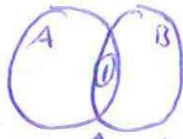


# Statistiek

H2 - events + probability

- Intersection

$$A \cap B, \textcircled{1}$$



↑  $\rightarrow$  stukje overlap is intersection ( $\cap$ )

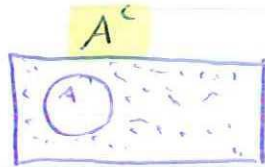
- Union

$$A \cup B$$



$\leftarrow$  Zoude A als B,  $A+B$ , is de union, ( $\cup$ )

- Complement



$\rightarrow$  het complement is het tegenovergesteld van A, alles behalve A

- De Morgan's law

$$(A \cup B)^c = A^c \cap B^c, \quad (A \cap B)^c = A^c \cup B^c$$

- Disjoint



$\rightarrow$  A en B hebben geen overlap.

- Probability function

$P(A)$  = kans op A

$\hookrightarrow$  Disjoint functions  $\rightarrow P(A \cup B) = P(A) + P(B)$

Vb. Hoe groot is de kans dat picht gekozen wordt uit de groep: {Picht, Jan, Klaas, Picket, Bart, Dirk}

-  $\int$  van de 6x wordt picht niet gekozen

$$P(\text{picht}) = \frac{1}{6}$$

- Probability of an union

- Probability of complement

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

$$P(A^c) = 1 - P(A)$$

### H3, Conditional probjt multibayes

#### • Conditional probability

vb. - lange maanden; (jan, mar, mei, jul, aug, oct, dec) = 7

- maanden met r erin; (jan, feb, mar, apr, sep, oct, nov, dec) = 8

$$\text{kans} - P(\text{lang}) = \frac{7}{12}$$

$$- P(r) = \frac{8}{12}$$

$$\rightarrow P(R|L) = \frac{4}{7}$$

↳ maanden met r die lang zijn

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

#### • multiplication rule

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

#### - Bayes rule

$$P(C_i|A) = \frac{P(A|C_i) \cdot P(C_i)}{P(A|C_1)P(C_1) + P(A|C_2)P(C_2) + \dots}$$

vb.

$$P(B|T) = \frac{P(T|B)}{P(T)} = \frac{P(T|B) \cdot P(B)}{P(T|B)P(B) + P(T|B^c)P(B^c)}$$

Onafhankelijk?

$$P(A|B) = P(A)$$

$$P(A) = P(A|B) \cdot P(B) + P(A|B^c) \cdot P(B^c)$$

Hu,

- Bernoulli

Ber(p)  $\rightarrow$  het gebeurt of het gebeurt niet

$$P(X=1) = \text{succes} = p \quad P(X=0) = \text{failure} = 1-p$$

- Binomial

Bin(n,p)  $\rightarrow$  het gebeurt of niet, in n keer gevallen

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

p = kans

n = aantal maal

k = gevraagd

- geometrisch

$$\text{Geo}(p) \quad P(X=k) = (1-p)^{k-1} \cdot p$$

$$\hookrightarrow P(X > n) = 1 - p^n$$

Vb. bin

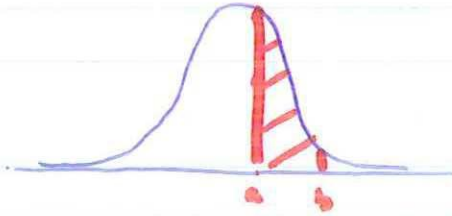
1000 lampen, kans dat lampje stuk is = 0,001 (p).

