

Tentamen Statistische methoden
4052STAMEY
11 juli 2013, 9:00–12:00

Studienummers: Vult u alstublieft op het meerkeuzevragenformulier uw *Delftse* studienummer in (tbv automatische verwerking); en op het open vragen formulier graag **beide**, naar volgend voorbeeld: 1234567(D), 7654321(L).

Bij dit examen is het gebruik van een (evt. grafische) rekenmachine toegestaan. Tevens krijgt u een formuleblad uitgereikt—na afloop inleveren alstublieft. Normering: De meerkeuzevragen tellen voor één derde en de open vragen voor twee derde van het cijfer. Bij de open vragen telt elk (vraag)onderdeel even zwaar.

Meerkeuzevragen

Toelichting: In het algemeen zijn niet altijd vijf van de zes alternatieven 100% fout, het juiste antwoord is het meest volledige antwoord. Maak op het bijgeleverde antwoordformulier het hokje behorende bij het door u gekozen alternatief zwart of blauw. Doorstrepen van een fout antwoord heeft geen zin: u moet het òf uitgummen, òf verwijderen met correctievloeistof òf een nieuw formulier invullen. Vergeet niet uw studienummer in te vullen èn aan te strepen.

1. Gegeven zijn $P(A) = \frac{1}{2}$, $P(A|B) = \frac{4}{5}$ en $P(A|B^c) = \frac{2}{5}$. Dan geldt $P(B) =$
a. $\frac{1}{4}$ b. $\frac{3}{10}$ c. $\frac{1}{3}$ d. $\frac{3}{8}$ e. $\frac{2}{5}$ f. $\frac{6}{5}$
2. Groep 1 bestaat uit 20 mannen en 30 vrouwen. Groep 2 telt 45 mannen en 15 vrouwen. We kiezen een willekeurig persoon uit groep 1 en vervolgens een willekeurig persoon uit groep 2, in de hoop een paar te vormen van twee personen van hetzelfde geslacht. Wanneer het geslacht van beide personen verschillend is, sturen we ze weg en herhalen de procedure; beide groepen zijn nu elk 1 kleiner geworden. We tellen het aantal malen X dat we de procedure moeten herhalen totdat we een paar van hetzelfde geslacht gekozen hebben. Dan heeft X een
a. $Bin(2, \frac{11}{20})$ verdeling b. $Geo(\frac{11}{20})$ verdeling
c. $Bin(45, \frac{11}{20})$ verdeling d. discrete verdeling met uitkomsten $1, 2, \dots, 46$
e. $Bin(50, \frac{11}{20})$ verdeling f. discrete verdeling met uitkomsten $1, 2, \dots, 50$
3. X en Y zijn onafhankelijke stochasten met kansverdelingen

$$\begin{array}{c|cc} x & 0 & 1 \\ \hline P(X=x) & 1/2 & 1/2 \end{array} \quad \text{resp.} \quad \begin{array}{c|ccc} y & -1 & 0 & 1 \\ \hline P(Y=y) & 1/4 & 1/4 & 1/2 \end{array}$$

De kansverdeling van $M = \max\{X, Y\}$ wordt dan gegeven door

- | | |
|---|---|
| a. $\begin{array}{c cc} m & 0 & 1 \\ \hline P(M=m) & 1/4 & 3/4 \end{array}$ | b. $\begin{array}{c ccc} m & -1 & 0 & 1 \\ \hline P(M=m) & 1/4 & 1/2 & 1/4 \end{array}$ |
| c. $\begin{array}{c ccc} m & -1 & 0 & 1 \\ \hline P(M=m) & 1/4 & 1/4 & 1/2 \end{array}$ | d. $\begin{array}{c ccc} m & -1 & 0 & 1 \\ \hline P(M=m) & 1/2 & 1/4 & 1/4 \end{array}$ |
| e. $\begin{array}{c cc} m & 0 & 1 \\ \hline P(M=m) & 1/3 & 2/3 \end{array}$ | f. $\begin{array}{c cc} m & 0 & 1 \\ \hline P(M=m) & 2/3 & 1/3 \end{array}$ |

