

Data-Science 1

kansen



Oefeningen huiswerk

- heeft het zin om een regressielijn te bepalen tussen schoenmaat en lengte?
- Kijk eens naar de correlatie tussen opwarming en zakgeld. Wat betekent dat? Heeft het zin om hier een regressielijn te bepalen?

Oefeningen huiswerk

- bereken een lineaire regressie die de schoenmaat voorspelt adh van de lengte
 - verwijder eerst de uitschieters
 - wat is de vergelijking van de rechte?
- welke schoenmaat voorspel je voor iemand van 180cm groot?
- wat is de gemiddelde fout op de voorspelling? Is dit veel of weinig?
- wat is de verklaarde variantie? Wat betekent dit?
- doe eens een logaritmische regressie. Welke vergelijking kom je uit? Geeft dit een beter of slechter resultaat dan een lineaire?

Oefeningen huiswerk

- In het leerboek “Edexcel GCSE Mathematics” van Keith Pledger wordt beweert dat (voor jongens) de lengte in cm gelijk is aan 5,3 maal (Engelse) schoenmaat plus 133. Welke formule bekom je voor de informatica studenten?

je kan EU maten omzetten naar UK met:

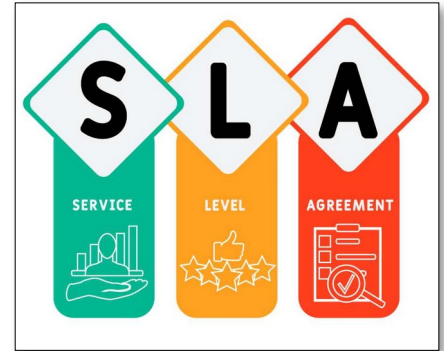
$$\text{UK} = \text{EU} - 34$$

Inhoud

- welk probleem willen we oplossen?
- kansen
- rekenen met kansen
- regel van Bayes

Welk probleem?

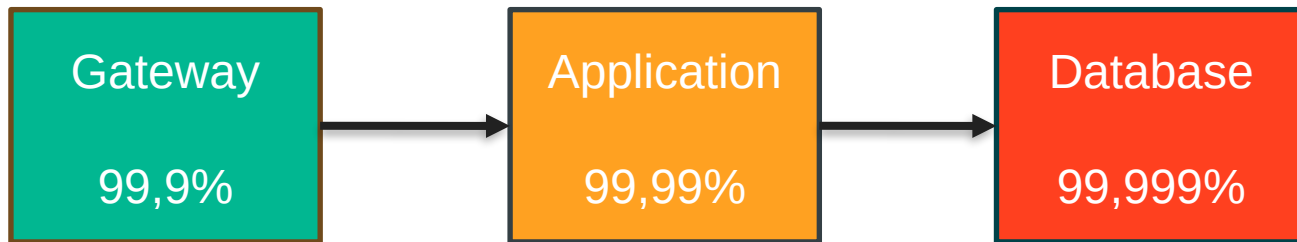
High availability server



- de beschikbaarheid van een server
 - reliability: de **kans** dat de server niet faalt
 - maintainability: de **kans** dat de server na falen succesvol hersteld wordt
 - availability: de **kans** dat een server op het gegeven moment niet faalt en niet wordt hersteld na een faling

High availability server

- de availability van een IT service hangt af van de availability van elk van de componenten
- Voorbeeld – Kubernetes cluster:

