

Funktiomuotoisen heritabiliteetin estimointi  
bayesiläisillä malleilla

Tilastotieteen pro gradu -tutkielma

Arttu Arjas

2319498

Matemaattisten tieteiden yksikkö

Oulun yliopisto

Kevät 2019

## Tiivistelmä

Heritabiliteetti eli perinnöllisyysaste on perinnöllisyystieteen tunnusluku, joka mittaa sitä, kuinka suuri osa jonkin piirteen vaihtelusta tietyssä biologisessa populaatiossa riippuu perintötekijöistä. Sen avulla voidaan muun muassa selvittää, kuinka nopeasti valinnalla voidaan muuttaa piirteen keskiarvoa tulevissa sukupolvissa. Heritabiliteetti ei välttämättä ole vakio, vaan se voi olla ajan funktio. Kirjallisuudessa ei kuitenkaan ole kovin paljon esimerkkejä heritabiliteettifunktion estimoinnista.

Tässä työssä esitellään kaksi bayesiläistä menetelmää funktiomuotoisen heritabiliteetin estimointiin. Ensimmäisessä, kaksivaiheisessa menetelmässä heritabiliteetti ja sen luottamusväli estimoidaan jokaisessa aikapisteessä erikseen lineaarista sekamallia ja bootstrap-menetelmää apuna käyttäen, jonka jälkeen muodostuneet aikasarjat silotetaan. Toisessa menetelmässä sekamalli on ikään kuin yleistetty pitkittäisaineistotilanteeseen, jolloin heritabiliteetti saadaan estimoitua jokaisessa aikapisteessä samaan aikaan yksivaiheisesti sopivien silottavien priorien avulla.

Tutkimuksessa havaittiin, että yksivaiheisessa menetelmässä estimaatit ovat melko täsmällisiä estimoitaville parametreille asetettavien priorien takia, koska ne muodostavat eri aikapisteiden välille tietynlaisen korrelaatorakenteen. Toisaalta tämä estimointitapa on kuitenkin laskennallisesti raskas. Kaksivaiheisessa menetelmässä estimaattien luottamusvälit ovat leveämmät, mutta menetelmä ei ole läheskään niin aikaavievä kuin yksivaiheinen estimointi.

# Sisältö

<b>1</b>	<b>Johdanto</b>	<b>1</b>
<b>2</b>	<b>Mallien ja menetelmien rakentaminen</b>	<b>4</b>
2.1	Bayesiläinen analyysi . . . . .	4
2.2	Sukulaisuusmatriisi . . . . .	6
2.3	Lineaarinen sekamalli . . . . .	8
2.4	Sekamallin yleistys usealle ajanhetkelle . . . . .	9
2.5	Silotus ja priorit . . . . .	10
2.6	Kaksivaiheinen menetelmä . . . . .	14
2.7	Yksivaiheinen menetelmä . . . . .	15
<b>3</b>	<b>Parametrien estimointimenetelmät</b>	<b>17</b>
3.1	Rajoitettu uskottavuusfunktio (REML) . . . . .	17
3.2	Bootstrap-menetelmä . . . . .	18
3.3	Monte Carlo -integrointi . . . . .	19
3.4	Markovin ketju Monte Carlo -menetelmät (MCMC) . . . . .	20
3.5	Metropolisin–Hastingsin algoritmi . . . . .	21
3.6	Gibbsin otanta . . . . .	22
3.7	Elliptinen viipaleotanta (Elliptical slice sampling) . . . . .	23
3.8	MCMC-algoritmien suppenemisen tarkastelu . . . . .	25
<b>4</b>	<b>Aineistot ja niiden analyysi</b>	<b>28</b>
4.1	Simuloitu aineisto . . . . .	28
4.2	Todellinen aineisto . . . . .	28
4.3	Algoritmit . . . . .	30
<b>5</b>	<b>Tulokset</b>	<b>32</b>
5.1	Simuloitu aineisto . . . . .	32
5.2	Todellinen aineisto . . . . .	33

<b>6</b>	<b>Pohdinta</b>	<b>36</b>
	<b>Viitteet</b>	<b>38</b>
<b>A</b>	<b>Kroneckerin matriisitulo</b>	<b>41</b>
<b>B</b>	<b>MCMC-algoritmit</b>	<b>42</b>
<b>C</b>	<b>Jälkikuviot</b>	<b>44</b>

# 1 Johdanto

Heritabiliteetti on tärkeä parametri perinnöllisyystieteessä. Se mittaa sitä, kuinka suuri osa jatkuva-asteikollisen piirteen vaihtelusta selittyy geneettisillä tekijöillä tietyssä populaatiossa. Se vaikuttaa muun muassa siihen, kuinka nopeasti valinnalla voidaan vaikuttaa ominaisuuden keskiarvoon tulevaisuudessa. Heritabiliteetin mittaaminen on vaikeaa klassisilla menetelmillä, koska niitä varten tarvitaan joko tietoa populaation sukupuusta tai risteytyskokeita (Sillanpää, 2011).

Heritabiliteetin määritelmä lähtee siitä, että jonkin piirteen tai fenotyypin havaittu arvo jaetaan perintötekijöistä riippuvaan osaan ja ympäristötekijöistä riippuvaan osaan:

$$\text{Fenotyyppi}(P) = \text{Geenit}(G) + \text{Ympäristö}(E). \quad (1)$$

Fenotyypin varianssi voidaan esittää geneettisen varianssikomponentin ja ympäristövarianssikomponentin summana, jos geneistä ja ympäristöstä riippuvat satunnaistermit oletetaan toisistaan riippumattomiksi

$$\sigma_P^2 = \sigma_G^2 + \sigma_E^2. \quad (2)$$

Yleensä riippumattomuusoletus perustuu ainoastaan siihen, että riippuvutta on vaikeaa tai mahdotonta mitata. (Visscher ym. 2008). Lopulta heritabiliteetti voidaan määritellä geneettisen varianssin osuutena kokonaisvarianssista

$$h^2 = \frac{\sigma_G^2}{\sigma_G^2 + \sigma_E^2}. \quad (3)$$

Varianssikomponentit ovat populaatioparametreja ja ne täytyy estimoida sopivasta aineistosta. Geneettisen varianssin estimointia varten täytyy

silloin olla informaatiota populaation sukulaisuusrakenteesta. Rakenteen voi ilmaista nk. sukulaisuusmatriisin avulla. Sukulaisuusmatriisin konstruoinnista on kerrottu alaluvussa 2.2. Heritabiliteetti ei välttämättä ole vakio, vaan se voi riippua ajasta. Esimerkiksi piirteeseen vaikuttavan geenin ilmentyminen saattaa vaihdella ajan suhteen (Bryois ym. 2017) ja ympäristö voi muuttua. Uusia automatisoituja menetelmiä eri fenotyyppien mittaamiseen on tullut käyttöön. Ne mahdollistavat esimerkiksi kasvien kasvuun ja kehitykseen liittyvien piirteiden täsmällisen mittaamisen yli ajan. Laadukas pitkittäisaineisto mahdollistaa biologisten ilmiöiden funktionaalisen mallintamisen, jossa voidaan hyödyntää sopivia tilastollisia malleja. Funktionaalisen mallintamisen iso hyöty on se, että siinä yhdistyy useiden ajanhetkien informaatio, jolloin on mahdollista, että estimointi tarkentuu verrattuna estimointiin vain yhdellä ajanhetkellä. (Li ja Sillanpää, 2015)

Automatisoidut aineistonkeräysmenetelmät luovat tarpeen tilastollisille menetelmille, jotka ottavat huomioon aineiston aikaulottuvuuden. Yleensä näiden menetelmien toivotaan myös olevan laskennallisesti nopeita (Kwak ym. 2014). Tässä tutkielmassa esitellään kaksi bayesiläistä menetelmää ajasta riippuvan heritabiliteetin mallintamiseen. Bayesiläisessä analyysissä (alaluku 2.1) voidaan luontevasti ottaa huomioon ennakkokäsitykset ja -tiedot tutkitavasta ilmiöstä nk. priorijakaumien avulla. Menetelmissä hyödynnettävät priorit perustuvat ajatukseen, että heritabiliteettifunktio (kuten oletettavasti moni muukin biologinen prosessi) on sileä (Pletcher ja Geyer, 1999). Eräs tällainen silottava priorin on multinormaalijakauma, joka muodostaa aikapisteen välille korrelaatorakenteen kovarianssimatriisin avulla, mikä tarkentaa estimointia verrattuna estimointiin vain yhdellä ajanhetkellä.

Ensimmäinen menetelmä on kaksivaiheinen. Siinä heritabiliteetti ja sen

luottamusväli estimoidaan jokaisella ajanhetkellä erikseen, jonka jälkeen saadut aikasarjat silotetaan. Samantyylistä aikapisteittäistä estimointia funktioarvoisille ilmiöille ehdottaa myös Kwak ym. (2014). Menetelmä on nopea ja se skaalautuu hyvin aineiston koon kasvaessa, mutta se ei ota huomioon aikapisteiden välistä riippuvutta, jolloin koko aineiston tilastollinen voima jää hyödyntämättä. Toinen menetelmä estimoi heritabiliteettifunktion yhdellä kertaa. Hyvää tässä menetelmässä on se, että se ottaa aikapisteiden välisen riippuvuuden huomioon ja tuottaa täsmällisempiä estimaatteja kuin kaksivaiheinen menetelmä. Se on kuitenkin huomattavasti hitaampi eikä skaalaudu hyvin aineiston koon kasvaessa.

Tämän tutkielman rakenne on seuraava: Ensin avataan yleisesti bayesiläistä viitekehystä tilastollisessa päättelyssä, jonka jälkeen rakennetaan malli fenotyypille yhdellä ajanhetkellä. Kaksivaiheinen menetelmä perustuu tähän malliin. Tämän jälkeen malli yleistetään tilanteeseen, jossa havaintoja on usealta ajanhetkeltä. Yksivaiheinen menetelmä perustuu tähän yleistettyyn malliin. Seuraavaksi perehdytään siihen, miten biologisiin prosesseihin liittyvät ennakkokäsitykset saadaan liitettyä osaksi malleja, ja miten mallien parametrit voidaan estimoida. Lopuksi esitellään malleilla analysoitavat aineistot, analyysiin käytettävät algoritmit ja tulokset.

## 2 Mallien ja menetelmien rakentaminen

Tässä luvussa kerrotaan yleisesti bayesiläisestä analyysistä, jonka jälkeen rakennetaan mallit, joiden avulla heritabiliteettifunktiota estimoidaan. Lisäksi esitellään sukulaisuusmatriisi, joka on keskeinen osa malleja.

### 2.1 Bayesiläinen analyysi

Bayesiläinen päättely perustuu havaintoaineiston ja tutkittavaan aiheeseen liittyvien ennakkokäsitysten eli priorin yhdistämiseen. Gelman ym. (2013, s. 3) jakaa bayesiläisen data-analyysiprosessin kolmeen vaiheeseen.

1. Parametrien yhteis(priori)jakauman asettaminen: yhteistodennäköisyysjakauma kaikille mallin tuntemattomille parametreille.
2. Posteriorijakauman tutkiminen: kun parametrien yhteisjakauma ehdollistetaan havaintoaineistolle, saadaan posteriorijakauma, johon tilastollinen päättely kohdistetaan.
3. Päättely ja mallidiagnostiikka: kuinka hyvin malli sopii dataan ja millaisia päätelmiä posteriorijakaumasta voidaan tehdä sekä kuinka herkkä posteriorijakauma on mallioletusten muutoksille?

Bayesiläisen ja frekventistisen lähestymistavan perustavanlaatuinen ero on se, että bayesiläisessä mallin parametrit ovat ikään kuin ei-havaittavia satunnaismuuttujia, joilla on todennäköisyysjakauma, kun taas frekventistisessä ne ovat kiinteitä arvoja. Tämä muuttaa tapaa, jolla tilastollisen analyysin tuloksia voidaan tulkita. Frekventistinen 95 %:n luottamusväli voidaan tulkita niin, että jos tehtäisiin useita samanlaisia satunnaiskokeita, 95 %:ssa niistä luottamusväli peittäisi parametrin oikean arvon. Bayesiläinen 95 %:n posterioriväli taas voidaan tulkita niin, että sillä on 95 %:n todennäköisyys sisältää tuntematon parametrin arvo. Posterioriväli sisältää



myös priori-informaatiota parametrasta, kun taas luottamusväli perustuu vain havaintoaineistoon ja havainnoille oletettuun malliin. Posteriorivälin tulkinta siis tuntuu hieman helpommalta ja intuitiivisemmalta kuin luottamusvälin tulkinta. Bayesiläistä analyysiä rajoittavat lähinnä käytännön ongelmat, jotka liittyvät monimutkaisten mallien rakentamiseen ja tietokoneiden laskentakapasiteetin rajallisuuteen.

Bayes-päättele perustuu Bayesin kaavaan. Olkoon  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_p)^T$  vektori, joka sisältää kaikki tilastollisen mallin parametrit ja  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  havaintoaineisto. Mallinnus aloitetaan havaintoaineiston ja parametrien yhteistodennäköisyystiheydestä  $p(\boldsymbol{\theta}, \mathbf{y})$ . Ehdollisen todennäköisyystiheyden kaavasta saadaan

$$p(\boldsymbol{\theta}, \mathbf{y}) = p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta}), \quad (4)$$

jossa  $p(\mathbf{y}|\boldsymbol{\theta})$  on tilastollinen malli havainnoille tai uskottavuusfunktio  $\boldsymbol{\theta}$ :n funktiona ja  $p(\boldsymbol{\theta})$  on parametrien reunatiheys, jota kutsutaan prioritiheydeksi. Nyt posterioritiheys  $p(\boldsymbol{\theta}|\mathbf{y})$  saadaan Bayesin kaavasta

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\boldsymbol{\theta}, \mathbf{y})}{p(\mathbf{y})} = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})}, \quad (5)$$

jossa jatkuva-asteikollisten parametrien tapauksessa  $p(\mathbf{y}) = \int p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$  on nk. normalisointivakio, joka varmistaa, että posterioritiheyden integraali yli parametriavaruuden on yksi. (Gelman ym. 2013, s. 6–7). Monesti (esimerkiksi MCMC-algoritmeissa) normalisointivakio ei ole kiinnostava eikä sitä tarvitse evaluoida. Tällöin on tyypillistä kirjoittaa kaava normalisoimattomalle posterioritiheydelle

$$p(\boldsymbol{\theta}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta}). \quad (6)$$

Monissa tilastollisissa malleissa on parametreja, jotka riippuvat toisistaan luonnollisella tavalla ongelman rakenteen kautta. Tällöin puhutaan hierarkkisista malleista. Esimerkiksi bayesiläisessä silotusongelmassa ajatellaan,

että sileä trendi  $\mathbf{f}$  on tuntematon parametrivektori jota estimoidaan, ja sille määrätään silottava priori. Tämä priorikin saattaa riippua jostakin tuntemattomasta skalaariparametrasta  $\lambda$ , joka voidaan myös estimoida.  $\lambda$  voi esimerkiksi kontrolloida sitä, kuinka sileä trendi on. Tällöin yhteisprioritiheys voidaan kirjoittaa muodossa

$$p(\mathbf{f}, \lambda) = p(\mathbf{f}|\lambda)p(\lambda). \quad (7)$$

Posterioritiheydeksi saadaan

$$p(\mathbf{f}, \lambda|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{f}, \lambda)p(\mathbf{f}, \lambda) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\lambda)p(\lambda). \quad (8)$$

Nyt uskottavuusfunktio  $p(\mathbf{y}|\mathbf{f}, \lambda)$  riippuu vain parametrivektorista  $\mathbf{f}$ .  $\lambda$  on nimeltään hyperparametri eli parametri, josta  $\mathbf{f}$ :n priori riippuu. Hierarkkinen mallinnus auttaa hahmottamaan ongelmaa ja se helpottaa Bayes-mallien tietokonelaskentaa. (Gelman ym. 2013, s. 107–108)

## 2.2 Sukulaisuusmatriisi

Sukulaisuusmatriisi  $\mathbf{A}$  on symmetrinen kovarianssimatriisi, jossa matriisin alkio  $A_{ij}$  kuvaa yksilöiden  $i$  ja  $j$  välistä sukulaisuuden astetta. Sukulaisuusmatriisin avulla havaintoaineistosta saadaan identifioitua geneettinen varianssikomponentti. Klassinen tapa konstruoida sukulaisuusmatriisi perustuu sukupuuhun. Sukupuuta ei ole kuitenkaan aina saatavilla. Uudempi tapa perustuu molekulaarisiin markkereihin, joilla selvitetään alleelin tyyppi eri kohdista kromosomia. Sukulaisuusmatriisi perustuu sukulaisuuskertoimeen (kinship coefficient), ja matriisin alkio  $A_{ij}$  saadaan kertomalla yksilöiden  $i$  ja  $j$  välisen sukulaisuuskertoimen arvo kahdella. Sukulaisuuskertoimen määrittelyssä taas tärkeitä termejä ovat syntyperältään identtiset (identical by descent, IBD) ja tilaltaan identtiset (identical by state, IBS) alleelit.

