
Verslag Python Challenge 2

Aiko Decaluwe, Fien Dewit, Viktor van Nieuwenhuize (Groep 4)

Opleiding: Tweede Bachelor Fysica en Sterrenkunde

Vak: Statistiek en Gegevensverwerking

Datum: 1 December 2020

1 Inleiding

De gegeven dataset bestaat uit waarden die uit een Gammaverdeling met een bepaalde k en θ zijn getrokken. Deze Gammaverdeling is gedefinieerd als:

$$f(y; k, \theta) = \frac{y^{k-1} e^{-y/\theta}}{\theta^k \Gamma(k)}$$

Om de schatters van k en θ te bepalen bestaan er verschillende methoden. We vergelijken de methode van de momenten (MM) en de Maximum Log Likelihood methode (MLLH). Verder zullen we ook bootstrappen en deze waarden in functie van het aantal getrokken waarden (N) uitzetten. Ten slotte wordt de gebootstrapte waarde voor N_{max} vergeleken met de geschatte waarden via de MM- en MLLH-schatters.

2 Berekening schatters en variantie

2.1 Methode van de momenten

We maken gebruik van de methode van de momenten om de schatters voor k en θ te bepalen in functie van de variabelen \bar{x} en $\overline{x^2}$. We berekenen dus de verwachtingswaarde van de gammaverdeling voor \bar{x} en $\overline{x^2}$:

$$\begin{aligned} M_1 = \bar{x} = \langle x \rangle &= \int_0^{+\infty} x \frac{x^{\hat{k}-1} e^{-x/\hat{\theta}}}{\hat{\theta}^{\hat{k}} \Gamma(\hat{k})} dx \\ &= \frac{1}{\hat{\theta}^{\hat{k}} \Gamma(\hat{k})} \int_0^{+\infty} x^{\hat{k}} \cdot (-\hat{\theta}) d e^{-x/\hat{\theta}} \\ &= \frac{\hat{\theta}}{\hat{\theta}^{\hat{k}} \Gamma(\hat{k})} \left(- \left[e^{-x/\hat{\theta}} x^{\hat{k}} \right]_0^{+\infty} + \int_0^{+\infty} e^{-x/\hat{\theta}} d x^{\hat{k}} \right) \\ &= \hat{k} \hat{\theta} \int_0^{+\infty} \frac{x^{\hat{k}-1} e^{-x/\hat{\theta}}}{\hat{\theta}^{\hat{k}} \Gamma(\hat{k})} dx \\ &= \hat{k} \hat{\theta} \\ M_2 = \overline{x^2} = \langle x^2 \rangle &= \int_0^{+\infty} x^2 \frac{x^{\hat{k}-1} e^{-x/\hat{\theta}}}{\hat{\theta}^{\hat{k}} \Gamma(\hat{k})} dx \\ &= \frac{1}{\hat{\theta}^{\hat{k}} \Gamma(\hat{k})} \int_0^{+\infty} x^{\hat{k}+1} \cdot (-\hat{\theta}) d e^{-x/\hat{\theta}} \\ &= \frac{\hat{\theta}}{\hat{\theta}^{\hat{k}} \Gamma(\hat{k})} \left(- \left[e^{-x/\hat{\theta}} x^{\hat{k}+1} \right]_0^{+\infty} + \int_0^{+\infty} (k+1) e^{-x/\hat{\theta}} x^{\hat{k}} dx \right) \\ &= \hat{\theta}(\hat{k}+1) \cdot M_1 \\ &= \hat{\theta}^2 \hat{k}(\hat{k}+1) \end{aligned}$$

Uit de integralen van de verwachtingswaarden voor \bar{x} en $\overline{x^2}$ bekomen we 2 vergelijkingen:

$$\begin{cases} \bar{x} &= \hat{k} \hat{\theta} \\ \overline{x^2} &= \hat{\theta}^2 \hat{k}(\hat{k}+1) \end{cases}$$

Deze kunnen we oplossen naar een schatter voor k en een schatter voor θ .