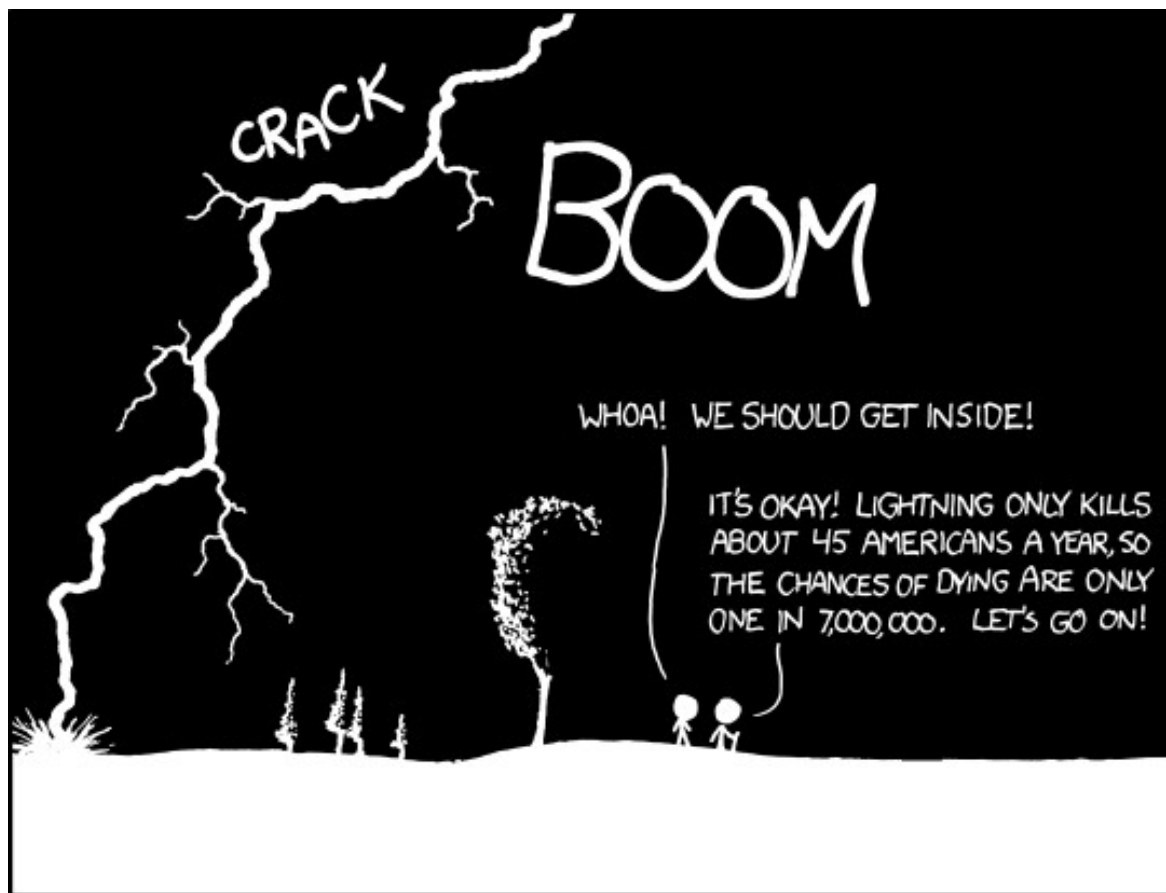


- $\text{availability cluster} = 0.999 \times 0.9999 \times 0.99999 = 0.9989$
- 99.89 % = 10 uur downtime / jaar

Kansen

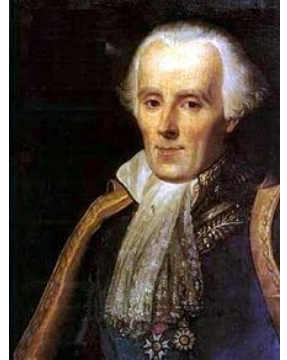
Kansen zijn raar...

- slaagkans
 - 2 mogelijkheden, dus ...
 - vorig jaar waren 70 van de 210 geslaagd, dus ...
 - voorbij 10 jaar was het slaagpercentage gemiddeld 25%, dus ...
 - als je KSO, ASO, TSO, BSO gevolgd hebt, is je slaagkans dan anders?
 - op het einde van het jaar is je slaagkans 100% of 0% (naargelang je resultaat)
- in dit hoofdstuk: enkel kansen in omstandigheden die duidelijk afgeleid zijn



THE ANNUAL DEATH RATE AMONG PEOPLE
WHO KNOW THAT STATISTIC IS ONE IN SIX.