4. Stel W, X, Y en Z zijn onafhankelijke standaard normaal verdeelde stochasten. Dan is

$$P(W + X + Y + Z \geq 1)$$

gelijk aan:

- a. 0.5000 b. 0.4013 c. 0.3446 d. 0.3085 e. 0.1587 f. 0.0228

5. Laat X en Y twee stochasten zijn met $\text{Var}(X) = 5$ en $\text{Var}(Y) = 3$ en correlatie $+1$. Dan is de variantie van de stochast $3X - 2Y - 1$ ongeveer:
- a. 8 b. 9 c. 11 d. 33 e. 57 f. 103

6. De stochasten X_1, X_2, \dots, X_{50} zijn onafhankelijke $\text{Exp}(5)$ -verdeelde stochasten. Met behulp van de *Centrale Limietstelling* zien we dat de kans $P(X_1 + X_2 + \dots + X_{50} > 12)$ bij benadering gelijk is aan:
- a. 0.01 b. 0.25 c. 0.04 d. 0.08 e. 0.11 f. 0.20

7. Beschouw de dataset: 1 2 2 3 3 4 4 4 4 4.
We trekken onafhankelijk van elkaar en met teruglegging 10 keer uit deze dataset. Wat is de kans dat de dataset die u zo krijgt precies zes keer een 4 bevat?
- a. 0 b. 0.205 c. 0.324 d. 0.315 e. 0.4 f. 0.5

8. De MAD van de onderstaande *Challenger* dataset

53 57 58 63 66 67 67 67 68 69 70 70
70 70 72 73 75 75 76 76 78 79 81

wordt gegeven door

- a. 3 b. 4 c. 7 d. 13 e. 28 f. 70

9. Jan koopt hooi van een bepaald merk voor zijn konijnen. Hij denkt echter dat er gemiddeld te weinig hooi in de verpakking zit. Hij meet het gewicht van het hooi in 9 zakken precies op, en hij vindt $\bar{x}_n = 0.977$ kg en een steekproefstandaarddeviatie $s_n = 0.030$ kg. Op de verpakking staat dat er gemiddeld 1 kg in een zak zit. We modelleren het gewicht van het hooi als een $N(\mu, \sigma^2)$ verdeling. Als nul-hypothese nemen we $H_0 : \mu = 1$, en als alternatief $H_1 : \mu < 1$. Verder kiezen we een significantieniveau van 5%. Welke van de volgende is de uitkomst voor p en welke conclusie hoort daarbij?
- a. $0.01 < p < 0.025$, verwerp H_0 . b. $0.01 < p < 0.025$, verwerp H_0 niet.
c. $0.025 < p < 0.05$, verwerp H_0 . d. $0.025 < p < 0.05$, verwerp H_0 niet.
e. $0.05 < p < 0.10$, verwerp H_0 . f. $0.05 < p < 0.10$, verwerp H_0 niet.

10. Een steekproef van omvang 5 uit een normale verdeling met onbekende μ en onbekende σ resulteert in de volgende dataset:

129.0 129.4 130.8 131.1 132.2

Het steekproefgemiddelde is 130.50 en de steekproefstandaarddeviatie is 1.3038.

Een symmetrisch tweezijdig 90%-betrouwbaarheidsinterval voor μ wordt gegeven door:

- a. $129.11 < \mu < 131.89$ b. $129.26 < \mu < 131.74$ c. $129.33 < \mu < 131.67$
d. $129.54 < \mu < 131.46$ e. $129.61 < \mu < 131.39$ f. $129.75 < \mu < 131.25$

Open vragen

Toelichting: Een antwoord alleen is *niet* voldoende: er dient een berekening, toelichting en/of motivatie aanwezig te zijn. Dit alles goed leesbaar en in goed Nederlands.