$$\hat{k} = \frac{\bar{x}^2}{x^2 - \bar{x}^2}$$

$$\hat{\theta} = \frac{\bar{x}^2 - \bar{x}^2}{\bar{x}}$$

De variantie van deze schatters wordt gegeven door de formule:

$$\text{cov}(\hat{\theta}_i, \hat{\theta}_j) = \sum_{l,m} \frac{\partial \hat{\theta}_i}{\partial \hat{e}_l} \frac{\partial \hat{\theta}_j}{\partial \hat{e}_m} \text{cov}(\hat{e}_l, \hat{e}_m)$$

Hierin worden de verschillende afgeleiden gegeven door:

$$\frac{\partial \hat{\theta}}{\partial \bar{x}} = \frac{-\bar{x}^2 - \bar{x}^2}{\bar{x}^2} \quad \frac{\partial \hat{\theta}}{\partial x^2} = \frac{1}{\bar{x}} \quad \frac{\partial \hat{k}}{\partial \bar{x}} = \frac{2\bar{x}\bar{x}^2}{(\bar{x}^2 - \bar{x}^2)^2} \quad \frac{\partial \hat{k}}{\partial x^2} = \frac{-\bar{x}^2}{(\bar{x}^2 - \bar{x}^2)^2}$$

De uitdrukkingen voor de covariantie van de schatters werden niet manueel uitgerekend. De vergelijkingen voor de afgeleiden van de schatters \hat{k} en $\hat{\theta}$ naar de momenten \bar{x} en x^2 werden in python geschreven en met deze uitkomsten vormden we de uitdrukking voor de covarianties. De uitdrukkingen voor de covariantie van de momenten werden ook apart uitgerekend.

$$\begin{aligned} \text{cov}(\hat{k}, \hat{k}) &= \left(\frac{\partial \hat{k}}{\partial \bar{x}} \right)^2 \text{cov}(\bar{x}, \bar{x}) + \left(\frac{\partial \hat{k}}{\partial x^2} \right)^2 \text{cov}(\bar{x}^2, \bar{x}^2) + 2 \left(\frac{\partial \hat{k}}{\partial \bar{x}} \right) \left(\frac{\partial \hat{k}}{\partial x^2} \right) \text{cov}(\bar{x}, \bar{x}^2) \\ \text{cov}(\hat{\theta}, \hat{\theta}) &= \left(\frac{\partial \hat{\theta}}{\partial \bar{x}} \right)^2 \text{cov}(\bar{x}, \bar{x}) + \left(\frac{\partial \hat{\theta}}{\partial x^2} \right)^2 \text{cov}(\bar{x}^2, \bar{x}^2) + 2 \left(\frac{\partial \hat{\theta}}{\partial \bar{x}} \right) \left(\frac{\partial \hat{\theta}}{\partial x^2} \right) \text{cov}(\bar{x}, \bar{x}^2) \\ \text{cov}(\hat{k}, \hat{\theta}) &= \left(\frac{\partial \hat{k}}{\partial \bar{x}} \right) \left(\frac{\partial \hat{\theta}}{\partial \bar{x}} \right) \text{cov}(\bar{x}, \bar{x}) + \left(\frac{\partial \hat{k}}{\partial x^2} \right) \left(\frac{\partial \hat{\theta}}{\partial x^2} \right) \text{cov}(\bar{x}^2, \bar{x}^2) + \left(\frac{\partial \hat{k}}{\partial \bar{x}} \right) \left(\frac{\partial \hat{\theta}}{\partial x^2} \right) \text{cov}(\bar{x}, \bar{x}^2) \\ &\quad + \left(\frac{\partial \hat{k}}{\partial x^2} \right) \left(\frac{\partial \hat{\theta}}{\partial \bar{x}} \right) \text{cov}(\bar{x}, \bar{x}^2) \end{aligned}$$

De covarianties van de momenten \bar{x} en \bar{x}^2 kan berekend worden via:

$$\text{cov}(\hat{e}_l, \hat{e}_m) = \frac{1}{N(N-1)} \sum_{i=1}^N (g_l(x_i) - \bar{g}_l)(g_m(x_i) - \bar{g}_m)$$