Kansen volgens Laplace



- "afgeleijnd experiment"
- gegeven: verzameling mogelijke uitkomsten en verzameling "gunstige" uitkomsten
- gevraagd: wat is de kans op een gunstige uitkomst?
- voorbeeld: zak knikkers met 20 blauwe en 30 rode. Ik kies een knikker. Wat is de kans op een rode knikker?

Kansen volgens Laplace

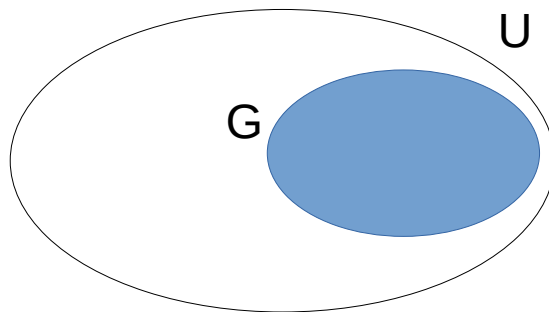
- $P(\text{gebeurtenis}) = \frac{\text{aantal gunstige uitkomsten}}{\text{aantal mogelijke uitkomsten}}$
 - notatie $P()$ voorbeeld: $P(\text{regen})$, $P(\text{ziek})$, ...
 - $0 \leq \text{kans} \leq 1$
 - $\text{kans} * 100 = \text{kans in percent}$ (zet er dan altijd % achter!)

Kansen volgens Laplace

- voorbeeld:
 - aantal studenten TI vorig jaar = 210
 - aantal geslaagde studenten TI vorig jaar = 70
 - wat is de kans dat een willekeurige student van vorig jaar geslaagd was?
 - wat is de kans dat een willekeurige student dit jaar geslaagd zal zijn?

Kansen volgens Laplace

- met verzamelingen

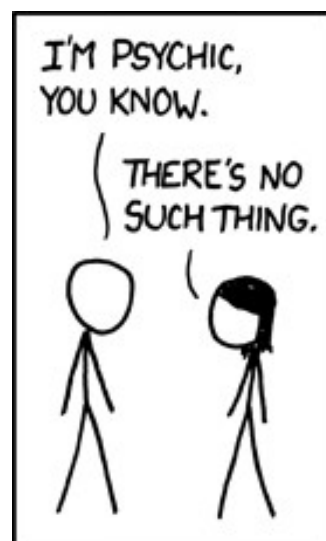


- $P(G) = \#G / \# U$

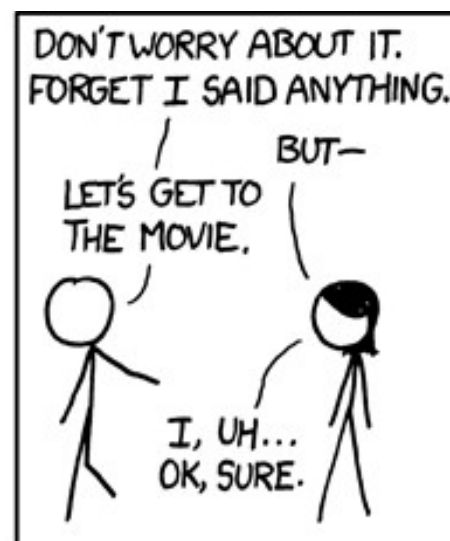
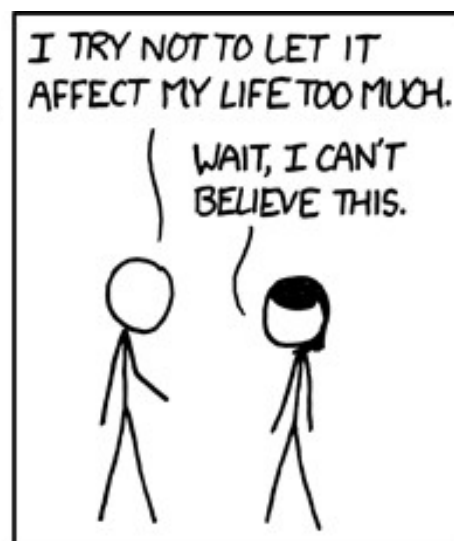
Voorbeeld

