$$

total sum of squares = mean sum of squares + error sum of sq.

totale variatie = var verhl. door model + onverhl. var.

→ hoe groter aandeel SSM in SST, hoe beter het model

$$\Rightarrow R^2 = \frac{SSM}{SST} = \frac{\sum_i (\hat{y}_i - \bar{y})^2}{\sum_i (y_i - \bar{y})^2}$$

• enkelvoudige lineaire regressie: $R^2 = r^2$, dus $0 \leq R^2 \leq 1$

• $R^2 = 1 \Rightarrow$ perfecte rechte

• $R^2 \approx 0 \Rightarrow \hat{\beta} = r \frac{s_y}{s_x} \approx 0$

→ andere methode: $f = \frac{MSM}{MSE} = \frac{\text{mean sq. model}}{\text{mean sq. error}} = \frac{\frac{SSM}{1}}{\frac{SSE}{n-2}} = \frac{\sum_i (y_i - \bar{y})^2}{\sum_i (y_i - \hat{y}_i)^2}$

• $MSE = s_e^2$

⇒ Samenvatting in **ANOVA-tabel**:

	vrijheidsgraden	SS	MS	F
Model	1	SSM	$MSM = \frac{SSM}{1}$	$f = \frac{MSM}{MSE}$
Error	$n-2$	SSE	$MSE = \frac{SSE}{n-2}$	
Total	$n-1$	SST		

Nagaan v.d. modelveronderstellingen

↳ lineair regressiemodel soms niet de beste optie

⇒ Wanneer is lineair goed?

↳ niet goed als R^2 klein

↳ onafhankelijkheid van ϵ_i en $N(0, \sigma_\epsilon^2)$ -verdeling van ϵ_i onderzoeken m.b.v. residu's $e_i = y_i - (\hat{\alpha} + \hat{\beta} x_i)$

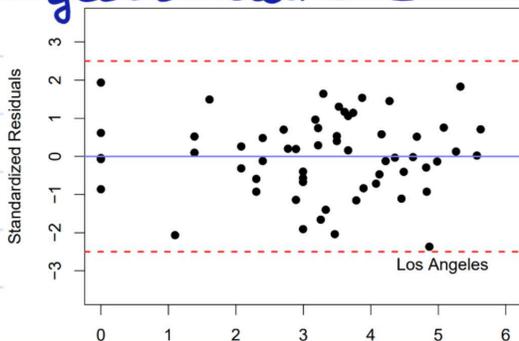
1) ϵ_i normaal verdeeld met $\mu = 0$: normale kwantielplot v.d. residu's e_i

• intercept ≈ 0

• normaliteit niet formeel testen

2) ϵ_i onderling onafhankelijk met dezelfde var. σ_ϵ^2

→ gestandaardiseerde residuplot (x_i, e_i)



→ soms ook (\hat{y}_i, e_i)

→ ongeordend lukraak persoon

→ gelijke spreiding

3) nagaan of er uitschieters zijn in de gestandaardiseerde residu's: $\frac{e_i}{s_{e_i}} \in [-2,5; 2,5]$?

↳ $E_i = Y_i - (\hat{\alpha} + \hat{\beta} x_i)$ toevalsvariabele

• $s_{E_i} = s_\epsilon \sqrt{1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{(n-1)S_x^2}}$

Interferentie omtrent regressieparameters

• kleinste kwadratschatters voor α en β : \hat{A} en \hat{B}

$$\rightarrow \text{Var}[\hat{A}] = \sigma_{\epsilon}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} \right) \Rightarrow \text{s.e.}(\hat{A}) = S_{\epsilon} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2}}$$

$$\text{Var}[\hat{B}] = \frac{\sigma_{\epsilon}^2}{(n-1)s_x^2} \Rightarrow \text{s.e.}(\hat{B}) = S_{\epsilon} \sqrt{\frac{1}{(n-1)s_x^2}}$$

$$\text{Cov}(\hat{A}, \hat{B}) = -\frac{\bar{x} \sigma_{\epsilon}^2}{(n-1)s_x^2}$$

$$\rightarrow \frac{\hat{A} - \alpha}{\text{s.e.}(\hat{A})} \sim t_{n-2}$$

$$\frac{\hat{B} - \beta}{\text{s.e.}(\hat{B})} \sim t_{n-2}$$

$$\rightarrow \text{BI voor } \alpha: \left[\hat{\alpha} \pm t_{n-2, \alpha/2} S_{\epsilon} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2}} \right]$$

$$\text{BI voor } \beta: \left[\hat{\beta} \pm t_{n-2, \alpha/2} \frac{S_{\epsilon}}{\sqrt{(n-1)s_x^2}} \right]$$

→ Als $\beta=0$ kan y niet lineair uit x voorspeld worden

↳ nagaan of lineaire regressie zinvol is met hypothesetest

$$\begin{cases} H_0: \beta = 0 \\ H_1: \beta \neq 0 \end{cases}$$

$$\bullet T = \frac{\hat{B} - \beta}{\text{s.e.}(\hat{B})} \sim t_{n-2}, \quad p = 2P(T \geq |t|)$$

of

$$\bullet F = \frac{MSM}{MSE} \sim F_{1, n-2}, \quad p = P(F > f)$$

$$\rightarrow t^2 = f$$

	Schatting	Standaardfout	Testwaarde	P-waarde
Intercept	$\hat{\alpha}$	$\text{s.e.}(\hat{A})$	$t = \frac{\hat{\alpha}}{\text{s.e.}(\hat{A})}$	$2P(T > t)$
Slope	$\hat{\beta}$	$\text{s.e.}(\hat{B})$	$t = \frac{\hat{\beta}}{\text{s.e.}(\hat{B})}$	$2P(T > t)$