1. Er zijn uitstekende tests voor zwangere vrouwen op het syndroom van Down, maar onfeilbaar zijn ze niet. Ten eerste is bekend dat 1 procent van de embryo's lijdt aan het syndroom. Als een baby het syndroom heeft, dan is er 90 procent kans dat de testuitslag positief is, maar bij een gezonde baby is er 1 procent kans dat de testuitslag toch positief is. Bereken de kans dat een vrouw die positief test een baby met het Downsyndroom in zich draagt.

N.B. Een positieve testuitslag betekent: volgens de test is er sprake van het Downsyndroom.

2. De continue stochast X heeft als verdelingsfunctie F , met

$$F(x) = \begin{cases} 0 & \text{als } x \leq 0 \\ \sqrt{x} & \text{als } 0 < x < 1 \\ 1 & x \geq 1. \end{cases}$$

- a. Bepaal $P\left(\frac{1}{4} \leq X \leq \frac{9}{16}\right)$.
 - b. Bepaal $E\left[3\sqrt{X} + 5\right]$.
 - c. Bepaal de verdeling van \sqrt{X} .
3. De hoeken van een gelijkbenige driehoek zijn θ , θ en γ . Ze worden elk apart gemeten (in radialen); dus $\theta + \theta + \gamma = \pi$ (rad). De metingen van de drie hoeken zijn X_1 , X_2 en X_3 . Gegeven is dat X_1 , X_2 en X_3 onafhankelijk zijn en zuiver zijn voor respectievelijk θ , θ en γ , met variantie σ^2 . We definiëren de volgende schatters voor θ :

$$S = \frac{1}{2}(X_1 + X_2) \quad \text{en} \quad T = \frac{1}{6}(2\pi + X_1 + X_2 - 2X_3)$$

- a. Ga voor elk van de schatters na of hij zuiver is voor θ .
 - b. Bereken de variantie voor beide schatters en geef aan welke van de twee schatters u zou prefereren.
4. Men beschikt over data die een realisatie vormen van een steekproef X_1, X_2, \dots, X_{25} uit een $N(\mu, 1)$ verdeling. Men ontwerpt een toets voor $H_0 : \mu = 0$ tegen $H_1 : \mu > 0$ met toetsingsgrootte \bar{X}_n , waar als beslissingsregel uitkomt: verwerp H_0 ten gunste van H_1 als voor de data geldt dat $\bar{x}_n \geq 0.3$.
 - a. Wat is het significantieniveau van de toets?
 - b. Wat is de kans op een type II fout als in werkelijkheid $\mu = 0.55$?
 - c. Is het mogelijk de steekproefgrootte n (in plaats van 25) zo te kiezen dat zowel het significantieniveau als de onder 4b bedoelde type II fout beide kleiner zijn dan 0.01? Indien ja, bepaal de kleinste mogelijke n ; zo nee, leg uit waarom het niet mogelijk is.

5. Schets het principe van de bootstrap en leg heel precies het verschil uit tussen de parametrische en de empirische bootstrap.

Antwoorden multiple choice:

1 a. Via $P(A) = P(B)P(A|B) + P(B^c)P(A|B^c) = P(B)P(A|B) + (1 - P(B))P(A|B^c)$:
 $\frac{1}{2} = \frac{4}{5}P(B) + \frac{2}{5}(1 - P(B)) = \frac{2}{5} + \frac{2}{5}P(B)$, dus $P(B) = (\frac{1}{2} - \frac{2}{5})/\frac{2}{5} = \frac{1}{4}$.