Jos tietyssä geenipaikassa, joka on määrätty paikka kromosomissa, yksilön

$i$  alleeli on fyysisesti sama kuin yksilöllä  $j$ , ja ne ovat molemmat kopioita yksilöiden yhteiseltä esi-isältä peritystä alleelistä, sanotaan, että alleelit ovat syntyperältään identtiset. Jos taas tietyssä geenipaikassa riippumatta alleelin alkuperästä yksilön  $i$  alleeli on samantyyppinen kuin yksilöllä  $j$ , sanotaan, että alleelit ovat tilaltaan identtiset.

Sukulaisuuskerroin pohjaa IBD-käsitteeseen. Se on todennäköisyys, jolla yksilöillä  $i$  ja  $j$  satunnaisesta geenipaikasta otetut alleelit ovat syntyperältään identtiset. Sukulaisuuskertoimen arvo kerrottuna kahdella antaa odotetun IBD-alleelien osuuden kahden yksilön välillä. Sukupuuhun perustuva sukulaisuusmatriisi sisältää kaikkien yksilöiden väliset odotetut IBD-alleelien osuudet. Jos sukupuuta ei ole saatavilla, voidaan odotetun IBD-alleelien osuuden sijasta laskea realisoitunut IBD-alleelien osuus kahden yksilön välillä, joka perustuu molekulaarisiin markkereihin. Näistä luvuista koostuvaa matriisia kutsutaan markkeripohjaiseksi sukulaisuusmatriisiksi. (Sillanpää, 2011)

Eräs tapa muodostaa markkeripohjainen sukulaisuusmatriisi on kuvattu seuraavassa. Tietyssä markkerissa genotyyppi voi olla joko AA, AB tai BB, jotka voidaan koodata numeroiksi järjestyksessä  $-1$ ,  $0$  ja  $1$ . Olkoon  $\mathbf{X}$   $N \times M$  -matriisi, jossa  $N$  viittaa yksilöiden määrään ja  $M$  markkerien määrään ja  $X_{ij} \in \{-1, 0, 1\}$  riippuen yksilön  $i$  genotyypistä markkerissa  $j$ . Lisäksi olkoon  $p_j$  alleelin A alleelifrekvenssi markkerissa  $j$ . Nyt voidaan muodostaa sukulaisuusmatriisi

$$\mathbf{A} = \frac{1}{c} \mathbf{Z} \mathbf{Z}^T, \quad (9)$$

jossa  $Z_{ij} = X_{ij} + 1 - 2p_j$  ja  $c = 2 \sum_j p_j(1 - p_j)$ . Tämä ei ole ainoa keino muodostaa markkeripohjaista sukulaisuusmatriisia. Tämän ja muita menetelmiä on listannut esimerkiksi VanRaden (2008).

## 2.3 Lineaarinen sekamalli

Jatkuva-asteikollista piirrettä voidaan mallintaa lineaarisen sekamallin avulla, joka on eräänlainen laajennus tavanomaiseen lineaariseen malliin. Sillä voidaan kätevästi ottaa huomioon vastevektorin alkioden välinen riippuvuus, joka tässä tapauksessa johtuu populaation yksilöiden välisestä sukulaisuudesta. Oletetaan, että mitataan jotain piirrettä  $N$ :ltä yksilöltä jossain populaatiossa. Lineaarinen sekamalli määritellään (Kang ym. 2008)

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}, \quad (10)$$

jossa  $\mathbf{y}$  on  $N \times 1$  -vektori vasteen havaittuja arvoja,  $\mathbf{X}$  on kiinteät vaikutukset sisältävä  $N \times P$  -matriisi,  $P \times 1$  -vektori  $\boldsymbol{\beta}$  sisältää kiinteiden vaikutusten regressiokertoimet,  $N \times Q$  -matriisi  $\mathbf{Z}$  liittää satunnaisvaikutukset  $\mathbf{u}$  ( $Q \times 1$  -vektori) siihen tasoon, jolla satunnaisvaikutus halutaan laskea, ja  $N \times 1$  -vektori  $\boldsymbol{\epsilon}$  on yksilökohtainen virhetermien vektori. Heritabiliteettia estimoidessa kiinteä vaikutus  $\boldsymbol{\beta}$  redusoituu pelkkään viitetasoon, joka on vektorin  $\mathbf{y}$  odotusarvo, ja  $\mathbf{X}$  on  $N \times 1$  -vektori, jossa jokainen alkio on 1. Lisäksi satunnaisvaikutus lasketaan jokaiselle yksilölle, joten  $\mathbf{Z} = \mathbf{I}$ . Oletetaan, että  $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \sigma_G^2 \mathbf{A})$ , jossa  $\mathbf{A}$  on sukulaisuusmatriisi, ja  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_E^2 \mathbf{I})$ . Kun lisäksi  $\mathbf{y}$  keskistetään, voidaan malli kirjoittaa muodossa

$$\tilde{\mathbf{y}} = \mathbf{u} + \boldsymbol{\epsilon}, \quad (11)$$

jossa  $\tilde{\mathbf{y}}$  on keskistetty havaintovektori ja joka muistuttaa kaavaa (1). Yhtälöä (11) ja sen oletuksia vastaava uskottavuusfunktio parametreille  $\sigma_E^2$  ja  $\sigma_G^2$  on

$$p(\tilde{\mathbf{y}} | \sigma_G^2, \sigma_E^2) = \frac{1}{\sqrt{(2\pi)^N |\mathbf{K}|}} \exp \left\{ -\frac{1}{2} \tilde{\mathbf{y}}^T \mathbf{K}^{-1} \tilde{\mathbf{y}} \right\}, \quad (12)$$

jossa  $\mathbf{K} = \sigma_G^2 \mathbf{A} + \sigma_E^2 \mathbf{I}$ .

## 2.4 Sekamallin yleistys usealle ajanhetkelle

Kun estimoidaan heritabiliteettifunktiota, aineistossa täytyy olla mittauksia kiinnostavasta piirteestä populaatiossa usealla eri ajanhetkellä. Kaava (2) voidaan tällöin kirjoittaa muodossa

$$\sigma_P^2(t) = \sigma_G^2(t) + \sigma_E^2(t), \quad t = 1, \dots, T, \quad (13)$$

jossa  $\sigma_x^2(t)$  viittaa varianssikomponentin  $x$  arvoon ajanhetkellä  $t$ . Tavoitteena on luoda malli, joka määrittelee yhteisjakauman havaintovektorille  $\tilde{\mathbf{y}}_T$ , joka sisältää jokaisen ajanhetken oman keskistetyyn havaintovektorin ladottuna peräkkäin. Määrittely onnistuu Kroneckerin matriisitulon avulla (liite A), ja tuloksena saadaan  $\tilde{\mathbf{y}}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_T)$ , jossa  $\mathbf{K}_T = \text{diag}(\sigma_G^2(1), \dots, \sigma_G^2(T)) \otimes \mathbf{A} + \text{diag}(\sigma_E^2(1), \dots, \sigma_E^2(T)) \otimes \mathbf{I}$ . Aikapistekohtaiset varianssikomponentit on siis kerätty diagonaalimatriiseihin ja  $\mathbf{A}$  on jälleen sukulaisuusmatriisi. Symboli  $\otimes$  viittaa Kroneckerin tuloon. Uskottavuusfunktio vektoriarvoisille parametreille  $\sigma_G^2$  ja  $\sigma_E^2$ , jotka sisältävät jokaisen aikapisteen varianssikomponentit, on

$$p(\tilde{\mathbf{y}}_T | \sigma_G^2, \sigma_E^2) = \frac{1}{\sqrt{(2\pi)^{NT} |\mathbf{K}_T|}} \exp \left\{ -\frac{1}{2} \tilde{\mathbf{y}}_T^T \mathbf{K}_T^{-1} \tilde{\mathbf{y}}_T \right\}. \quad (14)$$

Jakauman kovarianssimatriisin rakennetta voi olla vaikea hahmottaa. Se on seuraavanlainen:

$$\mathbf{K}_T = \begin{bmatrix} \sigma_G^2(1)\mathbf{A} + \sigma_E^2(1)\mathbf{I} & 0 & \dots & 0 \\ 0 & \sigma_G^2(2)\mathbf{A} + \sigma_E^2(2)\mathbf{I} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_G^2(T)\mathbf{A} + \sigma_E^2(T)\mathbf{I} \end{bmatrix}.$$

Rakenne on siis lohkodeagonaalinen ja diagonaalilla toistuvat aikapistekohtaiset kovarianssimatriisit. Lohkodeagonaalinen rakenne nopeuttaa estimointia, koska matriisin determinantin laskeminen ja matriisin kääntäminen helpottuu. Determinantti on diagonaalilla olevien alimatriisien determinanttien tulo, ja kääntämismatriisia varten täytyy kääntää vain alimatriisit erikseen.

## 2.5 Silotus ja priorit

Varianssifunktioiden priorit valitaan siis sillä perusteella, että funktioiden ajatellaan olevan sileitä. Tästä seuraa se, että heritabiliteettifunktiokin on sileä. Silottava priorit on bayesiläinen tapa tehdä silotusta. Silotus tarkoittaa sitä, että havaintoaineistoon sovitetaan funktio, joka kiteyttää aineiston tärkeät piirteet ja poistaa epäsäännöllisen kohinan. Silotusta käytetään esimerkiksi kuvankäsittelyssä ja aikasarjojen avulla ennustamisessa. Silotukseen on kehitetty monia eri menetelmiä. Klassisia tapoja ovat eksponentiaalinen silotus ja liukuva keskiarvo. Kehittyneempiä menetelmiä ovat frekventistisellä puolella rosoisuuden sakottaminen ja bayesiläisellä puolella vastaavat silottavat priorit (Fahrmeir ja Kneib, 2011, s. 18).

Fahrmeir ja Kneib (2011, s. 19) määrittelee klassisen silotusongelman seuraavasti: Olkoon  $\mathbf{y} = (y_1, \dots, y_T)^T$  vektori tasavälisiä havaintoja. Havainnoille oletetaan malli

$$\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon}, \quad (15)$$

jossa  $\mathbf{f} = (f_1, \dots, f_T)^T$  on sileä trendi ja  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_T)^T$  on joukko riippumattomia ja samoin jakautuneita virhetermejä niin, että  $\mathbb{E}(\epsilon_t) = 0$ . Trendin estimointiin on ehdotettu sakotetun neliösumman (PLS; Penalized Least Squares) minimointia:

$$\text{PLS}(\mathbf{f}) = \underbrace{\sum_{t=1}^T (y_t - f_t)^2}_{\text{mallin istuvuus}} + \lambda \underbrace{\sum_{t=3}^T (f_t - 2f_{t-1} + f_{t-2})^2}_{\text{[sakkotermi]}}. \quad (16)$$

Silotusparametri  $\lambda \geq 0$  kontrolloi tasapainotusta harhan ja varianssin välillä ja se pitää säätää sopivaksi jokaiselle aineistolle erikseen. (Fahrmeir ja Kneib, 2011, s. 19). Parametrin vaikutusta on havainnollistettu kuvassa 1.

Bayesiläisessä silotuksessa oletetaan, että tuntematon trendi on satunnaisvektori  $\mathbf{f}$ , jolla on prioritiheys  $p(\mathbf{f})$ . Yleisin virhetermeille oletettu jakauma on normaalijakauma, jolloin yhtälöä (15) vastaava malli on

$$\mathbf{y}|\mathbf{f} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{I}). \quad (17)$$

Kun  $\sigma^2$  oletetaan tunnetuksi, voidaan normalisointivakio sivuuttaa ja saada

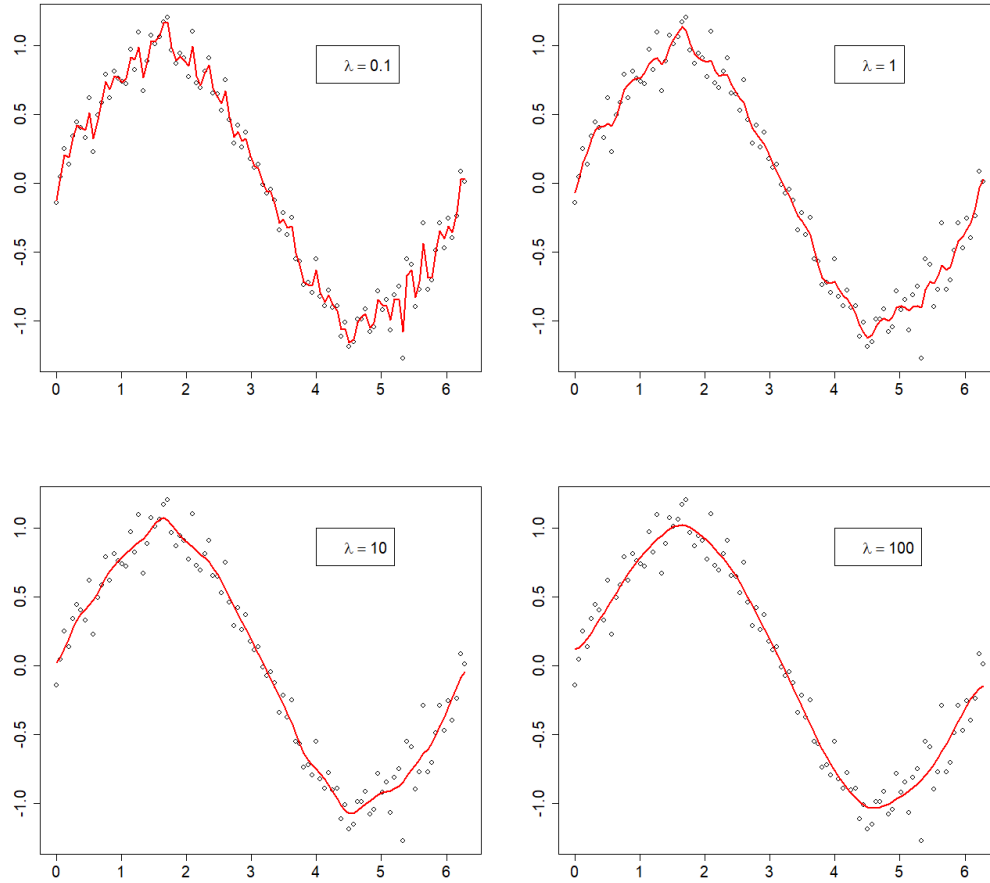
$$p(\mathbf{y}|\mathbf{f}) \propto \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{f})^T (\mathbf{y} - \mathbf{f}) \right\}. \quad (18)$$

Valittu prioriksi on yleensä gaussinen,  $\mathbf{f} \sim \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{P}^{-1})$ , jossa täsmällisyysmatriisina on jokin sakkomatriisi  $\mathbf{P}$ , jonka muoto riippuu siitä, millainen silotus halutaan tehdä. Tällöin prioritiheys  $\tau^2$ :n ollessa tunnettu on

$$p(\mathbf{f}) \propto \exp \left\{ -\frac{1}{2\tau^2} \mathbf{f}^T \mathbf{P} \mathbf{f} \right\}. \quad (19)$$

Kun sekä havaintomalli että prioriksi ovat gausset, voidaan trendin posterio-rijakauma johtaa analyttisesti. Kun merkitään  $\lambda = \sigma^2/\tau^2$ , se on

$$\mathbf{f}|\mathbf{y} \sim \mathcal{N}((\mathbf{I} + \lambda \mathbf{P})^{-1} \mathbf{y}, \sigma^2 (\mathbf{I} + \lambda \mathbf{P})^{-1}). \quad (20)$$



Kuva 1: Silotusparametrin vaikutus trendiin.

