

Antti Ilmari Penttilä

**Partikkelien muodon mallintaminen niiden  
2D-satunnaisprojektioista**

Pro gradu

4.11.2002

Tilastotieteen laitos

Helsingin yliopisto

Tiedekunta-Facultet-Faculty Valtiotieteellinen tiedekunta		Laitos-Institution-Department Tilastotieteen laitos	
Tekijä-Författare-Author Penttilä, Antti			
Työn nimi-Arbetets titel-Title Partikkelien muodon mallintaminen niiden 2D-satunnaisprojektiosta			
Oppiaine - Läroämne - Subject Tilastotiede			
Työn laji-Arbetets art-Level Pro gradu	Aika-Datum-Month and year 2002-11-04	Sivumäärä - Sidantal - Number of pages 53	
<p>Tiivistelmä - Referat - Abstract  Tutkielman tarkoituksena on muodostaa geometrinen muotomalli boorikarbidipartikkeleille (B4C), estimoida mallin parametrit partikkeleista otetuista kuvista, ja verrata mallin tuottamaa lineaarista polarisaatiota B4C-partikkelien mikrogravitaatiossa mitattuun polarisaatioon.</p> <p>B4C on yksi ranskalaisen PROGRA2-tutkimusryhmän mikrogravitaatiossa tutkimista partikkelityypeistä. Ryhmällä on käytössään polarisaation mittaukseen sopiva laitteisto parabolisiin lentoihin käytetyllä lentokoneella. Parabolisilla lennoilla koneen sisälle saadaan luotua lähes painottomat olosuhteet, jonka aikana polarisaatiomittaukset tehdään. Painovoima vaikuttaa partikkelien orientaatioon ja pakkaantumiseen, ja sitä kautta myös polarisaatioon. Tähtitieteessä mikrogravitaatiokohteita löytyy esimerkiksi tähtienvälisestä pölystä ja komeettojen pyrstöistä.</p> <p>Pienten partikkelien muotoa voidaan mallintaa muun muassa säännöllisillä muodoilla, vaikkapa ellipsoideilla tai sylintereillä, tai satunnaisesti deformoiduilla palloilla, kuten Gaussin palloilla. B4C-partikkelien muotomalliksi sopii kuitenkin paremmin satunnainen monitahokas. Tutkielmassa esitellään eräs sopiva malliproseduuri satunnaismonitahokaiden luomiseen. Mallissa on kaksi parametria, jotka estimoidaan partikkeleista otetusta kuvamateriaalista.</p> <p>Kuvamateriaalissa näkyy partikkelien 2D-satunnaisprojektiota. Kukin partikkeli on kuvattu vain yhdestä suunnasta, joten kuvista on mahdoton johtaa suoraan partikkelien kolmiulotteista muotoa. Kun partikkelien oletetaan kuitenkin noudattavan samaa muotomallia, voidaan kolmiulotteista muotoa estimoida tilastollisessa mielessä.</p> <p>Mallin realisaatioista voidaan myös ottaa satunnaisprojektiota, ja mitata samoja suureita kuin oikeista partikkeleista. Nämä suureet ovat satunnaismuuttujia, mutta muuttujien analyttisen jakauman johtaminen on hyvin vaikea tehtävä. Näin ollen mallin estimointiin ei voida käyttää suurimman uskottavuuden menetelmää. Malliproseduurin avulla saadaan kuitenkin simuloitua havaintoja tästä tuntemattomasta jakaumasta. Näistä havainnoista muodostettu ydinestimointi estimaatti tälle refraktioideksille vertaamalla mallin ja aitojen partikkelien polarisaatiokäyrien eroja refraktioidexsin reaali- ja imaginaariosien funktiona pienimmän neliösumman mielessä.</p> <p>Valonsirontatutkimuksessa halutaan usein arvioida sirottavan aineen ominaisuuksia sen valonsirontan perusteella. Kun ominaisuuksiin vaikuttaa kappaleen muoto, koko ja aineen refraktioidexsi, on inversion onnistumisen kannalta erittäin tärkeää, että kappaleen muotomalli on realistinen ja hyvin estimoitu. Tutkielmassa esiteltyä simuloitun uskottavuuden menetelmää voidaan käyttää erilaisten muotomallien estimointiin. Lisäksi menetelmää voidaan käyttää myös muissa estimointiongelmassa sovellusalasta riippumatta.</p>			
Avainsanat Nyckelord Keywords  pienien partikkelien mallintaminen simuloitu suurimman uskottavuuden menetelmä ydinestimointi valonsironta polarisaatio			
Säilytyspaikka - Förvaringsställe - Where deposited /(täytetään kirjastossa)			
Muita tietoja - Övriga uppgifter - Additional information			