$$P(X=0) = \binom{1000}{0} 0,001^0 (1-0,001)^{1000-0} = \dots$$

HS

-  $X$  is continuous

$$P(a \leq x < b) = \int_a^b f(x) dx$$



- uniform distribution-

$$\text{Interval} = [\alpha, \beta]$$

$$f(x) = \frac{1}{\beta - \alpha} \quad \text{for } \alpha \leq x \leq \beta$$

generator:  $U(\alpha, \beta)$

- exponential distribution-

$$f(x) = \lambda e^{-\lambda x} \quad \lambda = \text{parameter}$$

$$F(x) = 1 - e^{-\lambda x}$$

generator:  $\text{Exp}(\lambda)$

- normal distribution-

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

generator:  $N(\mu, \sigma^2)$

H7, expectation, variance

- expectations

discrete

$$E[X] = \sum a_i \cdot P(a_i)$$

vb.	a	2	4	6
	P(a)	0,2	0,4	0,2

$$E[X] = (2 \cdot 0,2) + (4 \cdot 0,4) + (6 \cdot 0,2) = 0,4 + 1,6 + 1,2 = 3,2$$

continuous

$$E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

vb. uniform distr. exp  $[2,5] \rightarrow E[X] = \frac{\alpha + \beta}{2}$

$$E[X] = \int_{-\infty}^{\infty} x \cdot f(x) dx = \int_{2}^{5} x \cdot \frac{1}{\beta - \alpha} dx = \int_{2}^{5} x \cdot \frac{1}{3} dx$$

$$\frac{1}{6} x^2 \Big|_{2}^{5} = \frac{25}{6} - \frac{4}{6} = \frac{21}{6} = 3\frac{1}{2}$$

vb. geom.

$$E[X] = \sum_{k=1}^{\infty} k \cdot p \cdot (1-p)^{k-1} = \frac{1}{p} \quad | \text{Geo}(p)$$

vb. exp

$$E[X] = \frac{1}{\lambda}$$

| EXP( $\lambda$ )

vb. normal

$$E[X] = \mu$$

| N( $\sigma^2, \mu$ )

- Variance

$$\text{Var}[X] = E[X^2] - E[X]^2$$

Ub. normal  $\text{Var}[X] = \sigma^2$

Ub. Exp  $\text{Var}[X] = 1/\lambda^2$

Ub. Uniform  $\text{Var}[X] = \frac{(\beta - \alpha)^2}{12}$

- mit Zanderungen:

$$E[rX + s] = r \cdot E[X] + s$$

$$\text{Var}[rX + s] = r^2 \cdot \text{Var}[X]$$



HP

- Change of units transformation

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

Stel  $P(X \leq 5)$  voor  $N(4, 25)$

Dan moet deze kans eerst omgeschreven worden naar een  $N(0,1)$ !

$$\text{Daarvoor: } Z = \frac{X - \mu}{\sigma}$$

$$Z = \frac{X - 4}{5} \rightarrow$$

$$P(X \leq 5) \rightarrow P\left(Z \leq \frac{5 - 4}{5}\right) = P(Z \leq 0,2) = 1 - 0,4207 = 0,5793$$

- Jensen Inequality

$$E[g(x)] = g(E[x])$$

echter niet bij convex functies,  $z^c$  afgl.  $> 0$

$$\text{dan } g(E[x]) \leq E[g(x)]$$

$$\text{vb. } g(x) = x^2 \rightarrow g'(x) = 2x, \text{ convex}$$

$$\text{vb. } g(x) = e^{-x} \rightarrow g''(x) = e^{-x}, \text{ convex}$$

H3

= Joint probability

$$P(a, b) = P(X=a, Y=b)$$

↳ Distribution

$$F(a, b) = P(X \leq a, Y \leq b)$$

↳ continuous distr.

$$P(a_1 \leq X \leq b_1, a_2 \leq Y \leq b_2) = \int_{a_1}^{b_1} \int_{a_2}^{b_2} f(x, y) dx dy$$