Waaruit dan volgt dat:

$$\begin{aligned} \text{cov}(\bar{x}, \bar{x}^2) &= \frac{1}{N(N-1)} \sum_{i=1}^N (x_i - \bar{x})(x_i^2 - \bar{x}^2) \\ \text{cov}(\bar{x}, \bar{x}) &= \frac{1}{N(N-1)} \sum_{i=1}^N (x_i - \bar{x})^2 \\ \text{cov}(\bar{x}^2, \bar{x}^2) &= \frac{1}{N(N-1)} \sum_{i=1}^N (x_i^2 - \bar{x}^2)^2 \end{aligned}$$

Vervolgens kunnen we de correlatie tussen $\hat{\theta}$ en \hat{k} bereken met behulp van de formule voor de variantie. Ook dit werd niet volledig met de hand uitgerekend, maar volgt gewoon uit de vorige uitkomsten.

$$\rho_{\hat{\theta}\hat{k}} = \frac{\text{cov}(\hat{\theta}, \hat{k})}{\sqrt{\text{cov}(\hat{\theta}, \hat{\theta}) \text{cov}(\hat{k}, \hat{k})}}$$

2.2 Maximum Log Likelihood Methode

De likelihoodvergelijking van de Gammaverdeling is:

$$\mathcal{L}(y; k, \theta) = \prod_{i=1}^N \frac{y_i^{k-1} e^{-y/\theta}}{\theta^k \Gamma(k)}$$

Indien men hiervan het natuurlijk logaritme neemt, verkrijgt men het volgende:

$$\begin{aligned} \ln(\mathcal{L}(y; k, \theta)) &= \sum_i^N \ln\left(\frac{y_i^{k-1} e^{-y/\theta}}{\theta^k \Gamma(k)}\right) \\ &= \sum_i^N \left[\ln(y_i^{k-1}) - \frac{y_i}{\theta} - \ln(\theta^k) - \ln(\Gamma(k)) \right] \end{aligned}$$

Om vergelijkingen voor de schatters \hat{k} en $\hat{\theta}$ te vinden berekenen we het maximum van deze log likelihood. Dit doen we door de partiële afgeleide naar de schatters te nemen en dit gelijk te stellen aan 0.

1) Afleiden naar k

$$\begin{aligned} \left. \frac{\partial}{\partial k} \left(\ln(\mathcal{L}(y; k, \theta)) \right) \right|_{k=\hat{k}} &= \sum_i^N \left[\ln(y_i) - \ln(\theta) - \frac{\Gamma'(k)}{\Gamma(k)} \right] \Big|_{k=\hat{k}} \\ \iff 0 &= \sum_i^N \left[\ln(y_i) - \ln(\hat{\theta}) - \frac{\Gamma'(\hat{k})}{\Gamma(\hat{k})} \right] \\ \iff 0 &= \overline{\ln(y)} - \ln(\hat{\theta}) - \psi^{(0)}(\hat{k}) \end{aligned}$$

2) Afleiden naar θ

$$\begin{aligned} \left. \frac{\partial}{\partial \theta} \left(\ln(\mathcal{L}(y; k, \theta)) \right) \right|_{\theta=\hat{\theta}} &= \sum_i^N \left[\frac{y_i}{\theta^2} - \frac{k\theta^{k-1}}{\theta^k} \right] \Big|_{\theta=\hat{\theta}} \\ \iff 0 &= \sum_i^N \left[\frac{y_i}{\hat{\theta}^2} - \frac{\hat{k}}{\hat{\theta}} \right] \\ \iff \frac{1}{N} \sum_i^N y_i &= \hat{k} \hat{\theta} \\ \iff \bar{y} &= \hat{k} \hat{\theta} \end{aligned}$$

Door af te leiden naar beide schatters wordt er een stelsel van 2 vergelijkingen verkregen. Uit de vergelijking voor de afgeleide naar $\hat{\theta}$ kon de schatter $\hat{\theta}$ worden gehaald. Dit was echter niet het geval voor de vergelijking die we verkregen nadat de log likelihood vergelijking afgeleid werd naar \hat{k} . Deze vergelijking werd opgelost met een numerieke methode.