Sakkomatriisin  $\mathbf{P}$  rooli on ikään kuin vetää vierekkäisiä havaintoja toisiinsa kohti tietyn funktiomuodon perusteella. Silotusparametrin  $\lambda$  arvo kontrolloi silotuksen astetta (l. trendin sileyttä), joka tarkoittaa sitä kuinka vahvasti havaintoja halutaan vetää toisiaan kohti edellämainitun  $\mathbf{P}$ :n sisällä. (Fahrmeir ja Kneib, 2011, s. 23–25)

Tässä tutkielmassa silotus tehdään priorin avulla, mutta sakkomatriisiin



$\mathbf{P}$  sijaan muodostetaan kovarianssimatriisi  $\mathbf{C}$  tietyn kovarianssifunktion mukaan. Kovarianssifunktio määrää jokaisen aikapisteparin välille kovarianssin niiden etäisyyden perusteella. Toisiaan lähellä olevilla aikapisteillä on suurempi kovarianssi kuin toisistaan kaukana olevilla aikapisteillä. Kovarianssifunktion muoto ja parametrit määräävät muun muassa sen, kuinka sileää estimoitu funktio on.

Merkitään  $d = |t_i - t_j|$  kahden ajanhetken välillä kuluvaa aikaa. Neliöeksponentiaalinen kovarianssifunktio määritellään

$$k_{\text{SE}}(d) = \tau^2 \exp \left\{ -\frac{d^2}{2l^2} \right\}, \quad (21)$$

jossa  $l > 0$  on pituusskaalaparametri ja  $\tau^2 > 0$  on voimakkuusparametri. Pituusskaalaparametri vaikuttaa estimoidun funktion sileyteen ja voimakkuusparametri sen vaihteluun. Priori tällä kovarianssifunktiolla tuottaa todella sileitä funktioestimaatteja. Se on yksi käytetyimmistä kovarianssifunktioista. Eksponentiaalinen kovarianssifunktio on

$$k_{\text{E}}(d) = \tau^2 \exp \left\{ -\frac{d}{l} \right\}, \quad (22)$$

jossa parametrit ovat samat kuin neliöeksponentiaalisessa kovarianssifunktiossa. Priorit eksponentiaalisella kovarianssifunktiolla tuottavat huomattavasti rosoisempia funktioestimaatteja kuin neliöeksponentiaalisella kovarianssifunktiolla. Nämä kaksi kovarianssifunktiota ovat itse asiassa erityistapauksia yleisemmästä Matérn-luokasta, joka määritellään

$$k_{\text{Matérn}}(d) = \tau^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left( \frac{\sqrt{2\nu}d}{l} \right)^\nu K_\nu \left( \frac{\sqrt{2\nu}d}{l} \right), \quad (23)$$

jossa parametri  $\nu > 0$  määrää funktion sileysominaisuuksia ja  $K_\nu$  on toisen lajin Besselin funktio. Eksponentiaalinen kovarianssifunktio saadaan, kun asetetaan  $\nu = \frac{1}{2}$  ja neliöeksponentiaalinen kovarianssifunktio raja-arvona, kun  $\nu \rightarrow \infty$ . (Rasmussen ja Williams, 2006, s. 83–86)

## 2.6 Kaksivaiheinen menetelmä

Kaksivaiheinen menetelmä hyödyntää pelkästään alaluvussa 2.3 kaavassa (11) määriteltyä mallia yhden aikapisteen vastevektorille. Siinä varianssikomponentit estimoidaan jokaisessa aikapisteessä erikseen maksimoimalla rajoitettu uskottavuusfunktio, josta kerrotaan alaluvussa 3.1. Menetelmä ei ole bayesiläinen, mutta se osoittautui yhden ajanhetken tapauksessa nopeamaksi ja täsmällisemmäksi tavaksi estimoida varianssikomponentit. Tämän lisäksi komponenteille estimoidaan luottamusvälit bootstrap-menetelmällä, josta kerrotaan alaluvussa 3.2. Tällä tavalla tuotettuja estimaatteja voidaan pitää aikasarjoina, jotka voivat sisältää kohinaa, ja ne silotetaan priorien avulla. Silotusmallina käytetään kaavassa (15) määriteltyä yhtälöä. Aikasarjat voidaan ilmaista muodossa

$$\log \mathbf{y} = \log \mathbf{f} + \boldsymbol{\epsilon}, \quad (24)$$

jossa  $\mathbf{y}$  on varianssikomponentin tai sen luottamusvälirajan aikasarja,  $\mathbf{f}$  on aikasarjaa kuvaava sileä funktio ja  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$  on virhetermi. Aikasarjat logaritmoidaan, jotta virhetermien normaalisuusoletus olisi järkevä. Logaritmoidulle funktiolle  $\mathbf{f}$  asetetaan prioriksi  $\log \mathbf{f} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_l)$ , jossa matriisin  $\mathbf{C}_l$  alkio  $C_{lij} = k_{\text{Matérn}}(|t_i - t_j|)$ . Alaindeksi  $l$  tarkoittaa sitä, että matriisin muoto riippuu pituusskaalasta  $l$ , joka myös estimoidaan. Sille ja virhevariانسsille  $\sigma^2$  asetetaan laakeat priorit, jotka sisältävät vain vähän informaatiota siitä, missä parametrien todennäköisimmät arvot ovat. Nämäkin parametrit logaritmoidaan, koska ne voivat saada vain positiivisia arvoja. Kovarianssifunktion muut parametrit  $\nu$  ja  $\tau^2$  täytyy kiinnittää identifioituvuusongelmien takia (Monterrubio-Gómez ym. 2018). Näillä parametreilla, kun  $\log \mathbf{y}$  keskis-

tetään, mallin voi kirjoittaa hierarkkisessa muodossa

$$\begin{aligned}
\log \mathbf{y} | \mathbf{f}, \sigma^2 &\sim \mathcal{N}(\log \mathbf{f}, \sigma^2 \mathbf{I}), \\
\log \mathbf{f} | l &\sim \mathcal{N}(\mathbf{0}, \mathbf{C}_l), \\
\log l &\sim u_{\log l}(\boldsymbol{\theta}_{\log l}), \\
\log \sigma^2 &\sim u_{\log \sigma^2}(\boldsymbol{\theta}_{\log \sigma^2}),
\end{aligned} \tag{25}$$

jossa  $u$  viittaa johonkin laakeaan priorijakaumaan ja  $\boldsymbol{\theta}$  sen parametreihin. Mallin parametrien normalisoimaton posterioritiheys on

$$\begin{aligned}
p(\log \mathbf{f}, \log l, \log \sigma^2 | \log \mathbf{y}) &\propto p_{\mathcal{N}}(\log \mathbf{y} | \log \mathbf{f}, \sigma^2 \mathbf{I}) p_{\mathcal{N}}(\log \mathbf{f} | \mathbf{0}, \mathbf{C}_l) \\
p_{u; \log l}(\log l | \boldsymbol{\theta}_{\log l}) &p_{u; \log \sigma^2}(\log \sigma^2 | \boldsymbol{\theta}_{\log \sigma^2}),
\end{aligned} \tag{26}$$

jossa  $p_{\mathcal{N}}$  viittaa multinormaalijakauman tiheysfunktioon ja  $p_u$  laakean priorijakauman tiheysfunktioon. Koska malli sisältää tuntemattomia hyperparametrejä, funktion  $\mathbf{f}$  posteriorijakaumaa ei voida johtaa analyttisesti, kuten kaavassa (20). Sen sijaan parametrien estimointiin täytyy käyttää numeerisia approksimointimenetelmiä, joista kerrotaan lisää luvussa 3.

## 2.7 Yksivaiheinen menetelmä

Yksivaiheinen menetelmä heritabiliteettifunktion estimointiin perustuu jokaisen aikapisteen havaintovektorin sisältävän vektorin  $\tilde{\mathbf{y}}_T$  yhteisjakaumaan, joka on johdettu alaluvussa 2.4. Logaritmoituille varianssikomponenttivektoreille  $\sigma_E^2$  ja  $\sigma_G^2$  asetetaan jälleen multinormaalipriorit:  $\log \sigma_E^2 \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_E)$  ja  $\log \sigma_G^2 \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_G)$ . Malli voidaan esittää seuraavasti:

$$\begin{aligned}
\tilde{\mathbf{y}}_T &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}_T), \\
\log \sigma_G^2 &\sim \mathcal{N}(\mathbf{0}, \mathbf{C}_G), \\
\log \sigma_E^2 &\sim \mathcal{N}(\mathbf{0}, \mathbf{C}_E).
\end{aligned} \tag{27}$$

Käytännössä havaittiin, että kovarianssifunktion pituuskaala ei identifoidu yhtä hyvin kuin kaksivaiheisella menetelmällä. Niinpä sen arvo kiinnitettiin

kaksivaiheisella menetelmällä estimoituun posteriorikeskiarvoon. Niinpä mallin parametreinä ovat ainoastaan varianssikomponenttivektorit. Mallia vastaava normalisoimaton posterioritiheys on

$$p(\boldsymbol{\sigma}_G^2, \boldsymbol{\sigma}_E^2 | \tilde{\mathbf{y}}_T) \propto p_{\mathcal{N}}(\tilde{\mathbf{y}}_T | \mathbf{0}, \mathbf{K}_T) p_{\mathcal{N}}(\log \boldsymbol{\sigma}_G^2 | \mathbf{0}, \mathbf{C}_G) p_{\mathcal{N}}(\log \boldsymbol{\sigma}_E^2 | \mathbf{0}, \mathbf{C}_E). \quad (28)$$

Tämänkin mallin parametrit estimoidaan numeerisilla menetelmillä, joista kerrotaan seuraavassa luvussa.

### 3 Parametrien estimointimenetelmät

Bayesiläisessä analyysissä tilastollinen päättely joudutaan usein perustamaan numeerisiin approksimaatioihin, koska posteriorijakauman tiheysfunktiolla ei välttämättä ole tunnettua parametrissa muotoa, josta parametrien arvot voisi päätellä. Tässä luvussa on lyhyt katsaus tutkielmassa hyödynnettäviin numeerisiin menetelmiin. Lisäksi esitellään kaksivaiheisessa menetelmässä hyödynnettävät klassiset tilastollisen päättelyn menetelmät rajoitettu uskottavuusfunktio ja bootstrap.

#### 3.1 Rajoitettu uskottavuusfunktio (REML)

Skalaarimuotoiset varianssikomponentit  $\sigma_E^2$  ja  $\sigma_G^2$  voidaan estimoida maksimoimalla lineaarista sekamallia (alaluku 2.3) vastaava rajoitettu uskottavuusfunktio. Se on määritelmältään täysi uskottavuusfunktio, josta on analyttisesti integroitu satunnaistekijät pois ja joka on lineaarisen muunnoksen jälkeen riippumaton ja kohtisuorassa kiinteisiin tekijöihin nähden. REML-estimoinnissa vältetään varianssin aliarviointi, koska siinä otetaan huomioon kiinteiden vaikutusten estimoinnista aiheutuva vapausasteiden menetys.

Yhtälöä (10) ja sen oletuksia vastaava log-uskottavuusfunktio ja rajoitettu log-uskottavuusfunktio ovat (Kang ym. 2008)

$$l_F(\mathbf{y}; \boldsymbol{\beta}, \sigma_G^2, \delta) = \frac{1}{2} \left[ -N \log(2\pi\sigma_G^2) - \log |\mathbf{H}| - \frac{1}{\sigma_G^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{H}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \quad (29)$$

ja

$$l_R(\mathbf{y}; \sigma_G^2, \delta) = l_F(\mathbf{y}; \hat{\boldsymbol{\beta}}, \sigma_G^2, \delta) + \frac{1}{2} [P \log(2\pi\sigma_G^2) + \log |\mathbf{X}^T \mathbf{X}| - \log |\mathbf{X}^T \mathbf{H}^{-1} \mathbf{X}|], \quad (30)$$

joissa  $\mathbf{H} = \mathbf{A} + \delta \mathbf{I}$ ,  $\delta = \frac{\sigma_E^2}{\sigma_G^2}$  ja  $\mathbf{A}$  on sukulaisuusmatriisi. Rajoitetun uskottavuusfunktion maksimoimiseen tarvitaan numeerisia optimointimenetelmiä,

joihin ei tässä työssä syvennytä tarkemmin, mutta niistä kertoo enemmän Kang ym. (2008). Tässä tutkielmassa varianssikomponenttien estimointiin yhdessä aikapisteessä käytettiin R-ohjelmiston paketin `rrBLUP` (Endelman, 2011) funktiota `mixed.solve`.

## 3.2 Bootstrap-menetelmä

Usein jonkin kiinnostavan tunnusluvun (asymptoottista) otantajakaumaa ei ole mahdollista johtaa analyttisesti. Tällöin suurelle voi olla vaikea laskea esimerkiksi luottamusvälejä. Tietokoneiden tehostuminen on kuitenkin luonut mahdollisuuden uudelleenkäyttää otosta niin, että kiinnostaville suureille voidaan simuloida empiirinen otantajakauma ja tarkastella sen käyttäytymistä. Tässä työssä bootstrap-menetelmää käytetään aikapistekohtaisten varianssikomponenttien luottamusvälien estimointiin.