H10, covariance

- 2 dimension. expectation  
discrete

$$E[g(x, y)] = \sum_i \sum_j g(a_i, b_j) p(x=a_i, y=b_j)$$

continuous

$$E[g(x, y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) dx dy$$

	a			
vb. ①	0	1	2	$p(x=b)$
0	0	$\frac{1}{4}$	0	$\frac{1}{4}$
1	$\frac{1}{4}$	0	$\frac{1}{4}$	$\frac{1}{2}$
2	0	$\frac{1}{4}$	0	$\frac{1}{4}$
$p(x=a)$	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$	1

$$E[xy] = (0 \cdot 0) \cdot 0 + (0 \cdot 1) \cdot \frac{1}{4} + (2 \cdot 0) \cdot 0 + \dots \dots \dots \text{enzw.}$$

$$= (a \cdot b) \cdot p(x=a, y=b)$$

Let op?  $E[x+y] = E[x] + E[y]$

- covariance

$$\text{COV}(x, y) = E[xy] - E[x] \cdot E[y]$$

vb ①, gebruik tabel ①

$$E[xy] = 1, \quad E[x] = 1, \quad E[y] = 1$$

$$\text{COV} = 1 - 1 = 0$$

$$\text{Var}(x+y) = \text{Var}(x) + \text{Var}(y) + 2\text{COV}(x, y)$$

$$\text{COV}(rx + s, ty + u) = r \cdot t \cdot \text{COV}(x, y)$$

- Density, coefficient

$$\rho(x, y) = \frac{\text{COV}(x, y)}{\sqrt{\text{Var}(x) \cdot \text{Var}(y)}}$$

H12, Poisson

$$P(X=k) = \frac{\mu^k}{k!} e^{-\mu}$$

$\mu = \lambda \cdot t \rightarrow$  aantal p. tijds eenheid

geverdeeld  $Pois(\mu)$

Uv. 12.6

er zit, in een koperen draad, om de 40cm een fout

Wat is de kans op 1 fout in een meter

$\mu =$  aantal verwachten fout in 1 meter  $\approx 2,5$

$k = 1$

$$P(X=1) = \frac{2,5^1}{1!} \cdot e^{-2,5} \approx 0,21$$

$$P(X=2) = \frac{2,5^2}{2!} e^{-2,5} \approx 0,26$$

= Expectatie, variance

$$E[X] = \mu, \text{ var}[X] = \mu$$

H13, law of big numbers

- average

$$\bar{X}_n = \frac{X_1 + X_2 + X_3 + \dots + X_n}{n} = \text{average}$$

- expectation, variance

$$E[\bar{X}_n] = \mu, \quad \text{var}[\bar{X}_n] = \frac{\sigma^2}{n}$$

H14, central limit

$$Z_n = \sqrt{n} \cdot \frac{\bar{X}_n - \mu}{\sigma} = \frac{\bar{X}_n - E[X_n]}{\sqrt{\text{Var}[X_n]}}$$

# Statistik

H15  
- Histogram

$$\text{bin width} = \frac{\text{number } x_i}{n} \quad n = \text{total elements}$$

$$\text{bin height} = \frac{\text{number } x_j}{n \cdot |B_i|} \quad |B_i| = \text{bin width}$$

H16

43, 43, 41, 41, 41, 42, 43, 50, 50, 41

- Sample mean

$$\text{average: } \bar{x}_n = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n} = 45$$

- Sample median

middle von alle  $\rightarrow 42$

- Sample variance

$$s_n^2 = \frac{1}{n} \cdot \sum (x_i - \bar{x}_n)^2$$

$\hookrightarrow$  Sample standard deviation

$$s_n = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x}_n)^2}$$

- mad

$$\text{mad} = \text{median}(|x_1 - \text{med}|, |x_2 - \text{med}|, \dots)$$

stel 3 u 50 g  $\rightarrow \text{med} = 5$

$$\text{mad} = \text{med}(2, 1, 0, 3, 4) = 2$$

~~H19~~ H19, estimators

- Schotker

$$T = e^{-\bar{x}_n}$$

- schotkers:  $\mu \rightarrow \bar{x}_n$   
 $\sigma^2 = s_n^2$  (zie H16)

H20, MSE

$$E[\bar{x}_n] = \frac{E[x_1] + E[x_2] + \dots + E[x_n]}{n}$$

- first estimator, based on sample mean

$$T_1 = 2\bar{x}_n + 1$$

$$E[T_1] = N = \text{number}$$

- second estimator, based on sample maximum

$$T_2 = \frac{n+1}{n} \cdot \underset{\text{max}}{m_n} - 1$$

$$E[T_2] = N$$

$$\text{vb. } 61 = 19 \cdot 56 \cdot 24 = 16$$

$$T_1 = 2\bar{x}_n + 1 \rightarrow \bar{x}_n = 35,2 \rightarrow T_1 = 71,4$$

$$T_2 = 72,2 = \frac{5+1}{5} \cdot 61 - 1$$

Hzi, likelihood.