Hiervoor hebben we de verschillende "Nonlinear solvers" getest op de performantie. De ene methode zal sneller zijn voor de ene vergelijking en een andere methode zal sneller zijn voor een andere vergelijking. Aangezien dit heel vaak wordt herhaald zal zelfs een kleine verbetering in snelheid een grote invloed hebben op de snelheid om het script uit te voeren. Uiteindelijk kwamen we uit op een gelijkspel tussen 'broyden1' en 'anderson'

$$\begin{cases} 0 &= \overline{\ln(y)} - \ln(\hat{\theta}) - \psi^{(0)}(\hat{k}) \\ \bar{y} &= \hat{k} \hat{\theta} \end{cases}$$

Wanneer de vergelijking voor $\hat{\theta}$ wordt ingevuld in vergelijking (1) krijgen we een vergelijking die enkel afhankelijk is van de schatter \hat{k} . Deze vergelijking kan dan numeriek opgelost worden.

$$\overline{\ln(y)} - \ln(\bar{y}) + \ln(\hat{k}) - \psi^{(0)}(\hat{k}) = 0$$

Om de variantie te bepalen van de schatters kan de volgende vergelijking worden gebruikt:

$$V(\hat{a})^{-1} = -\frac{d^2}{da^2} \left[\ln(\mathcal{L}) \right] \Big|_{a=\hat{a}}$$

1) Variantie \hat{k} :

$$\begin{aligned} V(\hat{k})^{-1} &= -\frac{\partial}{\partial k} \sum_i^N \left[\ln(y_i) - \ln(\hat{\theta}) - \frac{\Gamma'(k)}{\Gamma(k)} \right] \Big|_{k=\hat{k}} \\ &= N\psi^{(1)}(\hat{k}) \end{aligned}$$

2) De variantie van $\hat{\theta}$

$$\begin{aligned} V(\theta)^{-1} &= -\frac{\partial}{\partial \theta} \sum_{i=1}^N \left[\frac{y_i}{\theta^2} - \frac{\hat{k}}{\theta} \right] \Big|_{\theta=\hat{\theta}} \\ &= N \frac{\hat{k}}{\hat{\theta}^2} \end{aligned}$$

De covariantie kan berekend worden via:

$$\begin{aligned} \text{cov}^{-1}(\hat{k}, \hat{\theta}) &= -\frac{\partial^2}{\partial k \partial \theta} \left[\ln(\mathcal{L}) \right] \Big|_{\substack{k=\hat{k} \\ \theta=\hat{\theta}}} \\ &= -\frac{\partial}{\partial \theta} \sum_i^N \left[\ln(y_i) - \ln(\theta) - \frac{\Gamma'(\hat{k})}{\Gamma(\hat{k})} \right] \Big|_{\theta=\hat{\theta}} \\ &= \frac{N}{\hat{\theta}} \\ &= -\frac{\partial}{\partial \theta} \sum_{i=1}^N \left[\frac{y_i}{\hat{\theta}^2} - \frac{\hat{k}}{\hat{\theta}} \right] \Big|_{k=\hat{k}} \\ &= \text{cov}^{-1}(\hat{\theta}, \hat{k}) \end{aligned}$$

Hierna kunnen deze uitkomsten in een matrix worden gezet om de inverse covariantiematrix te bepalen.

$$V^{-1} = \begin{bmatrix} N\psi^{(1)}(\hat{k}) & \frac{N}{\hat{\theta}} \\ \frac{N}{\hat{\theta}} & N \frac{\hat{k}}{\hat{\theta}^2} \end{bmatrix}$$

Om de echte variantie te berekenen moet deze matrix worden geïnverteerd. Dit wordt volledig in python gedaan en niet expliciet uitgerekend in het verslag. Deze matrixinversie creëert ook een afhankelijkheid voor \hat{k} en $\hat{\theta}$ voor zowel de variantie van \hat{k} als de covariantie. Iets wat niet aanwezig is in de inverse variantie en covariantie.