Olkoon  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  havaintoaineisto, jossa kukin alkio  $y_i$  on reaalisuus samasta parametrisesta jakaumasta, jonka kertymäfunktio on  $F(y; \theta)$ . Oletetaan, että suurimman uskottavuuden estimaatti  $\hat{\theta}$  on kohtuullisen helposti laskettavissa, ja että jakaumasta on mahdollista simuloida satunnaislukuja millä tahansa  $\theta$ :n arvolla. Parametrisessa bootstrapissa (Efron ja Hastie, 2016, s. 169) simuloidaan otos  $\mathbf{Y}^* = (Y_1^*, Y_2^*, \dots, Y_n^*)$  jakaumasta, jonka kertymäfunktio on  $F(y; \hat{\theta})$ . Olkoon nyt  $T = t(Y)$  jokin havainnoista laskettava tunnusluku. Se voi olla vaikkapa suurimman uskottavuuden estimaatti  $\hat{\theta}$  tai joku sen muunnos. Tunnusluvun  $T$  otantajakaumaa simuloidaan toistamalla seuraavaa  $B$  kertaa:

1. Poimitaan otos  $\mathbf{Y}^{*(b)} = (Y_1^{*(b)}, Y_2^{*(b)}, \dots, Y_n^{*(b)})$  jakaumasta, jonka kertymäfunktio on  $F(y; \hat{\theta})$
2. Lasketaan  $T^{*(b)} = t(\mathbf{Y}^{*(b)})$

Parametrittomassa bootstrapissa (Efron ja Hastie, 2016, s. 159-160) otosten simulointi toteutetaan eri tavalla. Siinä otetaan alkuperäisestä otoksesta takaisinpanolla  $n$  havaintoa ja lasketaan näistä otoksista tunnusluvun  $T$  arvo. Nyt  $(T^{*(1)}, T^{*(2)}, \dots, T^{*(B)})$  on otos  $T$ :n otantajakaumasta. Otoksesta voidaan laskea esimerkiksi otoskeskiarvo ja sen keskivirhe

$$\bar{T}^* = \frac{1}{B} \sum_{b=1}^B T^{*(b)}, \quad (31)$$

$$SE(\bar{T}^*) = \sqrt{\frac{1}{B-1} \sum_{b=1}^B (T^{*(b)} - \bar{T}^*)^2}. \quad (32)$$

$T$ :n  $1 - \gamma$  prosentin luottamusväli saadaan bootstrap-otoksen empiirisistä fraktiileista. Alarajaksi valitaan  $100\gamma/2$  %:n fraktiili  $T_{\gamma/2}^*$  ja ylärajaksi  $100(1 - \gamma/2)$  %:n fraktiili  $T_{1-\gamma/2}^*$ . (Efron ja Hastie, 2016, s. 185-187)

### 3.3 Monte Carlo -integrointi

Yleinen ongelma bayesiläisessä analyysissä on posteriorijakauman parametrien odotusarvojen laskeminen. Jatkuvilla satunnaismuuttujilla odotusarvon lauseke on integraalimuotoinen. Integraalien analyttinen laskeminen on hankalaa jo muutamissa ulottuvuuksissa, joten niiden laskemisessa täytyy turvautua numeerisiin menetelmiin. Olkoon  $X$  satunnaismuuttuja, jonka tiheysfunktio on  $p_X$ . Muunnoksen  $g(X)$  odotusarvo määritellään

$$\mathbb{E}(g(X)) = \int_{\mathcal{X}} g(x)p_X(x)dx, \quad (33)$$

jossa  $\mathcal{X}$  on tiheysfunktion  $p_X$  kantaja. Monte Carlo -integroinnin periaate on tuottaa otos  $(X_1, X_2, \dots, X_n)$   $X$ :n jakaumasta ja approksimoida odotusarvoa

$$\mathbb{E}(g(X)) \approx \overline{g(X)}_n = \frac{1}{n} \sum_{i=1}^n g(X_i). \quad (34)$$

Suurten lukujen lain avulla voidaan osoittaa, että otoskeskiarvo suppenee vähintään stokastisesti kohti odotusarvoa, kun otoskoko kasvaa rajatta. Myös estimaattorin  $\overline{g(X)}_n$  varianssi voidaan estimoida otoksesta kaavalla

$$\text{var}(\overline{g(X)}_n) \approx \frac{1}{n-1} \sum_{i=1}^n (g(X_i) - \overline{g(X)}_n)^2. \quad (35)$$

(Robert ja Casella, 2009, s. 65)

### 3.4 Markovin ketju Monte Carlo -menetelmät (MCMC)

MCMC-menetelmät perustuvat siihen, että ne tuottavat otoksen parametrien posteriorijakaumasta. Tämä otos on matemaattisessa mielessä Markovin ketju, jossa alkiot riippuvat toisistaan. Vaikka riippuva otos on epäinformatiivisempi kuin riippumaton otos, MCMC-menetelmistä on niin paljon hyötyjä, että ne ovat erittäin suosittuja bayesiläisessä analyysissä. Ensinnäkään posteriorijakauman muodolle ei ole juuri rajoituksia. Lisäksi hankalatkin moniulotteiset ongelmat on helppo pilkkoa pienempiin osiin.

Markovin ketju on sellainen jono satunnaismuuttujia

$$X^{(0)}, X^{(1)}, \dots, X^{(t)}, \dots, \quad (36)$$

että  $X^{(t)}$ :n edellisistä satunnaismuuttujista ehdollinen todennäköisyysjakauma riippuu vain edellisestä satunnaismuuttujasta  $X^{(t-1)}$ . Ehdollista todennäköisyysjakaumaa

$$X^{(t+1)} | X^{(0)}, X^{(1)}, \dots, X^{(t)} \sim K(X^{(t)}, X^{(t+1)}) \quad (37)$$

kutsutaan ketjun siirtymäksi. Esimerkiksi satunnaiskulkuprosessi

$$X^{(t+1)} = X^{(t)} + \epsilon_t, \quad \epsilon_t \sim N(0, 1), \quad (38)$$



jossa virhetermit  $\epsilon_t$  ovat toisistaan riippumattomia, on Markovin ketju, jonka siirtymä on  $N(X^{(t)}, 1)$ . Yleisesti ottaen MCMC-menetelmissä toteutuu stationäärisyyssehto, eli jonon peräkkäiset alkiot ovat realisaatioita samasta jakaumasta. Tämä tarkoittaa sitä, että olipa aloituspiste  $X^{(0)}$  mikä tahansa, jonon  $\{X^{(t)}\}$  jakauma suppenee toivottua kohdejakaumaa, kun  $t$  kasvaa rajatta. Siispä satunnaismuuttujan  $X$  odotusarvon estimaattorina voidaan käyttää otoskeskiarvoa

$$\mathbb{E}(X) \approx \frac{1}{T} \sum_{t=1}^T X^{(t)}. \quad (39)$$

(Robert ja Casella, 2009, s. 168–169)

### 3.5 Metropolisin–Hastingsin algoritmi

Metropolisin–Hastingsin algoritmi on ehkäpä käytetyin MCMC-menetelmä. Siinä simuloidaan otos posteriorijakaumasta, jolla on tiheysfunktio  $p$ , simuloimalla satunnaislukuja jostain ehdotusjakaumasta, jolla on tiheysfunktio  $q$ , ja hyväksymällä ehdotusluku tietyllä todennäköisyydellä. Todennäköisyys riippuu tiheysfunktioiden arvoista. Oletetaan, että ollaan simuloitu jo ketju  $(X^{(0)}, X^{(1)}, \dots, X^{(t)})$ , ja ehdotusarvo on  $X'$ . Hyväksymistodennäköisyys on

$$\alpha_t = \min \left\{ \frac{p(X') q(X^{(t)}|X')}{p(X^{(t)}) q(X'|X^{(t)})}, 1 \right\}. \quad (40)$$

Todennäköisyydellä  $\alpha_t$  asetetaan ketjun seuraavaksi alkioksi  $X^{(t+1)} = X'$ , ja todennäköisyydellä  $1 - \alpha_t$  asetetaan  $X^{(t+1)} = X^{(t)}$ . Niinpä ketjussa voi olla useita samoja lukuja peräkkäin. On hyvä huomata, että jos ehdotusjakauma on symmetrinen edellisen hyväksytyyn alkion  $X^{(t)}$  suhteen,  $q(X^{(t)}|X') = q(X'|X^{(t)})$ , eli hyväksymistodennäköisyys on riippumaton ehdotusjakaumasta. Hyväksymisaste

$$\bar{\alpha} = \frac{1}{T} \sum_{t=1}^T \alpha_t \quad (41)$$

kuvaa sitä, kuinka pitkiä jaksoja samaa lukua ketjussa keskimäärin on. Hyväksymisastetta voi kontrolloida esimerkiksi ehdotusjakauman varianssia muuttamalla. (Robert ja Casella, 2009, s. 170–171)

Vaikka teoriassa ehdotusjakauman valinnalla ei ole juurikaan merkitystä, käytännössä se vaikuttaa huomattavasti algoritmin tehokkuuteen ja suppenemiseenopeuteen. Eräs luonnollinen ja paljon käytetty valinta ehdotusjakaumaksi on nk. satunnaiskulkuehdotusjakauma, joka on symmetrinen. Siinä ehdotusarvon  $X'$  lauseke on muotoa

$$X' = X^{(t)} + \epsilon_t, \quad (42)$$

jossa  $\epsilon_t$  on nollakeskinen satunnaisluku esimerkiksi tasa- tai normaalijakaumasta. Satunnaiskulkuehdotusjakauma yksinkertaistaa algoritmia siinäkin mielessä, että ehdotusjakaumaan liittyvä osa hyväksymistodennäköisyyden lausekkeessa (40) supistuu pois. (Robert ja Casella, 2009, s. 182)

### 3.6 Gibbsin otanta

Gibbsin otanta on hieman erityylinen MCMC-menetelmä kuin Metropolisin–Hastingsin algoritmi. Sen vahvuus piilee siinä, että sen avulla voidaan pilkkoa hankala moniulotteinen ongelma pienempiin, helpommin ratkaistaviin palasiin. Oletetaan, että satunnaismuuttuja  $\mathbf{X}$  koostuu useista komponenteista  $(X_1, X_2, \dots, X_p)$ . Yksittäinen komponentti voi olla yksi- tai moniulotteinen. Gibbsin otanta perustuu siihen, että jokainen  $\mathbf{X}$ :n komponentti ehdollistetaan erikseen kaikille muille komponenteille:  $X_i | X_1, X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_p$ . Oletetaan, että kaikkien näiden ehdollisten satunnaismuuttujien jakaumista osataan simuloida satunnaislukuja. Näitä jakaumia kutsutaan täysehdollisiksi jakaumiksi. Tämän jälkeen Gibbsin otanta -algoritmin rakentaminen on helppoa. Simuloidaan vuorotellen uusi arvo jokaiselle komponentille  $X_i$  omasta

täysehdoollisesta jakaumastaan niin, että kaikkien muiden komponenttien arvo on kiinnitetty niiden edelliseen simuloituun arvoon. Kun algoritmia on suoritettu tarpeeksi monta kierrosta, simuloitujen arvojen approksimoivat  $\mathbf{X}$ :n yhteisjakaumaa. (Robert ja Casella, 2009, s. 199–206)

Joskus täysehdoollisten jakaumien lausekkeita voi olla vaikeaa tai mahdotonta johtaa. Tällöin voidaan yhdistää Gibbsin otanta- ja Metropolisin–Hastingsin -algoritmit. Jos jostain täysehdoollisesta jakaumasta ei osata suoraan simuloida satunnaislukuja, komponentti  $X_i$  voidaan simuloida jostakin ehdotusjakaumasta ja laskea sille hyväksymistodennäköisyys kuten Metropolisin–Hastingsin algoritmilla. Jos hyväksymistodennäköisyyden lausekkeessa esiintyy jotain muita vektorin  $\mathbf{X}$  komponentteja, niiden arvot kiinnitetään viimeisimpään simuloituun arvoon.

### 3.7 Elliptinen viipaleotanta (Elliptical slice sampling)

Elliptinen viipaleotanta on melko uusi MCMC-menetelmä, jonka on huomattu toimivan hyvin bayesiläisissä piilomuuttujamalleissa, joissa käytetään multinormaaliprioria. Olkoon  $\mathbf{X} = (X_1, X_2, \dots, X_n)$  vektori piiloparametreja ja  $\mathbf{y}$  havaintoaineisto. Vektorille  $\mathbf{X}$  asetetaan gaussinen prior, jonka tiheysfunktio on

$$p(\mathbf{X}|\mathbf{0}, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp \left\{ -\frac{1}{2} \mathbf{X}^T \Sigma^{-1} \mathbf{X} \right\}. \quad (43)$$

Normalisoimaton posterioritiheys on

$$p(\mathbf{X}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{X})p(\mathbf{X}|\mathbf{0}, \Sigma), \quad (44)$$

jossa on  $p(\mathbf{y}|\mathbf{X})$  on havainnoille oletetun mallin uskottavuusfunktio ja jossa oletetaan, että  $\Sigma$  on tunnettu. Menetelmä lähtee ajatuksesta, että Metro-

lisiin–Hastingsin ehdotusarvona käytetään vektoria

$$\mathbf{X}' = \sqrt{1 - \rho^2} \mathbf{X}^{(t)} + \rho \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}), \quad (45)$$

jossa  $\rho \in [-1, 1]$  on askelkokoparametri. Kun  $\rho = 1$ , ehdotus on itse asiassa otos priorijakaumasta, kun taas  $\rho$ :n lähestyessä nollaa se on lähellä edellisellä iteraatiolla hyväksyttyä vektoria. Tässä ehdotusarvossa ongelmana on kuitenkin se, että askelkoko täytyy säätää aina uudelleen sopivaksi eri tilanteisiin. Olisi toivottavaa, että askelkoon valinta olisi jotenkin automatisoitu. (Murray ym. 2010)

Kun askelkoon annetaan vaihdella välillä  $[-1, 1]$ , ehdotusarvon lauseke (45) määrää puoliellipsin. Luonnollisempi parametointi ehdotusarvolle on

$$\mathbf{X}' = \mathbf{X}^{(t)} \cos \rho + \mathbf{v} \sin \rho, \quad (46)$$

joka määrää kokonaisen ellipsin. Askelkokoa siirretään automaattisesti jokaisella iteraatiolla kohti nollaa, kunnes tarpeeksi hyvä ehdotusarvo löytyy. (Murray ym. 2010)

Elliptisen viipaleotannan algoritmi on seuraavanlainen iteraatiolla  $t$ :

1. simuloi  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$
2. laske  $\kappa = \log p(\mathbf{y}|\mathbf{X}^{(t)}) + \log u$ ,  $u \sim \mathcal{U}(0, 1)$
3. simuloi  $\rho \sim \mathcal{U}(0, 2\pi)$
4. aseta  $(\rho_{\min}, \rho_{\max}) = (\rho - 2\pi, \rho)$
5. laske  $\mathbf{X}' = \mathbf{X}^{(t)} \cos \rho + \mathbf{v} \sin \rho$
6. jos  $\log p(\mathbf{y}|\mathbf{X}') > \kappa$ , aseta  $\mathbf{X}^{(t+1)} = \mathbf{X}'$

7. muutoin:
8. jos  $\rho < 0$ , aseta  $\rho_{\min} = \rho$ , muutoin aseta  $\rho_{\max} = \rho$
9. simuloi  $\rho \sim \mathcal{U}(\rho_{\min}, \rho_{\max})$
10. palaa kohtaan 5