Zie Smoker p. 313

$$L(p) = C \cdot p(x_1=1)^{29} \cdot p(x_2=2)^{16} \cdot p(x_3=3)^{17} \cdot \dots \cdot p(x_{12}=12)^{42}$$

geometric distribution

~~$$L(p) = C \cdot p^{29} \cdot (1-p)^{16} \cdot p^{16} \cdot (1-p)^{17} \cdot \dots \cdot p^{42} \cdot (1-p)^{42}$$~~

$$L(p) = C \cdot p^{93} \cdot ((1-p)p)^{16} \cdot ((1-p)p^2)^{17} \cdot \dots \cdot ((1-p)^{12})^7$$

$$= C \cdot p^{93} \cdot (1-p)^{322}$$

max likelihood

$$L'(p) = C \cdot p^{92} (1-p)^{321} (93 - 415p) = 0$$

$$\Rightarrow p=0, p=1, \underline{p=0,224}$$



$$L = C \cdot P(x=0,1)^{440} \cdot P(x=2)^{93} \cdot P(x=3)^{9} \cdot P(x=4)^7 \\ - P(x=7)^1 \\ = C \cdot \left( \frac{\mu^2}{2!} e^{-\mu} \right)^{440}$$

H22

Kleinst - quadrat - methode

Data: paar:  $(x_1, y_1), (x_2, y_2), (x_n, y_n)$

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

$$\beta = \frac{n \cdot (\sum x_i y_i) - (\sum x_i) \cdot (\sum y_i)}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\alpha = \bar{y}_n - \beta \cdot \bar{x}_n$$

H23

Betrachtbarkeits interval

$$N(0,1) \rightarrow z = \frac{\bar{x}_n - \mu}{\sigma_x / \sqrt{n}}$$

$$P(-1,96 \leq z \leq 1,96) = 0,95 = 95\% \text{ betrachtbarkeit}$$

$$P(-1,96 \leq \frac{\bar{x}_n - \mu}{\sigma_x / \sqrt{n}} \leq 1,96)$$

$$P(\bar{x}_n - 1,96 \cdot \frac{\sigma_x}{\sqrt{n}} \leq \mu \leq \bar{x}_n + 1,96 \cdot \frac{\sigma_x}{\sqrt{n}})$$

Interval  $(\bar{x}_n - 1,96 \frac{\sigma_x}{\sqrt{n}}, \bar{x}_n + 1,96 \frac{\sigma_x}{\sqrt{n}})$  heist ca 95% betrachtbarkeits interval.

$\sigma$  niet bekend:

$$T_{\text{(scheltes)}} = \frac{\bar{X}_n - \mu}{s_n / \sqrt{n}} = \text{steekende gemiddelt}$$

H2c.

hypothese:

$$N(\mu, \sigma)$$

$$\sigma = 2$$

$H_0$  = nulhypothese

$H_1$  = alternatieve hypothese

vb. snelheidslimiet = 120 km/h

↳ steekproef 5% is vals bevestig.