- twee dobbelstenen: wat is de kans om 7 te gooien?

Gebeurtenis A		#A
2	(1,1)	1
3	(1,2); (2,1)	2
4	(1,3); (2,2); (3,1)	3
5	(1,4); (2,3); (3,2); (4,1)	4
6	(1,5); (2,4); (3,3); (4,2); (5,1)	5
7	(1,6); (2,5); (3,4); (4,3); (5,2); (6,1)	6
8	(2,6); (3,5); (4,4); (5,3); (6,2)	5
9	(3,6); (4,5); (5,4); (6,3)	4
10	(4,6); (5,5); (6,4)	3
11	(5,6); (6,5)	2
12	(6,6)	1
#U (= TOTAAL)		36



HOLY SHIT!



THIS TRICK MAY ONLY WORK 1% OF THE TIME,
BUT WHEN IT DOES, IT'S TOTALLY WORTH IT.

Kansen als relatieve frequentie

- gegeven: zak met 20 blauwe en 30 rode knikkers
- gevraagd: als ik 1000 keer een knikker neem (en terug leg), wat is dan de verwachte relatieve frequentie?

Kans als relatieve frequentie

Steekproef: 10, 100 en 1000 keer:

steekproef lengte	rel. freq. blauw	rel. freq. rood
10	0,7	0,4
100	0,44	0,56
1000	0,392	0,608
10000	0,4019	0,5981



Kans is dus de relatieve frequentie van een (theoretische) oneindige steekproef

Voorwaardelijke kansen

- kansen kunnen afhangen van bepaalde kennis
- voorbeeld: wat is de kans dat een smartphone nog werkt als ik weet dat die in het water is gevallen?
- notatie: $P(\text{werkt} \mid \text{water})$
- als $P(\text{werkt}) = P(\text{werkt} \mid \text{water})$ dan zeggen we dat de gebeurtenissen “onafhankelijk” zijn