Algoritmin yksi iteraatio tuottaa yhden realisaation parametrivektorin  $\mathbf{X}$  posteriorijakaumasta. (Murray ym. 2010). Vektorin  $\mathbf{v}$  tuottaminen onnistuu helposti. Jos  $\Sigma = \mathbf{L}\mathbf{L}^T$ , jossa  $\mathbf{L}$  on Choleskyn hajotelman yläkolmiomatriisi, niin  $\mathbf{v} = \mathbf{L}\mathbf{d}$ , jossa  $\mathbf{d} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

### 3.8 MCMC-algoritmien suppenemisen tarkastelu

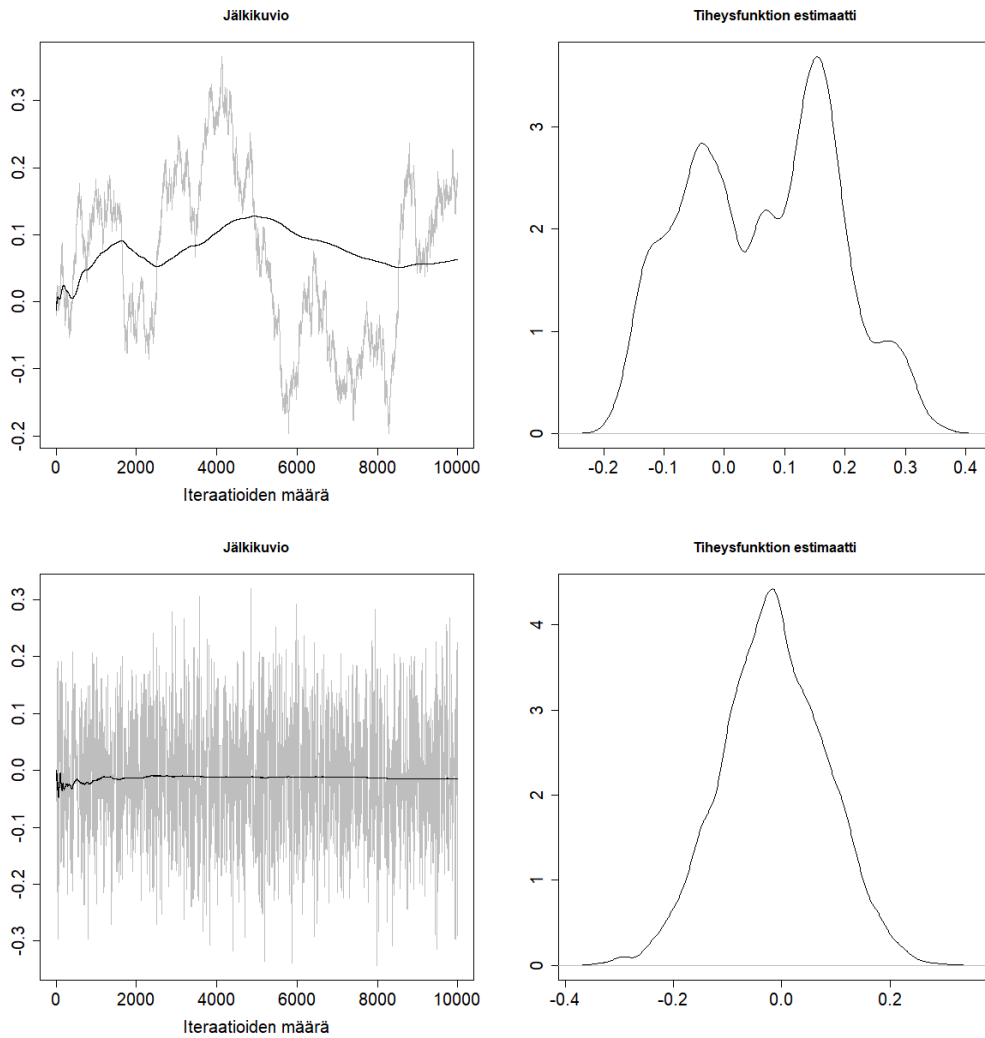
Koska MCMC-menetelmät perustuvat satunnaislukujen simulointiin, on perusteltua kysyä kuinka paljon satunnaislukuja on tarpeeksi. Liian pieni määrä ei edusta posteriorijakaumaa tarpeeksi hyvin ja tietyn pisteen jälkeen simuloinnin jatkaminen ei enää merkittävästi muuta estimaatteja.

Tärkein kriteeri MCMC-algoritmin suppenemiselle on se, että tuotetun ketjun jakauman tulisi olla sama kuin posteriorijakauma. On kuitenkin vaikeaa todentaa että näin on, koska posteriorijakauman muotoa ei välttämättä edes tunneta. On kuitenkin mahdollista tutkia, kuinka riippumaton tai riippuva ketju on aloituspisteestä eri ketjujen realisaatioita vertaamalla. Eräs yksinkertainen tapa on verrata eri ketjujen keskiarvoja. Niiden tulisi luonnollisesti olla lähellä toisiaan. (Robert ja Casella, 2009, s. 238–240)

Suosittu tapa tutkia suppenemista on piirtää kuvia parametrien Markovin ketjuista. Eräs tällainen kuva on nk. jälkikuvio, jossa on erotettu viivalla parametrin arvo kullakin iteraatiolla. Jälkikuvion pitäisi näyttää sta-

tionääriseltä, eli sen pitäisi vaihdella tasaisesti saman keskiarvon ympärillä ja vaihteluvälin tulisi olla suurin piirtein vakio. (Robert ja Casella, 2009, s. 242–244). Kuvassa 2 on esitetty jälkikuvio ja tiheysfunktion ydinestimaatti Gaussin ytimellä kahdelle eri Markovin ketjulle. Ylärivillä näkyy huonolaatuinen jälkikuvio, jossa kumulatiivinen keskiarvo ei vakiinnu. Lisäksi tiheysfunktio on lievästi kaksihuippuinen. Alarivillä on hyvälaatuinen jälkikuvio. Erot kuvissa saatiin aikaan Metropolisin–Hastingsin algoritmin ehdotusjakauman varianssia muuttamalla.

MCMC-algoritmeissa on myös tyypillistä jättää osa Markovin ketjun alusta käyttämättä analyysissä, koska voidaan ajatella, että ketju ei ole vielä ehtinyt päästä kohdejakauman todennäköisimmälle alueelle. Tätä alkuosaa kutsutaan burn-in jaksoksi.



Kuva 2: Kahden eri Markovin ketjun jälkikuvio harmaalla, kumulatiivinen keskiarvo mustalla viivalla ja tiheysfunktion estimaatti. Molemmissa on estimoitu normaalijakauman odotusarvoa samasta aineistosta, joka on simuloitu standardinormaalijakaumasta.

## 4 Aineistot ja niiden analyysi

Tässä luvussa esitellään aineistot, joihin sovelletaan alaluvuissa 2.6 ja 2.7 esiteltyjä menetelmiä. Lisäksi esitellään MCMC-algoritmit, joiden avulla tilastollinen päättely tehdään.

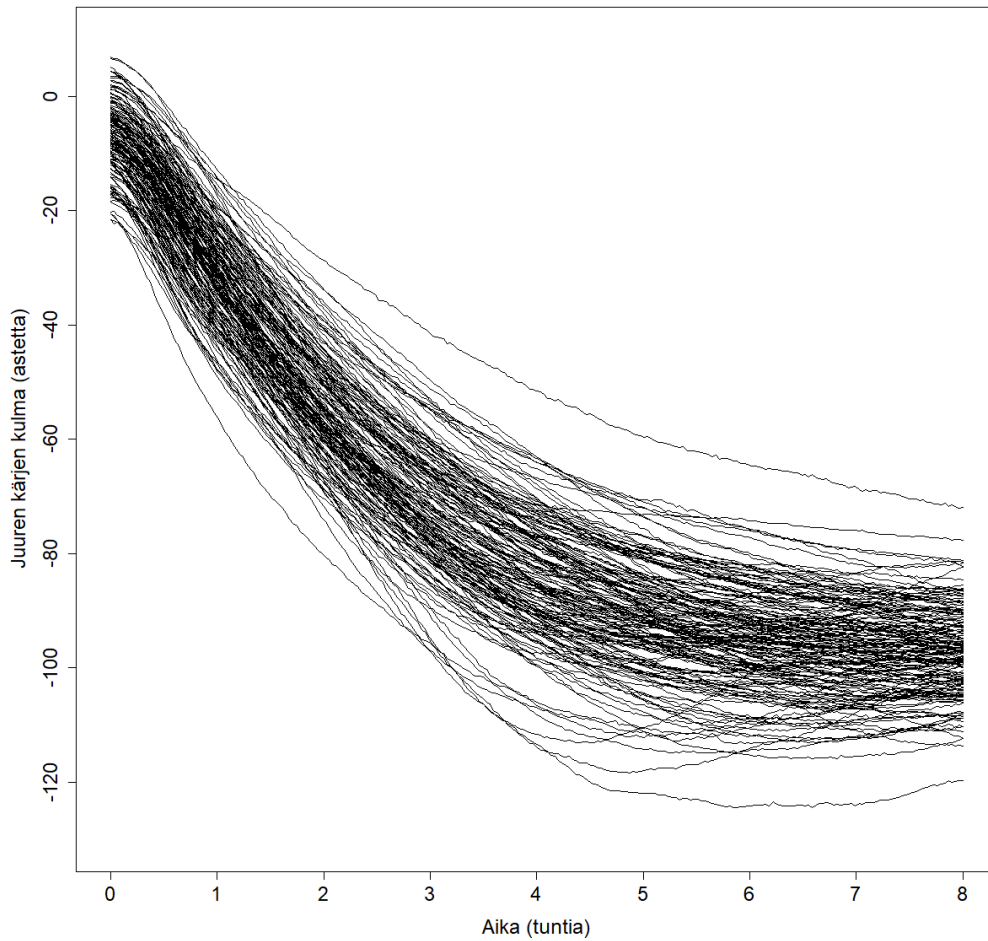
### 4.1 Simuloitu aineisto

Simuloitu aineisto luotiin sen takia, että estimointimenetelmiä voitaisiin testata ja tuloksia verrata aineiston generoineen mallin parametrien arvoihin. Simuloidussa aineistossa oli  $N = 100$  yksilöä, joilla kullakin oli  $T = 50$  havaintoa tasavälein välillä  $[0, 24]$  tuntia. Aineistoon simuloitiin ajasta riippuva geneettinen varianssi  $\sigma_G^2(t) = \cos \frac{2\pi}{24}t + 2$  ja ympäristövarianssi  $\sigma_E^2(t) = \sin \frac{2\pi}{24}t + 2$ . Aineisto voisi vertautua oikeassa elämässä mittauksiin jostakin vuorokauden- tai vuodenajasta syklisesti riippuvasta fenotyypistä. Aikapistekohtainen kovarianssimatriisi oli  $\mathbf{K}(t) = \sigma_E^2(t)\mathbf{I} + \sigma_G^2(t)\mathbf{A}$ , jossa sukulaisuusmatriisin  $\mathbf{A}$  alkio  $A_{ij} = \frac{1}{2^{|i-j|}}$ ,  $i, j = 1, \dots, N$ . Matriisille  $\mathbf{K}(t)$  tehtiin Choleskyn hajotelma  $\mathbf{K}(t) = \mathbf{L}(t)\mathbf{L}^T(t)$ . Vektorit  $\mathbf{y}(t)$  simuloitiin jakaumasta  $\mathcal{N}(\mathbf{0}, \mathbf{K}(t))$  niin, että  $\mathbf{y}(t) = \mathbf{L}(t)\mathbf{d}$ , jossa  $\mathbf{d} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .

### 4.2 Todellinen aineisto

Toinen analysoitava aineisto koostuu lituruohon (*Arabidopsis thaliana*) siemenmittauksista. Siinä on mitattu 162 lituruohon siemenistä kasvavan juuren päähän syntyvän kulman suuruutta. Siemenet olivat petrimaljoissa, jotka asetettiin niin, että niiden kulma maahan nähden oli 90 astetta. Näin painovoima vaikutti kasvavaan juureen kohtisuorasti. Juuria kuvattiin kahdeksan tunnin ajan kahden minuutin välein, joten aineistossa on yhteensä 162 aikasarjaa, joissa kussakin on 241 aikapistettä. Koska aikapisteitä oli näin paljon, valittiin analyysiin joka kuudes eli 40 nopeuttamaan estimointia. Kulmat mi-





Kuva 3: Kaikkien yksittäisten siementen havaintoaikasarjat visualisoituna.

tattiin automaattisesti tietokoneen avulla. Mittaustapaa kuvailee tarkemmin Moore ym. (2013). Aineisto sisältää myös genotyypitiedot 234 eri markerista, joten yksilöiden välinen sukulaisuusmatriisi  $\mathbf{A}$  on myös konstruoitavissa.

Ilmiötä, jossa kasvin juuret tai versot muuttavat kasvusuuntaansa paino-

voimavektoriin nähden, kutsutaan gravitropismiksi. Kasvin kyky aistia painovoima mahdollistaa versojen kasvun valoa kohti, jolloin yhteyttämistä tapahtuu enemmän, ja juurten kasvun syvemmälle maahan, jolloin kasvi on tukevammin pystyssä ja saa paremmin vettä ja mineraaleja. (Masson ym. 2002). Muun muassa tähän ja moneen muuhun tutkimukseen hyödynnetään lituruohoa, joka on erittäin suosittu kasvi biologien keskuudessa. On huomattu, että sen avulla voidaan tutkia monien muidenkin kasvien biologista rakennetta ja toimintaa. Lituruohon hyviä puolia ovat muun muassa sen nopea kasvu ja pieni genomi, joka mahdollistaa yksityiskohtaisen molekulaarisen analyysin. (Meinke ym. 1998)

### 4.3 Algoritmit

Aineistojen analyysiin kaksivaiheisella menetelmällä käytettiin algoritmia 1 (liite B), joka on muunneltu versio Monterrubio-Gómez ym. (2018) algoritmista 1. Algoritmissa hyödynnetään Gibbsin otantaa, eli siinä simuloidaan yksi parametri tai parametrien osajoukko kerrallaan niin, että muiden parametrien arvot on kiinnitetty. Simulaatioita ei kuitenkaan tehdä parametrien täysehdoollisista jakaumista, vaan sileän funktion  $f$  arvot simuloidaan käyttämällä elliptistä viipaleotantaa ja skalaariparametrit  $\sigma^2$  ja  $l$  simuloidaan Metropolisin–Hastingsin algoritmilla. Ehdotusjakaumana on satunnaiskulkunormaalijakauma, jonka keskipiste on parametrin edellinen hyväksytty arvo. Jakauman varianssia mukautetaan niin, että ehdotusparametrien hyväksymisaste on noin 0.44. Algoritmia ajettiin 150 000 MCMC-iteraatiota, jossa burn-in -jakso oli 50 000 iteraatiota. Posteriorijakaumista laskettiin parametrien posteriorikeskiarvot ja algoritmin suppenemista tarkasteltiin Markovin ketjujen jälkikuvioista.

Algoritmi 2 (liite B) on yksivaiheisen menetelmän analyysialgoritmi. Se

on hyvin samankaltainen kuin kaksivaiheisessa menetelmässä. Se on Gibbs-henkinen siinä mielessä, että se simuloi toisen varianssikomponenttivektorin kun toinen on kiinnitetty edelliseen hyväksytyyn arvoonsa. Vektoreiden simulointiin käytetään elliptistä viipaleotantaa. Algoritmia ajettiin 150 000 MCMC-iteraatiota, ja burn-in -jaksona oli ensimmäiset 50 000 iteraatiota. Posteriorijakaumista laskettiin parametrien posteriorikeskiarvot ja niiden 95 %:n posteriorivälit. Algoritmin suppenemista tutkittiin Markovin ketjujen jälkikuvioista. Molemmat esiteltyt algoritmit on implementoitu C++-ohjelmointikielellä.