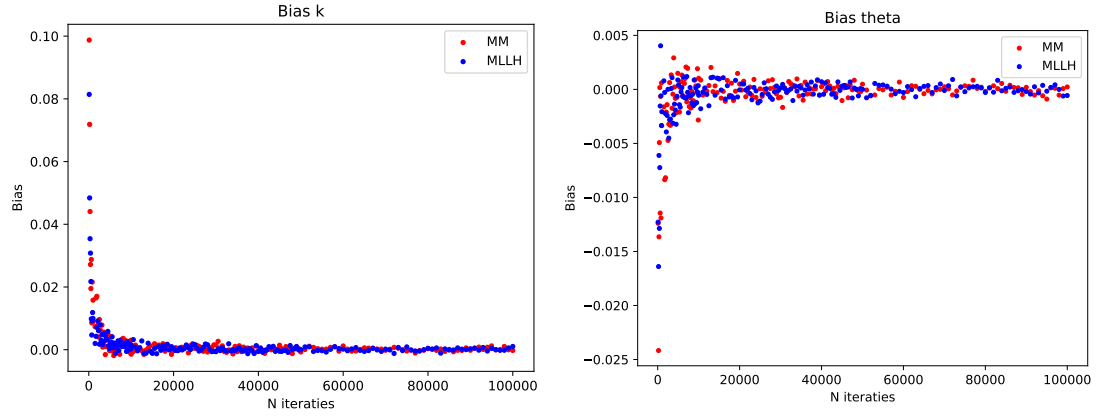
Verder werd er ook gebruik gemaakt van de polygamma functie $\psi^{(n)}(x)$. Waarbij $\psi^{(0)}(x)$ de digamma functie is.

3 Vergelijken van MM en MLLH-methode

Nu bestuderen we de eigenschappen van \hat{k} en $\hat{\theta}$ als functie van het aantal getrokken samples. Dit wordt gedaan aan de hand van bootstrappen. Hierbij nemen we N willekeurige samples uit de set metingen en passen we de schatters \hat{k} en $\hat{\theta}$ toe. Dit werd dan M keer herhaald. op basis van de M resultaten die we uitkwamen, werd de bias, $V(\hat{k})$, $V(\hat{\theta})$ en $\rho(\hat{k}, \hat{\theta})$ bepaald. Wij kozen $M = 1000$.

Dit werd herhaald voor enkele metingen tussen 1 en de hoeveelheid data N_{max} . Hoe groter de waarde voor N , hoe kleiner het verschil tussen N en $N - 1$ werd. Hierdoor was het mogelijk om de stappen te vergroten naarmate N groter werd, wat ons wat winst in snelheid opleverde.

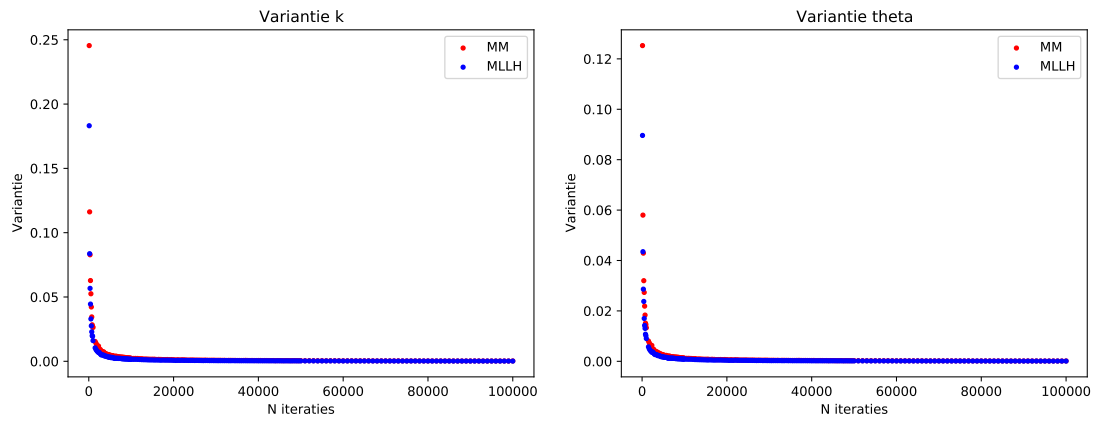
In figuur 1, 2 en 3 zijn de plots te zien voor de bias, de variantie en de correlatiecoëfficiënt in functie van N .