2 d. Nadat we de eerste keer een tweetal personen uit de groepen gekozen hebben is de samenstelling veranderd, want ze gaan niet weer terug. De kans op een succesvolle keuze is de tweede keer anders dan de eerste keer en hangt bovendien ook nog af van de eerste uitkomst. Om deze reden is zowel de binomiale als de geometrische verdeling uit te sluiten, want die zijn gebaseerd op een rij *onafhankelijke* 'experimenten' met 'constante' succeskans. Zo blijven d. en f. over en de vraag is dus hoe lang we maximaal bezig kunnen zijn. Daarvoor moet telkens bij een man uit groep 1 een vrouw uit groep 2 worden gekozen en omgekeerd. Voor de man-vrouw combinaties kan dit hoogstens 15 keer, want dan zijn de vrouwen in groep 2 op. Bij het omgekeerde is het aantal vrouwen in groep 1 de limiet: 30. Het kan dus maximaal $15 + 30 = 45$ maal fout gaan. Er zijn dan in beide groepen alleen mannen over, dus de 46ste keer is het dan raak.

3 a. Er geldt $P(M = 0) = P(X = 0, Y = -1) + P(X = 0, Y = 0) = P(X = 0)P(Y = -1) + P(X = 0)P(Y = 0) = 1/8 + 1/8 = 1/4$. Omdat M alleen de waarde 0 of 1 kan aannemen volgt direct dat $P(M = 1) = 1 - P(M = 0) = \frac{3}{4}$.

4 d. $S = W + X + Y + Z$ is normaal verdeeld met parameters $\mu = 0$ en $\sigma^2 = 1 + 1 + 1 + 1 = 4$. Dus $P(S \geq 1) = P(S/2 \geq 0.5) = 0.3085$.

5 c. Voor de variantie maakt de verschuiving over 1 niet uit: $\text{Var}(3X - 2Y - 1) = \text{Var}(3X - 2Y)$. Verder is

$$\text{Var}(3X - 2Y) = \text{Var}(3X) + \text{Var}(2Y) - 2\text{Cov}(3X, 2Y).$$

Omdat X en Y correlatie 1 hebben, geldt er dat $\text{Cov}(X, Y) = \rho \cdot \sigma_x \cdot \sigma_y = \sqrt{5 \cdot 3}$, en dus dat $\text{Cov}(3X, 2Y) = 6\text{Cov}(X, Y) = 6\sqrt{15} = 23.24$. Alles tezamen vinden we:

$$\text{Var}(3X - 2Y - 1) = 3^2\text{Var}(X) + 2^2\text{Var}(Y) + 12\text{Cov}(X, Y) = 9 \cdot 5 + 4 \cdot 3 + 12\sqrt{15} \approx 45 + 12 - 46.48 = 10.52.$$

6 d. Uit het gegeven dat elke X_i een $Exp(5)$ -verdeling heeft volgt direct dat $\mu = E[X_i] = 1/5$ en $\sigma^2 = \text{Var}(X_i) = 1/25$. Voor de som $S = X_1 + X_2 + \dots + X_{50}$ geldt dus $E[S] = 50 \cdot 1/5 = 10$ en $\text{Var}(S) = 50 \cdot 1/25 = 2$. Er geldt dan bij benadering

$$P(X_1 + X_2 + \dots + X_{50} > 12) = P(S > 12) \approx P\left(Z > \frac{12 - 10}{\sqrt{2}}\right) = P(Z > \sqrt{2}) \approx 0.0793.$$

7 b. Deze kans is $\binom{10}{6}(1/2)^{10} = 210/1024 = 0.205$.

8 b. De steekproefmediaan is het 12de getal in volgorde naar grootte: 70. De absolute afwijkingen t.o.v. 70 zijn

17	13	12	7	4	3	3	3	2	1	0	0
0	0	2	3	5	5	6	6	8	9	11	

of in volgorde naar grootte

0	0	0	0	1	2	2	3	3	3	3	4	5	5	6	6	7	8	9	11	12	13	17
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	----	----	----

De steekproefmediaan van deze absolute verschillen is 4.

9 c. De toetsingsgrootheid wordt gegeven door

$$t = \frac{\bar{x}_n - 1}{s_n/\sqrt{n}} = -2.3.$$

Er geldt dat $t_{8,0.05} < |t| < t_{8,0.025}$ (merk op dat T een t -verdeling heeft met $n - 1 = 8$ vrijheidsgraden!). Aangezien alleen nog negatievere waarden van T wijzen op de alternatieve hypothese H_1 , geldt dus dat $0.025 < p < 0.05$. Aangezien p kleiner is dan het significantieniveau, verwerpen we H_0 .