## 5 Tulokset

Tässä luvussa esitellään analyysien tulokset, verrataan menetelmien estimaatteja ja tutkitaan algoritmien suppenemista Markovin ketjujen jälkikuviosta. Lisäksi siemenaineistosta saatavia tuloksia verrataan kirjallisuudessa aikaisemmin esitettyihin tuloksiin.

### 5.1 Simuloitu aineisto

Kuvassa 4 on rinnakkain kaksivaiheisen ja yksivaiheisen menetelmän estimaatit simuloitun aineiston varianssikomponentti- ja heritabiliteettifunktioille. Kuvista huomaa selkeän eron kaksivaiheisen menetelmän luottamusvälien ja yksivaiheisen menetelmän posteriorivälien välillä: luottamusvälit ovat huomattavasti leveämpiä kuin posteriorivälit. Ero johtuu siitä, että yksivaiheinen menetelmä ottaa kaikkien aikapisteiden informaation huomioon priorin kautta, jolloin estimoinnin epävarmuus pienenee. Välit ovat kuitenkin määritelmällisesti eri asia, joten niitä ei voi täysin verrata toisiinsa. Näyttää siltä, että näin pienen aineiston perusteella heritabiliteettifunktiosta ei voi sanoa kaksivaiheisen menetelmän perusteella juuri mitään varmaa: luottamusväli täyttää melkein koko parametrin arvoalueen. Simuloitusta aineistosta estimoituja funktioita voidaan verrata oikeisiin funktioihin, koska ne tiedetään. Kaksivaiheisen menetelmän estimaatit näyttävät seuraavan parametrien oikeita arvoja huonommin kuin yksivaiheisen menetelmän estimaatit.

Kaksivaiheisen menetelmän Markovin ketjujen jälkikuviot (liite C) näyttävät kauttaaltaan hyvin käyttäytyviltä: kumulatiivinen keskiarvo vakiintuu nopeasti ja vaihteluväli pysyy samana. Yksivaiheisen menetelmän jälkikuviotkin näyttävät hyväksyttäviltä, mutta niistä huomaa että algoritmia olisi

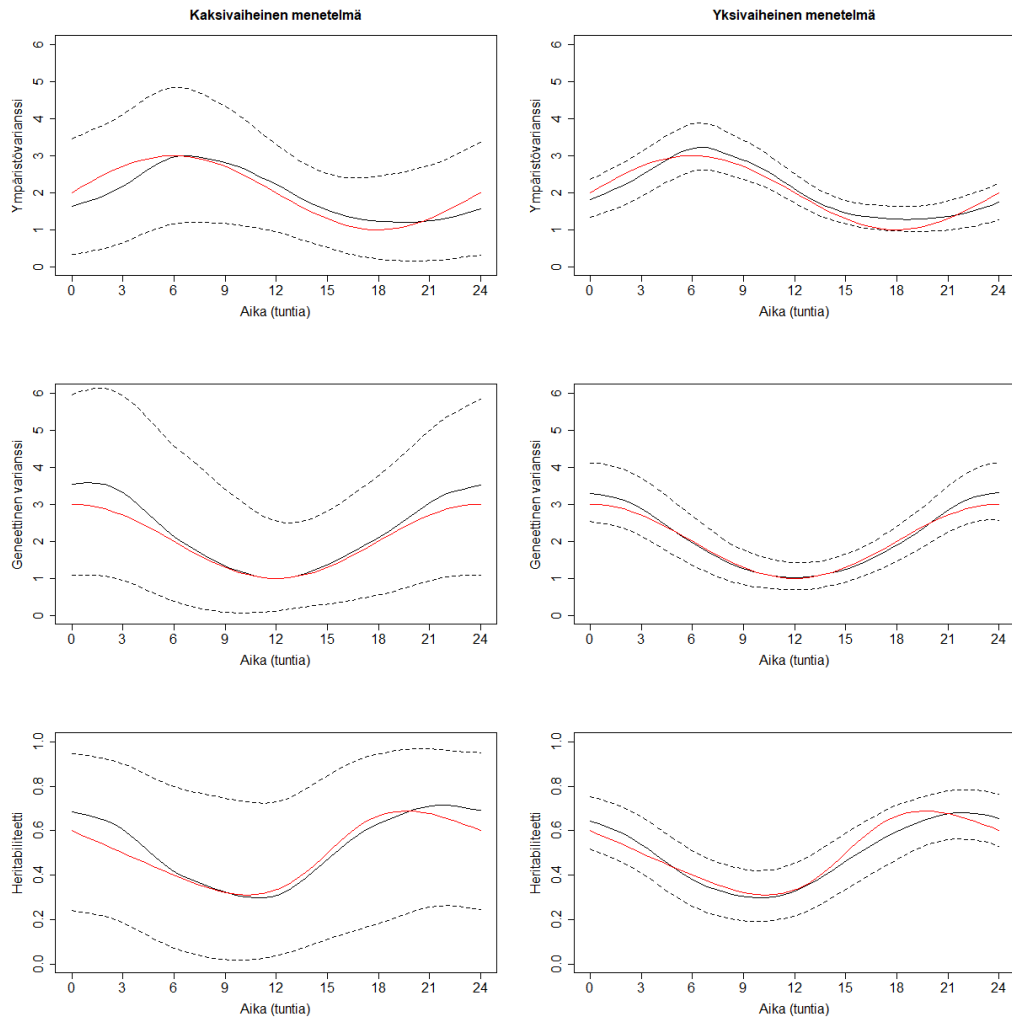
voinut ajaa kauemminkin.

## 5.2 Todellinen aineisto

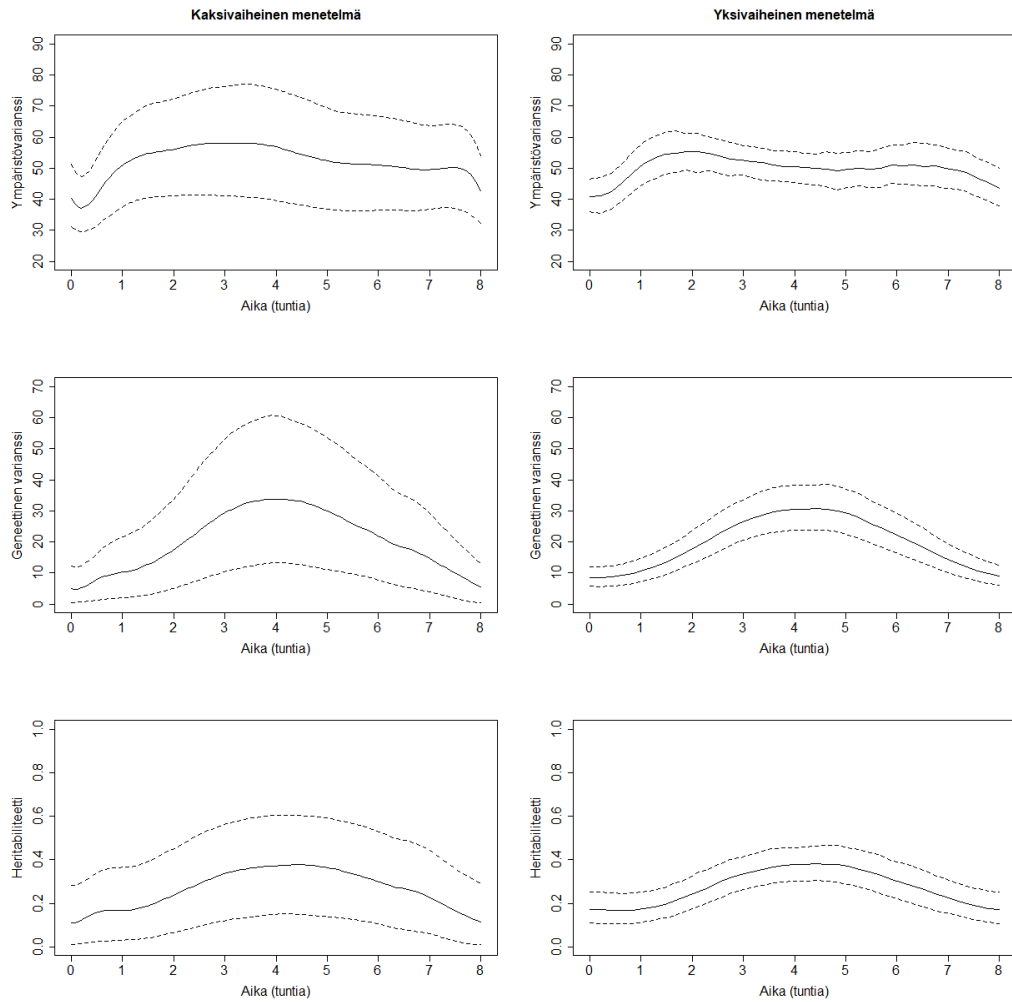
Kuvassa 5 on rinnakkain kaksivaiheisen ja yksivaiheisen menetelmän estimaatit siemenaineiston varianssikomponentti- ja heritabiliteettifunktioille. Luottamusvälien ja posteriorivälien ero on jälleen huomattava, mutta simuloitua aineistoa suuremman yksilömäärän ansiosta kaksivaiheisen menetelmän heritabiliteettiestimaatti on täsmällisempi. Keskiarvofunktioissa ei menetelmien välillä ole suurta eroa. Ympäristövarienssifunktio näyttää melko tasaiselta, kun taas geneettinen varianssi on koholla kokeen keskivaiheilla. Tämä aiheuttaa kummun myös heritabiliteettifunktioon.

Kaksivaiheisen menetelmän jälkikuviot (liite C) näyttävät hyvin käyttäytyviltä. Yksivaiheisen menetelmän jälkikuvioista huomaa jälleen, että MCMC-iteraatioiden määrä olisi voinut olla suurempi.

Samaa siemenaineistoa on analysoitu aikaisemminkin kirjallisuudessa, mutta pääpainona ei ole ollut heritabiliteettifunktion estimointi. Funktiosta on kuitenkin piirretty kuvia, joihin tämän tutkielman tuloksia voi verrata. Muun muassa Vanhatalo ym. (2019) on analysoinut samaa aineistoa. Siinä heritabiliteettifunktion estimaatti on hyvin lähellä tämän työn estimaatteja. Tosin siinä funktiolle ei ole estimoitu posteriori- tai luottamusvälejä. Myös Moore ym. (2013) on analysoinut siemenaineistoa. Siinä varianssikomponentti- ja heritabiliteettifunktiot ovat samanmuotoisia kuin tässä tutkielmassa, mutta niiden saamat arvot ovat erikoisella tavalla aivan eri luokkaa. Tässäkään artikkelissa ei ole posteriori- tai luottamusvälejä funktioille.



Kuva 4: Simuloidusta aineistosta eri menetelmillä estimoidut ympäristövarianssi, geneettinen varianssi ja heritabiliteetti mustalla viivalla ja niiden 95 %:n luottamusvälit tai posteriorivälit katkoviivalla. Parametrien oikeat arvot on piirretty punaisella viivalla.



Kuva 5: Todellisesta aineistosta eri menetelmillä estimoidut ympäristövarianssi, geneettinen varianssi ja heritabiliteetti mustalla viivalla ja niiden 95 %:n luottamusvälit tai posteriorivälit katkoviivalla.

## 6 Pohdinta

Tässä työssä esiteltiin kaksi bayesiläistä menetelmää funktiomuotoisen heritabiliteetin estimoimiseen pitkittäisaineistosta: kaksivaiheinen menetelmä ja yksivaiheinen menetelmä. Menetelmät ovat siinä mielessä joustavia, että ne eivät aseta juurikaan rajoituksia heritabiliteettifunktion muodolle. Sen tulee olla ainoastaan riittävän sileä. Lisäksi menetelmät yleistyvät kaikkiin sellaisiin tilanteisiin, joissa biologisen populaation sukulaisuusmatriisi tiedetään, oli se sitten konstruoitu molekulaarisen markkeriaineiston tai sukupuun perusteella. Yksivaiheisen menetelmän vahvuuksiin vaikuttaa kuuluvan myös estimoinnin täsmällisyys. Algoritmeja testattaessa huomattiin, että heritabiliteettiestimaatti vain yhdelle aikapisteelle, riippumatta muista pisteistä, on melko epätasainen. Kun taas heritabiliteetti estimoidaan useille aikapisteille yhtä aikaa, niin priorin kautta tuleva lisäinformaatio helpottaa tehtävää huomattavasti. Ilmiön voi huomata vertaamalla menetelmien luottamus- ja posteriorivälejä tulokset-luvussa. Kaksivaiheisen menetelmän vahvuus taas on sen hyvä skaalautuvuus. Sillä saa isoillekin aineistoille estimaatteja kohtuullisessa ajassa. Tässä työssä analysoiduilla aineistoilla kaksivaiheisella menetelmällä kesti noin tunnin saada tuotettua kaikki estimaatit. Yksivaiheisella menetelmällä siinä kesti simuloidulle aineistolle noin vuorokausi ja todelliselle aineistolle noin kolme vuorokautta.

Hyvä jatkotutkimuksen aihe voisi olla yksivaiheisen menetelmän kehittäminen nopeammaksi. Suurin pullonkaula algoritmissa on  $N \times N$  -kokoisen matriisin determinantin laskeminen ja kääntäminen, joka tehdään joka MCMC-iteraatiolla  $T$  kertaa. Matriisi on yleisemmässä muodossa  $a\mathbf{A} + b\mathbf{I}$ , jossa pienet kirjaimet ovat skalaareja ja isot matriiseja. Kaavassa päivittyy aina toinen skalaarikomponentti, kun toinen pysyy samana. Voisi yrittää selvittää, onko o.o. matriisin determinantille ja käänteismatriisille olemassa käyttökelpoista



päivityskaavaa, joka nopeuttaisi laskentaa.

Toinen kehityskohde on kiinnitettävien hyperparametrien arvojen valinta. Tällä hetkellä ne täytyy valita kokeilemalla, johon voi kulua paljon aikaa. Yksi vaihtoehto on niiden estimointi muiden parametrien tapaan, mutta käytännössä niiden identifioituvuus havaintoaineistosta voi osoittautua ongelmalliseksi.