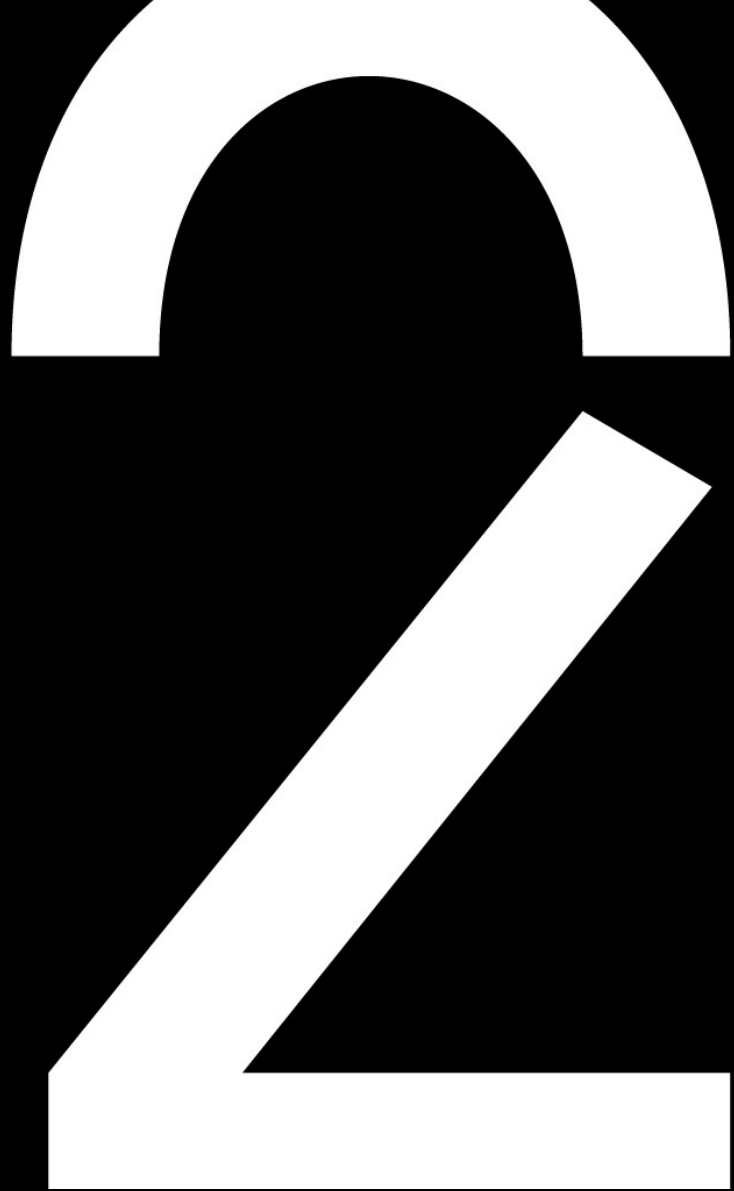
Kans als relatieve frequentie

Je kan voorwaardelijke kansen soms ook gemakkelijk uit een kruistabel aflezen:

	Wit merk (White label)	Geen wit merk (Private label)	
Slechte koeling	1498	1513	3011
Goede koeling	504	6485	6989
	2002	7998	10000

Wat is de kans dat een PC van een wit merk uit deze steekproef een slechte koeling heeft?

Rekenen met kansen



De inverse

- $P(\text{niet } G) = 1 - P(G)$
- bv: kans dat het regent is 55%, wat is de kans dat het niet regent?

$$P(\text{regent niet}) = 1 - P(\text{regent}) = 1 - 0,55$$

De productregel

- voorbeeld: groep kinderen, 40% meisjes, 60% jongens, 10% draagt een bril
- dus:
 - $P(\text{jongen}) =$
 - $P(\text{meisje}) =$
 - $P(\text{bril}) =$
- gevraagd: kans op een brildragend meisje
 - $P(\text{meisje EN bril}) = ?$

De productregel

- als er evenveel meisjes als jongens een bril dragen:

$$P(\text{meisje EN bril}) = P(\text{meisje}) * P(\text{bril})$$

- maar... wat als alle brildragers jongens zijn?

De productregel

- meer algemene formule:

$$P(\text{meisje EN bril}) = P(\text{meisje}) * P(\text{bril} \mid \text{meisje})$$

- wat is $P(\text{bril} \mid \text{meisje})$???
- stel dus dat 20% van de meisjes een bril draagt, dan ...
- opmerking: je kan het ook zo vinden (afhankelijk van de gegevens die je hebt):

$$P(\text{meisje EN bril}) = P(\text{bril}) * P(\text{meisje} \mid \text{bril})$$

De productregel

- met kruistabel: stel dat er 100 kinderen zijn

	jongens	meisjes	totaal
bril			10
geen bril			
totaal	60	40	100

20% van de meisjes draagt een bril...

De productregel

- met kruistabel: stel dat er 100 kinderen zijn

	jongens	meisjes	totaal
bril	2	8	10
geen bril	58	32	90
totaal	60	40	100

20% van de meisjes draagt een bril...