↳ 3 metingen

$$H_0 = \mu = 120$$

$$H_1 = \mu \geq 120$$

$$T = \frac{\bar{X}_n - 120}{2/\sqrt{3}} \rightarrow T \sim N(0, 1)$$

$$z_{0,05} = 1,645 = T$$

$$H_0 = \mu = 120$$

$$H_1 = \mu \geq 120 \text{ c.15\%} \rightarrow 1,645$$

stel  $\mu = 123$

$$P(T < 1,645 | \mu = 123) = \frac{120 - 123}{2/\sqrt{3}} + 1,645$$

$$= \Phi(-0,95) =$$

stel  $\mu = 121$

$$P(T < 1,645 | \mu = 121) = \frac{120 - 121}{2/\sqrt{3}} + 1,645 = \Phi(-0,70)$$

unbiased estimator

↳ expectation  $\mu$ :  $\bar{x}_n = \frac{x_1 + \dots + x_n}{n} \rightarrow E[x_i]$

↳ variance  $\sigma^2$ :  $s_n^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x}_n)^2 \rightarrow \text{Var}(x_i)$

Zuivere Schötter

↳ gemiddelde waarden

$$\bar{x}_n = \frac{x_1 + \dots + x_n}{n}$$

(T<sub>1</sub>)

$$E[\bar{x}_n] = \frac{1}{2}(n+1) \rightarrow \text{Var}[T] = \text{Var}[2 \cdot \bar{x}_n - 1]$$
$$\rightarrow T = 2 \cdot \bar{x}_n - 1$$
$$= 4 \cdot \text{Var}[\bar{x}_n]$$
$$= \frac{(n+1)(n-1)}{3n}$$

↳ max genome waarde

$$\bar{m}_n = \max[x_1, x_2, \dots, x_n]$$

(T<sub>2</sub>)

$$E[\bar{m}_n] = \frac{n}{n+1} (n+1) \rightarrow \text{Var}[T] = \frac{(n+1)(n-1)}{n(n+2)}$$

$$T = \frac{n+1}{n} \bar{m}_n - 1$$

---

$$\frac{\text{Var}(T_1)}{\text{Var}(T_2)} = \frac{n+2}{3}$$

T<sub>2</sub> = zuiverder wanneer er <sup>groter</sup> ~~er~~ gelid > dan 1 uit komk.

Mean squared error

$$MSE = \text{Var}(T) - (E(T) - \theta)^2$$

$$= \text{variatie} - \text{afwijk}^2$$

Likelihood

bar  $\rightarrow$  p of  $1-p$

	1	2	3	4	5
kop	0	0	1	1	1
mont	3	3	2	2	2

$$P(\bar{x}_n) = \cancel{p^0} (1-p)^{- (1-p)}$$

$$= p^0 \cdot ((1-p)p)^0 \cdot ((1-p)^2)^1 \cdot ((1-p)^3 p)^1 \cdot ((1-p)^4 p)^1$$

$$= 1 - 1 - (1-p)^2 p - (1-p)^3 p - (1-p)^4 p$$

$$= (1-p)^2 p$$

$$= (1-p)^2 p^3$$

MAD = media van median afw.

Stel: 42, 48, 50, 51, 57, 59, 62, 72

media = 54

afw: 12, 6, 4, 3, 3, 4, 8, 18

: 3, 3, 4, 4, 6, 8, 12, 18

↓

media van media afw

$$MAD = \frac{1}{2}(4+6) = 5$$

Stekproef afw.

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

$$\bar{x}_n = 55$$

$$s_8^2 = \frac{1}{7} \sum_{i=1}^8 (x_i - 55)^2$$

$$= \frac{1}{7} \cdot ((42-55)^2 + (48-55)^2 + \dots)$$

$$= \frac{1}{7} \cdot 610 = 87,1428$$

$$s = 9,34$$

empirische verdelings functie

$$F_8(60) = \frac{\text{aantal punten} \leq 60}{n} = \frac{6}{8} = 0,75$$



Kans-massa-funktion  $S = a + b = x + y$

$P(a,b)$	b					a+b			
	0	1	2	3		0	1	2	3
a									
-1	0,05	0,15	0,05	0,1	-1	-1	0	1	+2
0	0,05	0,15	0,15	0,05	→ 0	0	1	2	3
2	0,05	0,2	0,05	0,1	2	2	3	4	5

S	0	1	2	3	4	5
P(S)	0,05	0,15	0,3	0,2	0,2	0,1