## Viitteet

- Abadir, Karim M. ja Magnus, Jan R. (2005). *Matrix Algebra*. Cambridge University Press, Cambridge.
- Bryois, Julien, Buil, Alfonso, Ferreira, Pedro G., Panousis, Nikolaos I., Brown, Andrew A., Viñuela, Ana, Planchon, Alexandra, Bielser, Deborah, Small, Kerrin, Spector, Tim ja Dermitzakis, Emmanouil T. (2017). Time-dependent genetic effects on gene expression implicate aging processes. *Genome Research* 27: 542–552.
- Efron, Bradley ja Hastie, Trevor (2016). *Computer Age Statistical Inference: Algorithms, Evidence, and Data Science*. Institute of Mathematical Statistics Monographs. Cambridge University Press, Cambridge.
- Endelman, J. B. (2011). Ridge regression and other kernels for genomic selection with R package rrBLUP. *Plant Genome* 4: 250–255.
- Fahrmeir, Ludwig ja Kneib, Thomas (2011). *Bayesian Smoothing and Regression for Longitudinal, Spatial and Event History Data*. Oxford University Press, Oxford.
- Gelman, Andrew, Carlin, John B., Stern, Hal S., Dunson, David B., Vehtari, Aki ja Rubin, Donald B. (2013). *Bayesian Data Analysis*. CRC Press LLC, Boca Raton.
- Kwak, Il-Youp, Moore, Candace R., Spalding, Edgar P. ja Broman, Karl W. (2014). A Simple Regression-Based Method to Map Quantitative Trait Loci Underlying Function-Valued Phenotypes. *Genetics* 197 (4): 1409–1416.
- Kang, Hyun Min, Zaitlen, Noah A., Wade, Claire M., Kirby, Andrew, Heckerman, David, Daly, Mark J. ja Eskin, Eleazar (2008). Efficient control

- of population structure in model organism association mapping. *Genetics* 178 (3): 1709–1723.
- Li, Zitong ja Sillanpää, Mikko J. (2015). Dynamic Quantitative Trait Locus Analysis of Plant Phenomic Data. *Trends in Plant Science* 20 (12): 822–833.
- Masson, Patrick H., Tasaka, Masao, Morita, Miyo T., Guan, Changhui, Chen, Rujin ja Boonsirichai, Kanokporn (2002). Arabidopsis thaliana: A Model for the Study of Root and Shoot Gravitropism. *The Arabidopsis book / American Society of Plant Biologists* 1: e0043.
- Meinke, David W., Cherry, J. Michael, Dean, Caroline, Rounsley, Steven D. ja Koornneef, Maarten (1998). Arabidopsis thaliana: A Model Plant for Genome Analysis. *Science* 282: 662, 679–82.
- Monterrubio-Gómez, Karla, Roininen, Lassi, Wade, Sara, Damoulas, Theo ja Girolami, Mark (2018). Posterior Inference for Sparse Hierarchical Non-stationary Models. *arXiv e-prints*, arXiv:1804.01431: arXiv:1804.01431. arXiv: 1804.01431 [stat.CO].
- Moore, Candace R., Johnson, Logan S., Kwak, Il-Youp, Livny, Miron, Broman, Karl W. ja Spalding, Edgar P. (2013). High-Throughput Computer Vision Introduces the Time Axis to a Quantitative Trait Map of a Plant Growth Response. *Genetics* 195 (3): 1077–1086.
- Murray, Iain, Adams, Ryan Prescott ja MacKay, David J.C. (2010). Elliptical Slice Sampling. *Proceedings of Machine Learning Research* 9: 541–548.
- Pletcher, Scott D. ja Geyer, Charler J. (1999). The Genetic Analysis of Age-Dependent Traits: Modeling the Character Process. *Genetics* 153: 825–835.

- Rasmussen, Carl Edward ja Williams, Christopher K. I. (2006). *Gaussian Processes for Machine Learning*. MIT Press, Cambridge.
- Robert, Christian P. ja Casella, George (2009). *Introducing Monte Carlo Methods with R*. Springer, New York.
- Sillanpää, Mikko J. (2011). On statistical methods for estimating heritability in wild populations. *Molecular Ecology* 20: 1324–1332.
- Vanhatalo, Jarno, Li, Zitong ja Sillanpää, Mikko J. (2019). A Gaussian process model and Bayesian variable selection for mapping function-valued quantitative traits with incomplete phenotypic data. *Bioinformatics*. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btz164. eprint: <http://oup.prod.sis.lan/bioinformatics/advance-article-pdf/doi/10.1093/bioinformatics/btz164/28355391/btz164.pdf>. URL: <https://doi.org/10.1093/bioinformatics/btz164>.
- VanRaden, Paul M. (2008). Efficient Methods to Compute Genomic Predictions. *Journal of Dairy Science* 91 (11): 4414–4423.
- Visscher, Peter M., Hill, William G. ja Wray, Naomi R. (2008). Heritability in the genomics era – concepts and misconceptions. *Nature Reviews Genetics* 9: 255–266.

## A Kroneckerin matriisitulo

Tässä työssä tilastollisen mallin määrittelyssä tarvitaan Kroneckerin tuloa. Hieman harvinaisempaan operaatioon sitä on syytä avata. Kahden reaaliarvoisen matriisin  $\mathbf{A}_{n \times m}$  ja  $\mathbf{B}_{p \times q}$  Kroneckerin tulo määritellään

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \dots & a_{1m}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \dots & a_{2m}\mathbf{B} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1}\mathbf{B} & a_{n2}\mathbf{B} & \dots & a_{nm}\mathbf{B} \end{bmatrix}. \quad (47)$$

Tulomatriisi on kooltaan  $np \times mq$ . Kroneckerin tulolle on johdettu monia ominaisuuksia, joista todettakoon tässä työssä tärkeät

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{B})^{-1} &= \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}, \\ |\mathbf{A} \otimes \mathbf{B}| &= |\mathbf{A}|^n |\mathbf{B}|^m, \quad \text{kun } \mathbf{A} \in \mathbb{R}^{n \times n} \text{ ja } \mathbf{B} \in \mathbb{R}^{m \times m}. \end{aligned} \quad (48)$$

(Abadir ja Magnus, 2005, s. 273-279)

## B MCMC-algoritmit

---

### Algoritmi 1

---

**Alkuarvot:**  $\mathbf{f}^{(1)}, \sigma^{2(1)}, l^{(1)}$

- 1:  $t = 1:(T - 1)$
  - 2: simuloi  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$
  - 3: simuloi  $u \sim \mathcal{U}(0, 1)$
  - 4: laske  $\kappa = \log p_{\mathcal{N}}(\mathbf{y}|\mathbf{f}^{(t)}, \sigma^{2(t)}) + \log u$
  - 5: simuloi  $\rho \sim \mathcal{U}(0, 2\pi)$
  - 6: aseta  $(\rho_{\min}, \rho_{\max}) = (\rho - 2\pi, \rho)$
  - 7: laske  $\mathbf{f}' = \mathbf{f}^{(t)} \cos \rho + \mathbf{v} \sin \rho$
  - 8: **jos**  $p_{\mathcal{N}}(\mathbf{y}|\mathbf{f}', \sigma^{2(t)}) > \kappa$  **niin**
  - 9:     asetta  $\mathbf{f}^{(t+1)} = \mathbf{f}'$
  - 10: **muutoin**
  - 11:     **jos**  $\rho < 0$  **niin**
  - 12:         asetta  $\rho_{\min} = \rho$
  - 13:     **muutoin**
  - 14:         asetta  $\rho_{\max} = \rho$
  - 15:     **lopeta jos**
  - 16:         simuloi  $\rho \sim \mathcal{U}(\rho_{\min}, \rho_{\max})$
  - 17:         palaa kohtaan 7.
  - 18:     **lopeta jos**
  - 19:         simuloi  $\log \sigma^{2'} \sim N(\log \sigma^{2(t)}, s_1)$
  - 20:         laske  $\alpha_{\sigma^2} = \min \left\{ 1, \frac{p_{\mathcal{N}}(\mathbf{y}|\mathbf{f}^{(t+1)}, \sigma^{2'}) p_{u;\log \sigma^2}(\log \sigma^{2'})}{p_{\mathcal{N}}(\mathbf{y}|\mathbf{f}^{(t+1)}, \sigma^{2(t)}) p_{u;\log \sigma^2}(\log \sigma^{2(t)})} \right\}$
  - 21:         todennäköisyydellä  $\alpha_{\sigma^2}$  aseta  $\log \sigma^{2(t+1)} = \log \sigma^{2'}$ , muutoin aseta  
         $\log \sigma^{2(t+1)} = \log \sigma^{2(t)}$
  - 22:         mukauta  $s_1$
  - 23:         simuloi  $\log l' \sim N(\log l^{(t)}, s_2)$
  - 24:         laske  $\alpha_l = \min \left\{ 1, \frac{p_{\mathcal{N}}(\mathbf{f}^{(t+1)}|\mathbf{0}, \mathbf{C}_l) p_{u;\log l}(\log l')}{p_{\mathcal{N}}(\mathbf{f}^{(t+1)}|\mathbf{0}, \mathbf{C}_{l^{(t)}}) p_{u;\log l}(\log l^{(t)})} \right\}$
  - 25:         todennäköisyydellä  $\alpha_l$  aseta  $\log l^{(t+1)} = \log l'$ , muutoin aseta  
         $\log l^{(t+1)} = \log l^{(t)}$
  - 26:         mukauta  $s_2$
  - 27: **lopeta**
-

---

**Algoritmi 2**

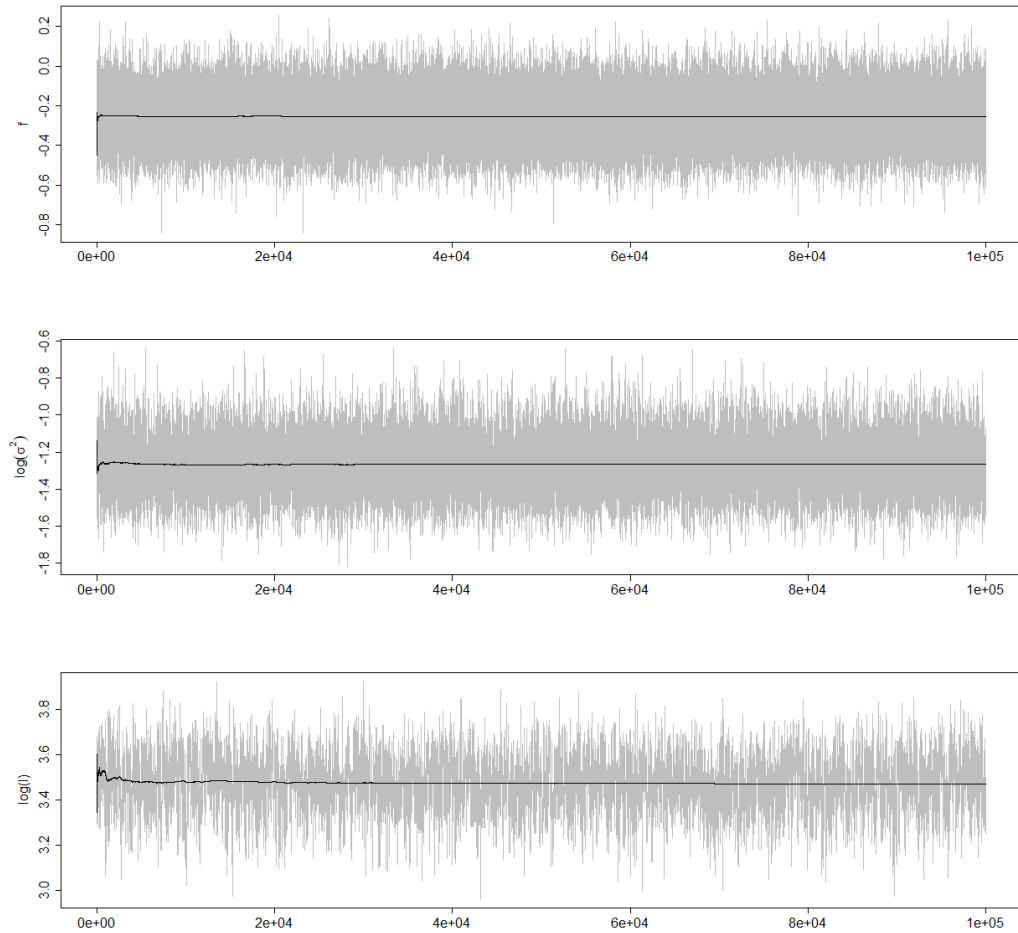
---

**Alkuarvot:**  $\sigma_G^{2(1)}, \sigma_E^{2(1)}$

```
1: t = 1:(T - 1)
2:   simuloi  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_G)$ 
3:   simuloi  $u \sim \mathcal{U}(0, 1)$ 
4:   laske  $\kappa = \log p_{\mathcal{N}}(\tilde{\mathbf{y}}_T | \mathbf{0}, \mathbf{K}_T) + \log u$ 
5:   simuloi  $\rho \sim \mathcal{U}(0, 2\pi)$ 
6:   aseta  $(\rho_{\min}, \rho_{\max}) = (\rho - 2\pi, \rho)$ 
7:   laske  $\sigma_G^{2'} = \sigma_G^{2(t)} \cos \rho + \mathbf{v} \sin \rho$ 
8:   päivitä  $\mathbf{K}_T$ 
9:   jos  $\log p_{\mathcal{N}}(\tilde{\mathbf{y}}_T | \mathbf{0}, \mathbf{K}_T) > \kappa$  niin
10:     aseta  $\sigma_G^{2(t+1)} = \sigma_G^{2'}$ 
11:   muutoin
12:     jos  $\rho < 0$  niin
13:       aseta  $\rho_{\min} = \rho$ 
14:     muutoin
15:       aseta  $\rho_{\max} = \rho$ 
16:     lopeta jos
17:     simuloi  $\rho \sim \mathcal{U}(\rho_{\min}, \rho_{\max})$ 
18:     palaa kohtaan 7.
19:   lopeta jos
20:   simuloi  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_E)$ 
21:   simuloi  $u \sim \mathcal{U}(0, 1)$ 
22:   laske  $\kappa = \log p_{\mathcal{N}}(\tilde{\mathbf{y}}_T | \mathbf{0}, \mathbf{K}_T) + \log u$ 
23:   simuloi  $\rho \sim \mathcal{U}(0, 2\pi)$ 
24:   aseta  $(\rho_{\min}, \rho_{\max}) = (\rho - 2\pi, \rho)$ 
25:   laske  $\sigma_E^{2'} = \sigma_E^{2(t)} \cos \rho + \mathbf{v} \sin \rho$ 
26:   päivitä  $\mathbf{K}_T$ 
27:   jos  $\log p_{\mathcal{N}}(\tilde{\mathbf{y}}_T | \mathbf{0}, \mathbf{K}_T) > \kappa$  niin
28:     aseta  $\sigma_E^{2(t+1)} = \sigma_E^{2'}$ 
29:   muutoin
30:     jos  $\rho < 0$  niin
31:       aseta  $\rho_{\min} = \rho$ 
32:     muutoin
33:       aseta  $\rho_{\max} = \rho$ 
34:     lopeta jos
35:     simuloi  $\rho \sim \mathcal{U}(\rho_{\min}, \rho_{\max})$ 
36:     palaa kohtaan 25.
37:   lopeta jos
38: lopeta
```