10 b. De grenzen van het betrouwbaarheidsinterval worden gegeven door $\bar{x}_n \pm t_{n-1, \alpha/2} s_n / \sqrt{n}$, met $\alpha = 0.1$, en $n = 5$. Dit geeft $130.50 \pm 2.132 \cdot 1.3038 / \sqrt{5}$.

Antwoorden open vragen:

1 Voer in: D : een baby heeft Down, en T : de testuitslag is positief. De gegevens zijn: $P(D) = 0.01$, $P(T | D) = 0.9$ en $P(T | D^c) = 0.01$. Gevraagd wordt $P(D | T) = \frac{P(D \cap T)}{P(T)}$. $P(D \cap T) = P(D) P(T | D) = 0.009$, en $P(T) = P(D \cap T) + P(D^c \cap T) = 0.009 + P(D^c) P(T | D^c) = 0.009 + 0.99 \cdot 0.01 = 0.0189$. Op elkaar delen geeft het antwoord: $P(D | T) = 0.009/0.0189 = 0.4762$.

2a Omdat X een continue stochast is geldt dat

$$P\left(\frac{1}{4} \leq X \leq \frac{9}{16}\right) = F\left(\frac{9}{16}\right) - F\left(\frac{1}{4}\right) = \sqrt{\frac{9}{16}} - \sqrt{\frac{1}{4}} = \frac{3}{4} - \frac{1}{2} = \frac{1}{4}.$$

2b Omdat de kansdichtheid $f(x)$ van X gelijk is aan $F'(x)$, geldt er dat:

$$f(x) = \begin{cases} \frac{1}{2\sqrt{x}} & \text{als } 0 < x < 1 \\ 0 & \text{als } x \leq 0 \text{ of } x \geq 1. \end{cases}$$

Maar dan geldt er wegens de ‘‘Change of variable formula’’ dat

$$E[\sqrt{X}] = \int_{-\infty}^{\infty} \sqrt{x} f(x) dx = \int_0^1 \sqrt{x} \frac{1}{2\sqrt{x}} dx = \int_0^1 \frac{1}{2} dx = \frac{1}{2}.$$

2c Schrijf $Y = \sqrt{X}$. Omdat $0 \leq X \leq 1$ gelden dezelfde grenzen voor Y . Voor $0 \leq y \leq 1$ geldt dus:

$$P(Y \leq y) = P(\sqrt{X} \leq y) = P(X \leq y^2) = \sqrt{y^2} = y.$$

Dit betekent dat Y een $U(0, 1)$ -verdeling heeft.

3a Omdat X_1 en X_2 beide zuiver zijn voor θ , is $E[X_1] = \theta$ en $E[X_2] = \theta$. Dit betekent dat

$$E[S] = \frac{1}{2}(E[X_1] + E[X_2]) = \frac{1}{2}(\theta + \theta) = \theta.$$

Dus S is zuiver voor θ . Op een zelfde manier, gebruikmakend van het feit dat $\pi = 2\theta + \gamma$, is

$$\begin{aligned} E[T] &= \frac{1}{6}(2\pi + E[X_1] + E[X_2] - 2E[X_3]) \\ &= \frac{1}{6}(2\pi + \theta + \theta - 2\gamma) \\ &= \frac{1}{6}(2(2\theta + \gamma) + \theta + \theta - 2\gamma) \\ &= \theta. \end{aligned}$$

Dus ook T is zuiver voor θ .