De productregel

- gebeurtenissen zijn "onafhankelijk" als $P(A|B)=P(A)$
- voorbeeld: kans dat iemand een rode T-shirt draagt en kans dat iemand goed is in wiskunde
- $P(\text{rood EN wiskunde})$
= $P(\text{rood}) * P(\text{wiskunde} \mid \text{rood})$
= $P(\text{rood}) * P(\text{wiskunde})$

De somregel

- voorbeeld: wat is de kans dat een willekeurig kind een meisje is OF een bril draagt (of allebei)?

	jongens	meisjes	totaal
bril	2	8	10
geen bril	58	32	90
totaal	60	40	100

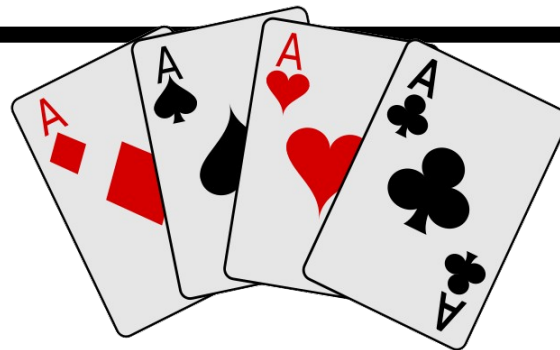
De somregel

- $P(\text{meisje OF bril})$
 $= P(\text{meisje}) + P(\text{bril}) - P(\text{meisje EN bril})$
 $= 0,4 + 0,1 - 0,08 = 0,42 = 42\%$

- algemeen:

$$P(A \text{ of } B) = P(A) + P(B) - P(A \text{ en } B)$$

Voorbeeld kaarten



- gegeven: boek kaarten, trek er 1 uit
- wat is de kans dat die kaart een aas is en een harten?

$P(\text{aas en harten})$

$= P(\text{aas}) * P(\text{harten} \mid \text{aas})$

$= 4/52 * 1/4 = 1/52$

- wat is de kans dat die kaart een aas is of een harten?

$P(\text{aas of harten})$

$= P(\text{aas}) + P(\text{harten}) - P(\text{aas en harten})$

$= 4/52 + 13/52 - 1/52 = 16/52$

Voorbeeld gaming

- 3 spelers proberen vijand uit te schakelen
 - speler1: $1/2$ kans om te raken
 - speler2: $1/3$ kans om te raken
 - speler3: $1/4$ kans om te raken
- wat is de kans dat de vijand overleeft?



Voorbeeld gaming

- kans op overleven = speler 1 mist EN speler 2 mist EN speler 3 mist
- spelers schieten 1 maal, onafhankelijk van elkaar
- $P(\text{overleven}) =$
 $P(\text{speler1 mist}) * P(\text{speler2 mist}) * P(\text{speler3 mist})$
- $P(\text{overleven}) = 1/2 * 2/3 * 3/4 = 0,25$
- dus...