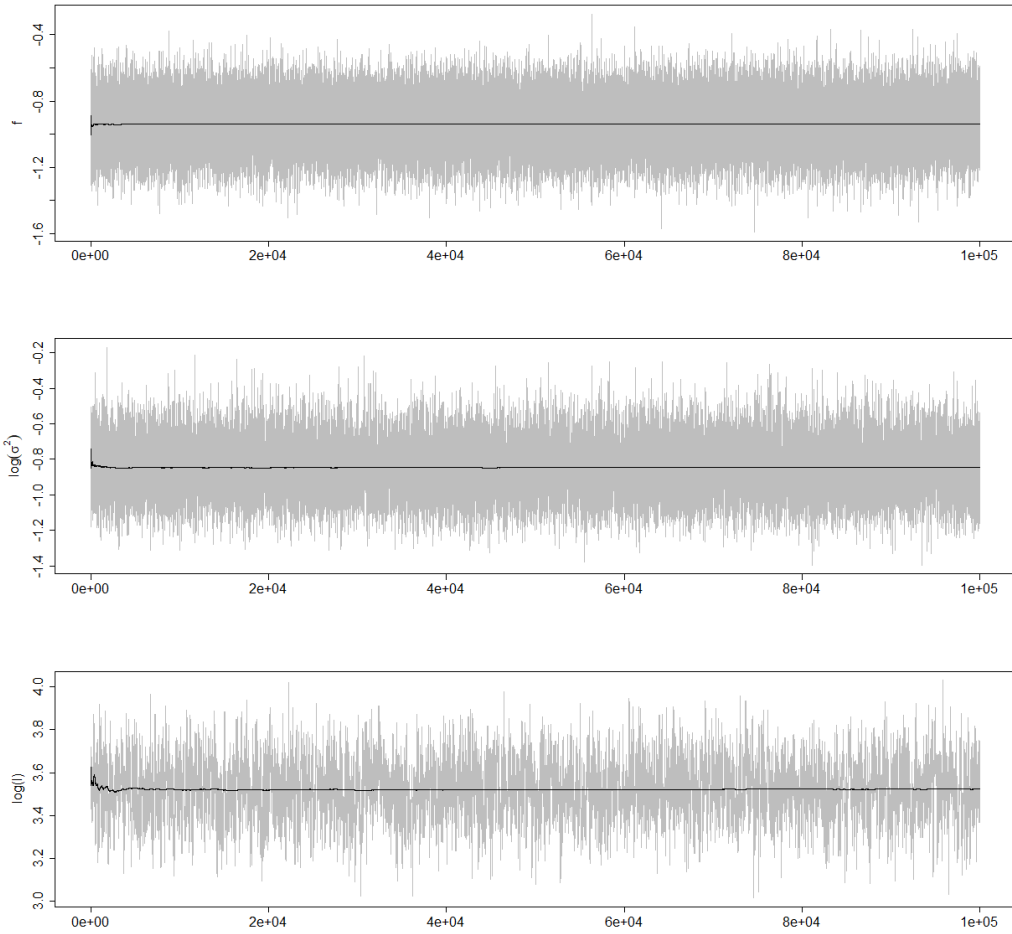
---

## C Jälkikuviot

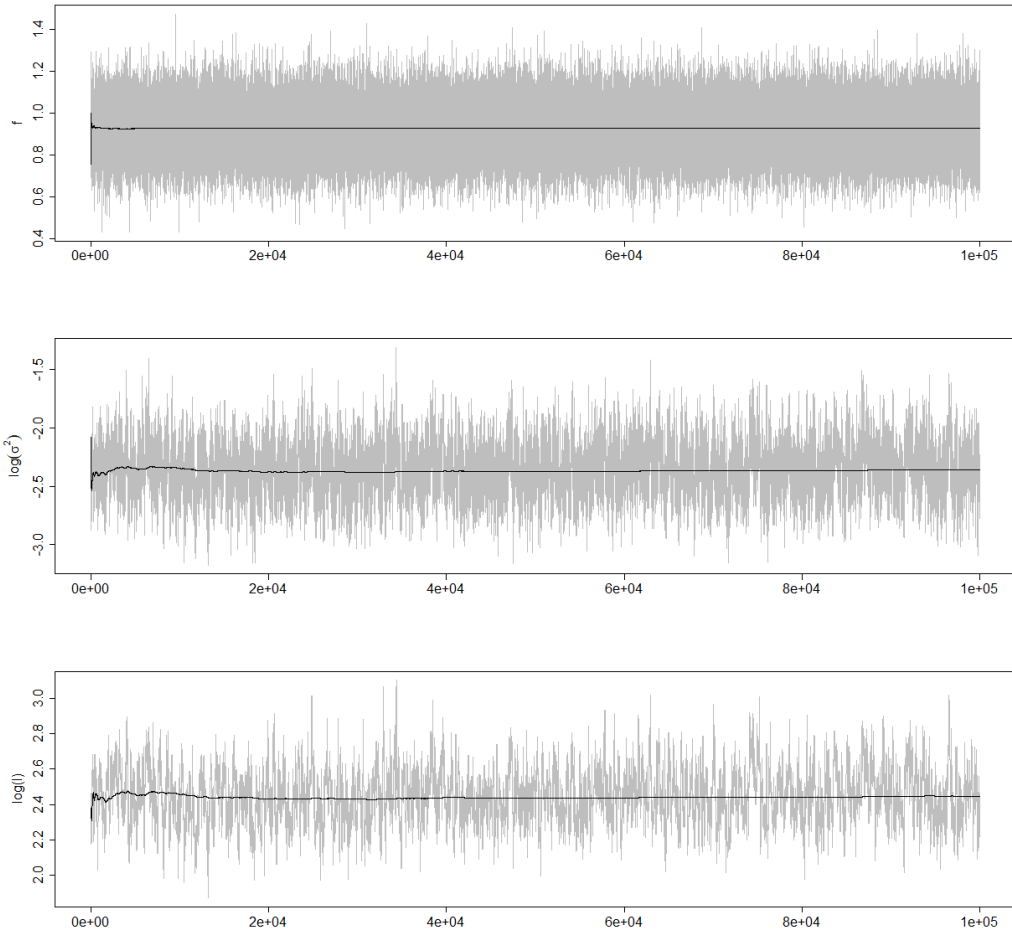


Kuva 6: Kaksivaiheisen menetelmän simuloitun aineiston ympäristövarianssin silotusmallin parametrien Markovin ketjujen jälkikuviot ja kumulatiiviset keskiarvot. Funktion  $f$  komponenttiparametri on satunnaisesti valittu.

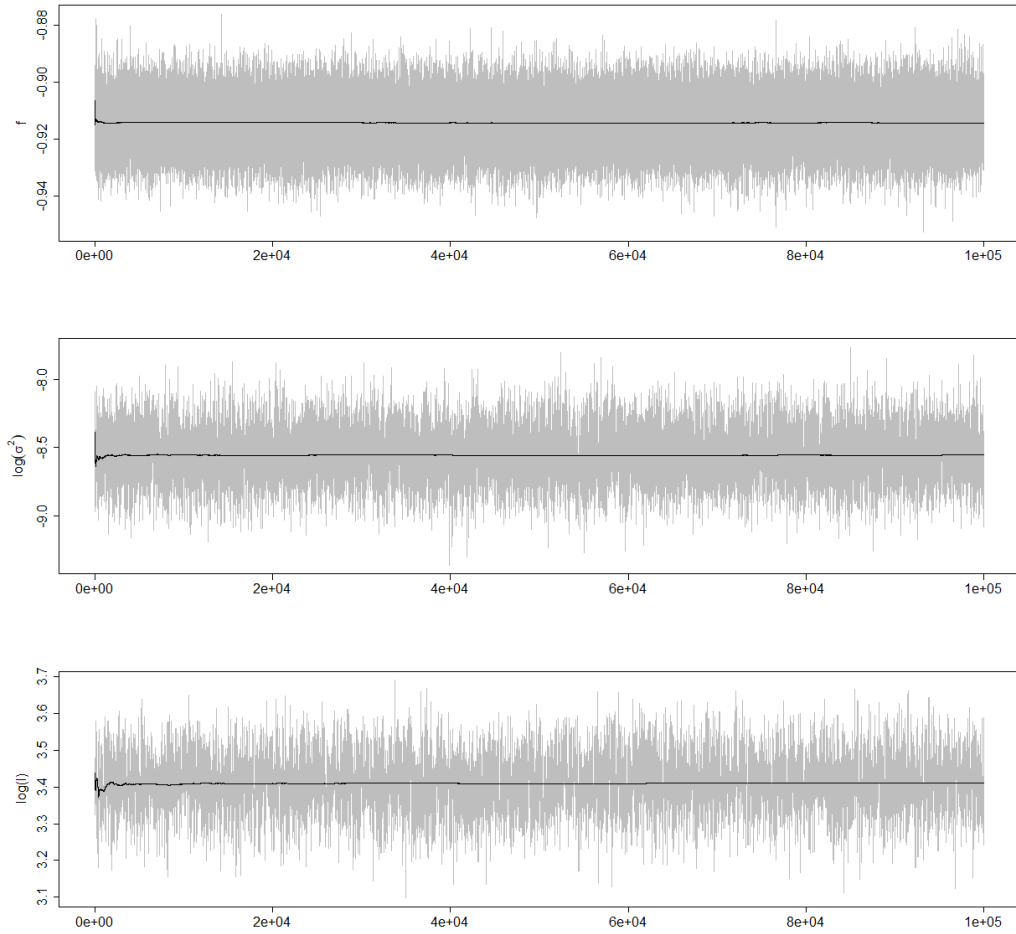




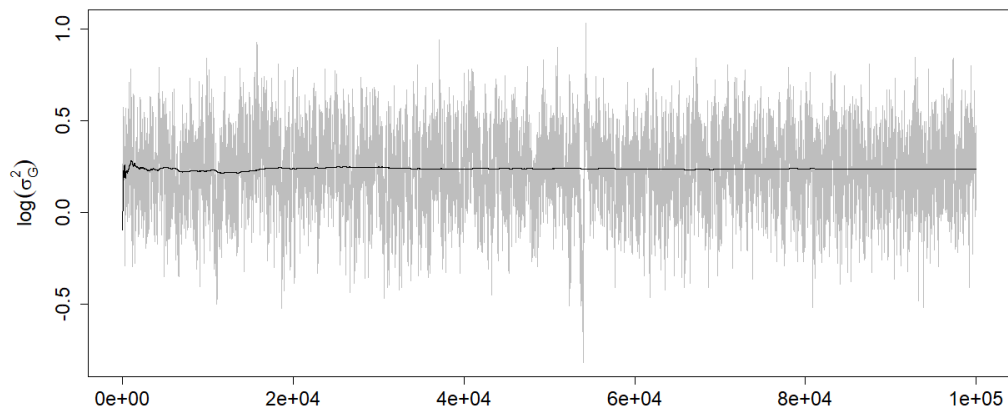
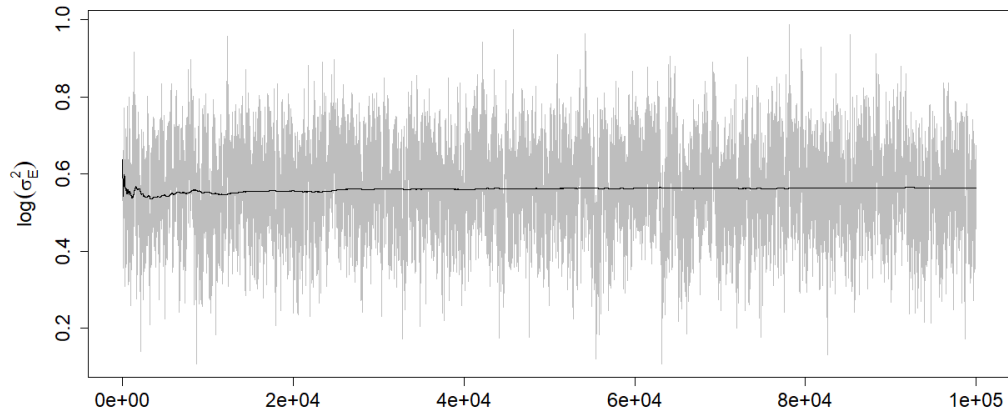
Kuva 7: Kaksivaiheisen menetelmän simuloidun aineiston geneettisen variaanssin silotusmallin parametrien Markovin ketjujen jälkikuviot ja kumulatiiviset keskiarvot. Funktion  $f$  komponenttiparametri on satunnaisesti vaillu.



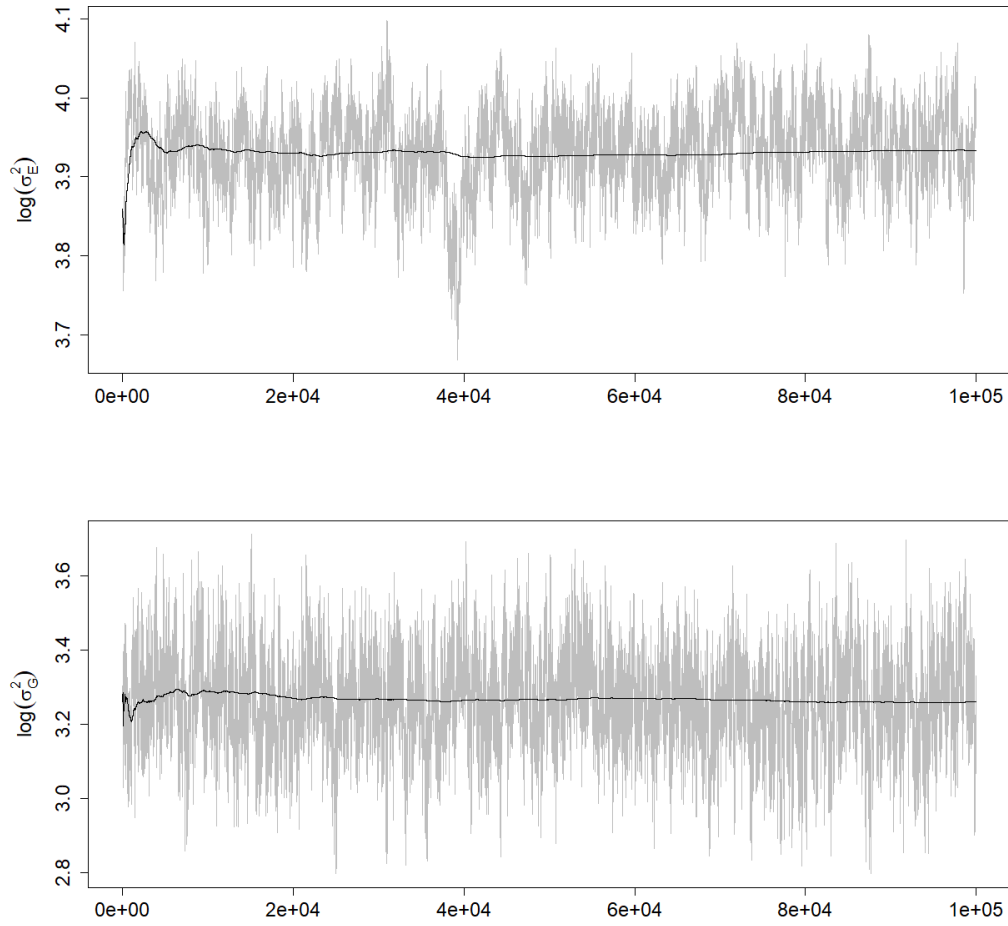
Kuva 8: Kaksivaiheisen menetelmän siemenaineiston ympäristövariانسin silotusmallin parametrien Markovin ketjujen jälkikuviot ja kumulatiiviset keskiarvot. Funktion  $f$  komponenttiparametri on satunnaisesti valittu.



Kuva 9: Kaksivaiheisen menetelmän siemenaineiston geneettisen varianssin silotusmallin parametrien Markovin ketjujen jälkikuviot ja kumulatiiviset keskiarvot. Funktion  $f$  komponenttiparametri on satunnaisesti valittu.



Kuva 10: Yksivaiheisen menetelmän simuloitun aineiston ympäristö- ja geneettisen varianssin satunnaisen komponenttiparametrin Markovin ketjujen jälkikuviot ja kumulatiiviset keskiarvot.



Kuva 11: Yksivaiheisen menetelmän siemenaineiston ympäristö- ja geneettisen varianssin satunnaisen komponenttiparametrin Markovin ketjujen jälkikuvio ja kumulatiiviset keskiarvot.