3b Voor S krijgen we (vanwege de onafhankelijkheid van X_1 en X_2)

$$\text{Var}(S) = \text{Var}\left(\frac{1}{2}(X_1 + X_2)\right) = \frac{1}{4}[\text{Var}(X_1) + \text{Var}(X_2)] = \frac{1}{4}(\sigma^2 + \sigma^2) = \frac{\sigma^2}{2}.$$

Voor T krijgen we

$$\text{Var}(T) = \text{Var}\left(\frac{1}{6}(2\pi + X_1 + X_2 - 2X_3)\right) = \frac{1}{36}[\text{Var}(X_1) + \text{Var}(X_2) + 4\text{Var}(X_3)] = \frac{1}{36}(\sigma^2 + \sigma^2 + 4\sigma^2) = \frac{\sigma^2}{6}.$$

Omdat beide schatters zuiver zijn geven we de voorkeur aan degene met de kleinste variantie, dus aan T .

4a Merk op dat \bar{X}_{25} een $N(\mu, 1/25)$ verdeling heeft. De kans op een type I fout is daarom gelijk aan

$$P(\bar{X}_n \geq 0.3 \mid \mu = 0) = P\left(\frac{\bar{X}_n}{1/5} \geq \frac{0.3}{1/5}\right) = P(Z \geq 1.5) \approx 0.0668.$$

4b Dit is

$$P(\bar{X}_n < 0.3 \mid \mu = 0.55) = P\left(\frac{\bar{X}_n - 0.55}{1/5} < \frac{0.3 - 0.55}{1/5}\right) = P(Z < -1.25) = P(Z > 1.25) \approx 0.1056.$$

4c Wanneer we de berekeningen bij de vorige onderdelen overdoen voor steekproefgrootte n , waarvoor geldt dat \bar{X}_n een $N(\mu, 1/n)$ verdeling heeft, dan vinden we voor de type I fout $P(Z \geq 0.3 \cdot \sqrt{n})$ en voor de type II fout $P(Z \leq -0.25 \cdot \sqrt{n}) \equiv P(Z \geq 0.25 \cdot \sqrt{n})$. De tweede kans is altijd groter; het is dus voldoende te zorgen dat die onder de 0.01 blijft. In de tabel vinden we $z_{0.01} = 2.326$, dus er moet gelden: $0.25 \cdot \sqrt{n} \geq 2.326$ ofwel $n \geq 86.56$. Derhalve: $n = 87$ is de kleinste waarde die aan de eisen voldoet.

5 Dit refereert aan hoofdstuk 18 in het boek van Dekking et al (2005). Het principe van de bootstrap (pagina 270) is toe te passen op elke steekproefgrootte $h(X_1, X_2, \dots, X_n)$, afgeleid van een *random sample* X_1, X_2, \dots, X_n :

1. Bepaal uit de gegeven data x_1, x_2, \dots, x_n een schatting \hat{F} van de verdelingsfunctie F .
2. Vervang het random sample X_1, X_2, \dots, X_n door het (bootstrap) random sample $X_1^*, X_2^*, \dots, X_n^*$ dat wordt getrokken uit \hat{F} .
3. Benader de verdeling van $h(X_1, X_2, \dots, X_n)$ door die van $h(X_1^*, X_2^*, \dots, X_n^*)$.

In de praktijk wordt het principe vaak toegepast door een groot aantal trekkingen van $h(X_1^*, X_2^*, \dots, X_n^*)$ met behulp van simulatie uit te voeren, en hiermee de verdeling van $h(X_1, X_2, \dots, X_n)$ te schatten. Soms is deze stap overbodig omdat de verdeling van $h(X_1^*, X_2^*, \dots, X_n^*)$ uit \hat{F} te bepalen is.

Het verschil tussen de parametrische en empirische bootstrap wordt bepaald door de wijze waarop de schatting van de verdelingsfunctie, \hat{F} , tot stand komt. Bij de parametrische bootstrap wordt de verdeling of de verdelingsfunctie gegeven met eventueel meerdere parameters. Deze parameters moeten uit de bestaande steekproef geschat worden. Bij de empirische bootstrap is deze informatie niet beschikbaar en maken we gebruik van de empirische verdelingsfunctie direct afgeleid uit de data. Dit laatste geval is equivalent aan het trekken van een bootstrap monster uit de bestaande data, met teruglegging.