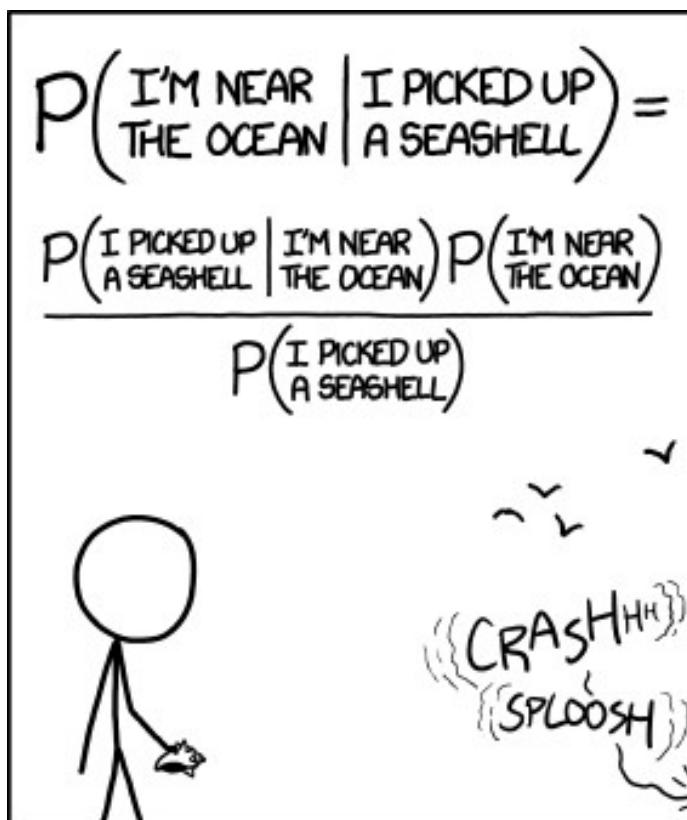


Regel van Bayes

De regel van Bayes

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

- voorbeeld:
 - kans dat student slaagt = 80%
 - kans dat student zich bezat als hij gebuisd is = 90%
 - kans dat student zich bezat = 60%
- we zien een zatte student: wat is de kans dat hij gebuisd is?



STATISTICALLY SPEAKING, IF YOU PICK UP A SEASHELL AND DON'T HOLD IT TO YOUR EAR, YOU CAN PROBABLY HEAR THE OCEAN.

Machine learning met Bayes

- vooral interessant bij analyseren van tekst
 - is een bepaalde tekst verdacht of niet?
 - is een bepaalde tekst spam of niet?
 - bevat een tekst een positieve of negatieve emotie?
 - ...

Voorbeeld: spam



- gegeven een tekst, wat is de kans dat deze spam is?

$P(\text{spam} \mid \text{woorden})$

$= P(\text{woorden} \mid \text{spam}) * P(\text{spam}) / P(\text{woorden})$

- gegeven: teksten uit het verleden met aanduiding of ze spam waren of niet

Voorbeeld spam

- we benaderen de kansen met relatieve frequenties (“waarschijnlijkheid” of “likelihood”):
 - $P(\text{spam}) \sim \# \text{spam} / \# \text{mails}$
 - $P(\text{woorden}) \sim P(\text{woord1}) * P(\text{woord2}) * \dots * P(\text{woordn})$
 - $P(\text{woorden} \mid \text{spam}) \sim P(\text{woord1} \mid \text{spam}) * P(\text{woord2} \mid \text{spam}) * \dots * P(\text{woordn} \mid \text{spam})$

Samenvatting



Samenvatting

- $P(\text{niet } A) = 1 - P(A)$
- $P(A \text{ en } B) = P(A) * P(B|A)$
- $P(A \text{ of } B) = P(A) + P(B) - P(A \text{ en } B)$
- $P(A|B) = P(B|A) * P(A) / P(B)$

Oefeningen

Oefeningen

- Klassikaal
 - 3 politici hebben resp. 30%, 20% en 50% kans dat ze verkozen worden (V_1 , V_2 , V_3)
 - de kans dat de politici de belastingen verlagen is resp. 50%, 40% en 30%
 - wat is de kans dat de belastingen verlaagd worden?

Oefeningen

$$\begin{aligned} P(B_{\text{verlaagd}}) &= P(V_1) \cdot P(B_{\text{verlaagd}} | V_1) \\ &+ P(V_2) \cdot P(B_{\text{verlaagd}} | V_2) \\ &+ P(V_3) \cdot P(B_{\text{verlaagd}} | V_3) \end{aligned}$$

Oefeningen

- Klassikaal
 - na de verkiezingen worden de belastingen verlaagd
 - hoeveel kans is er dat dat komt doordat politicus 3 verkozen werd?

Oefeningen

$$P(V_3|B_{\text{verlaagd}}) = \frac{P(V_3) \cdot P(B_{\text{verlaagd}}|V_3)}{P(B_{\text{verlaagd}})}$$

Oefeningen

- Zie Canvas