(a) *Bias van \hat{k} in functie van N_{max}*

(b) *Bias van $\hat{\theta}$ in functie van N_{max}*

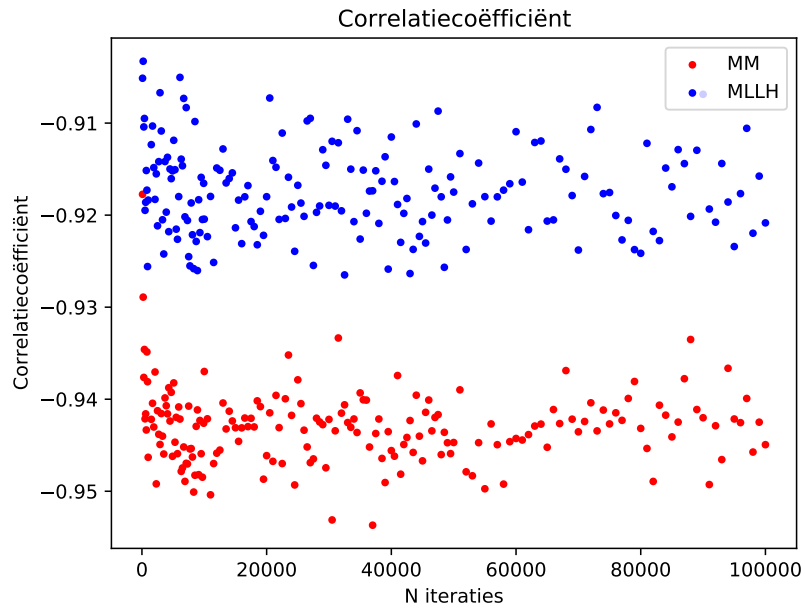
Figuur 1: *Bias voor \hat{k} en $\hat{\theta}$ voor MM en MLLH*



(a) *Variantie van \hat{k} in functie van N_{max}*

(b) *Variantie van $\hat{\theta}$ in functie van N_{max}*

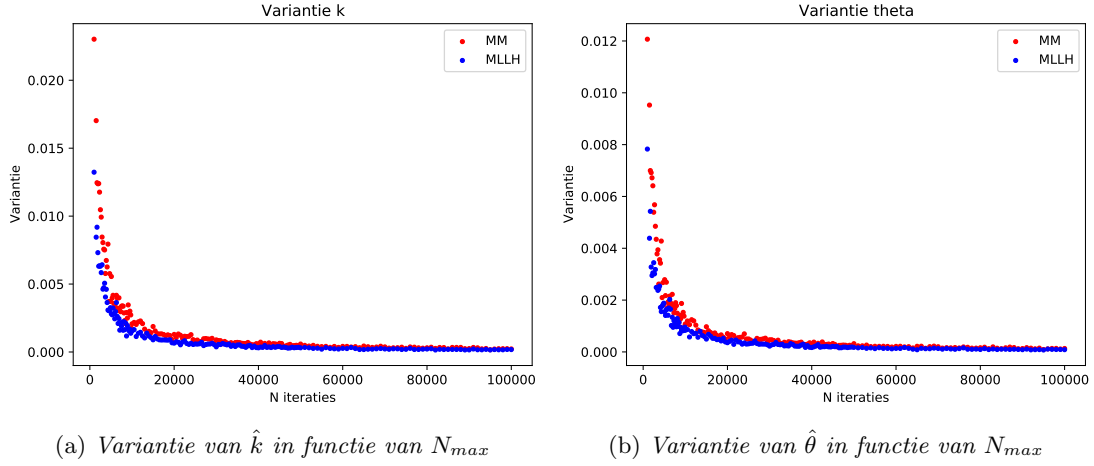
Figuur 2: *Variantie voor \hat{k} en $\hat{\theta}$ voor MM en MLLH*



Figuur 3: *Correlatiecoëfficiënt voor \hat{k} en $\hat{\theta}$ in functie van N_{max}*

In figuur 1 zien we de biassen van \hat{k} en $\hat{\theta}$. Er is weinig verschil op te merken tussen de bias voor de schatters bij de methode van de momenten en de maximum likelihood methode.

Voor de variantie is er echter vrij weinig te zien op de plots in figuur 2. Hiervoor maakten we een extra plotje waar de punten pas vanaf $N = 1000$ werden geplot om een duidelijker beeld te krijgen van wat er gebeurd. Dit is te zien in figuur 4.



Figuur 4: Variantie voor \hat{k} en $\hat{\theta}$ voor MM en MLLH

Deze figuren tonen duidelijk dat de variantie van beide schatters groter is als de methode van de momenten wordt toegepast. Dit betekent dat de MLLH-methode toch wat efficiënter is dan de MM.

Ten slotte zien we bij de correlaties een duidelijke scheiding. De correlaties zijn niet afhankelijk van het aantal getrokken waarden wat te verwachten is. We zien een bijna discrete scheiding tussen de correlaties. Ook is de absolute waarde van de correlatiecoëfficiënt van de schatters \hat{k} en $\hat{\theta}$ kleiner bij de MLLH-methode dan bij de MM. De MLLH-schatters zijn dus minder gecorreleerd.

De MLLH-methode is efficiënter dan de MM. Er moet echter wel vermeld worden dat het bepalen van de schatters hierbij een stuk ingewikkelder was. \hat{k} moest namelijk numeriek uitgerekend worden. Dit zorgde ervoor dat ook het bootstrappen bij de MLLH-methode een stuk langer duurde. Verder waren de MLLH-schatters minder gecorreleerd dan de MM-schatters. Hierdoor komen we tot dezelfde conclusie als de syllabus. De MM is een stuk eenvoudiger maar iets minder precies en dus goed voor een snelle schatting. Wanneer men preciezere resultaten wilt, is de MLLH-methode beter, maar dit is een pak ingewikkelder en duurt ook een stuk langer.

Indien we de methode slechts 1 keer moeten toepassen, is de tijd een minder groot probleem. Maar in toepassingen zoals deze, waar de methodes heel veel keer worden toegepast (in ons geval ongeveer 17,6 miljard keer), is het verschil in tijd van veel groter belang.

Het duurde 152 seconden om de berekeningen voor de MM uit te voeren. Voor de MLLH-methode was dit 595 seconden. Wat net geen 4 keer de tijd voor de MM is.

4 Vergelijken van de gebootstrapte varianties en de schattingen van de varianties

4.1 Methode van de Momenten

	bootstrap	schatting
$V(\hat{k})$	$0.231 \cdot 10^{-3}$	$0.241 \cdot 10^{-3}$
$V(\hat{\theta})$	$0.119 \cdot 10^{-3}$	$0.127 \cdot 10^{-3}$
$\rho_{\hat{k}, \hat{\theta}}$	-0.941	-0.943

Tabel 1: Vergelijking tussen de bootstrap en de schatting voor MM

Bij beide varianties en de correlatie kunnen we zien dat de waarden voor de bootstrap kleiner zijn dan die van de schatting. Het verschil is wel relatief klein, vooral bij de covariantie.

4.2 Maximum Log Likelihood Methode

	bootstrap	schatting
$V(\hat{k})$	$0.160 \cdot 10^{-3}$	$0.159 \cdot 10^{-3}$
$V(\hat{\theta})$	$88.6 \cdot 10^{-6}$	$86.8 \cdot 10^{-6}$
$\rho_{\hat{k}, \hat{\theta}}$	-0.920	-0.918

Tabel 2: Vergelijking tussen de bootstrap en de schatting voor MLLH

In tegenstelling tot de methode van de momenten is de bootstrap altijd iets groter dan de schatting. Zowel voor de varianties als de correlatie. Het verschil is echter zeer minimaal. Het is ook een stuk kleiner dan het verschil dat we zien bij de methode van de momenten. De variantie voor $\hat{\theta}$ is zelfs een grote orde kleiner bij de MLLH-methode dan bij de MM. Ook dit resultaat versterkt de stelling die in deel 3 werd gemaakt. Daar stelden we dat de MLLH-methode preciezer is dan de MM.