# Sisältö

<b>1 Johdanto</b>	<b>4</b>
1.1 Valonsironta ja työn tarkoitus . . . . .	5
1.2 Työn rakenne . . . . .	7
<b>2 PROGRA<sup>2</sup>-ohjelman partikkelien mallinnus</b>	<b>8</b>
2.1 PROGRA <sup>2</sup> -ohjelma . . . . .	8
2.2 Teoreettinen malli PROGRA <sup>2</sup> -ohjelman partikkeleille . . . . .	10
2.2.1 Mallin konstruointi . . . . .	11
<b>3 Mallin parametrien estimointi</b>	<b>14</b>
3.1 B <sub>4</sub> C-partikkeleista saatu havaintoaineisto . . . . .	14
3.2 Parametrien estimointi havaintoaineistosta . . . . .	15
3.2.1 Suurimman uskottavuuden menetelmä . . . . .	16
3.2.2 Simuloitu suurimman uskottavuuden menetelmä . . . . .	17
3.2.3 Tiheysfunktion ydinestimointi . . . . .	18
3.2.4 Malliparametrien estimointi B <sub>4</sub> C-partikkeleille . . . . .	20
<b>4 Mallin verifiointi polarisaatiomittausten avulla</b>	<b>30</b>
4.1 Polarisaatiomittaukset . . . . .	31
4.2 Valonsirontasimulaatiot . . . . .	32
4.2.1 Simuloinnin toteutus . . . . .	32
4.2.2 Refraktioindeksin määrittäminen . . . . .	33
<b>5 Päätelmät</b>	<b>38</b>
<b>Lähteet</b>	<b>40</b>

<b>Liitteet</b>	<b>43</b>
<b>A Tasainen jakauma pallokoordinaateissa</b>	<b>43</b>
<b>B Lähdekoodi monitahokasmallin realisaatioiden arpomiseen</b>	<b>45</b>
<b>C Rosenblatin ehtojen todistus ydinestimaatille</b>	<b>46</b>
<b>D Simuloidun suurimman uskotavuuden estimoinnin toteutus</b>	<b>48</b>

# 1 Johdanto

Tässä työssä tarkastellaan pienten partikkelien muodon mallintamista hyödyntämällä informaatiota niiden 2D-satunnaisprojektioista. Pieniksi partikkeleiksi kutsutaan läpimitaltaan muutamasta kymmenestä nanometristä aina noin sataan mikrometriin vaihtelevia hiukkasia. Näitä hiukkasia esiintyy joka puolella, muun muassa maapallon ilmakehässä, tähtienvälisessä pölyssä, teollisuuden pigmentteinä paperin päällysteessä tai maalissa jne.

Lumme ym. (1995) on ehdottanut, että partikkelimuodot voitaisiin luokitella kolmeen pääluokkaan: satunnaisiin monitahokkaisiin, satunnaisesti deformaatioihin palloihin ja satunnaisagregaatteihin. Monitahokkaat ovat tasomaisista tahkoista koostuvia kappaleita, joissa tahkojen välillä on terävä kulma. Pallomaisten kappaleiden pinta taas vaihtelee sileästi, käyttökelpoisimpana mallina voisi pitää ns. Gaussin palloa, jossa pallon säde on lognormaalisti jakautunut satunnaisprosessi, jonka määrittelevät säteen odotusarvo, hajonta ja autokorrelaatiofunktio. Satunnaisagregaatit koostuvat useasta yhteen liittyneestä partikkelista, vaikkapa samankokoisista palloista.

Partikkeleista saatava kuvamateriaali on useimmiten 2D-projektioita partikkelin aidosta kolmiulotteisesta muodosta. Fotogrammetria ja stereoskopia ovat aloja, joissa tutkitaan kolmiulotteisen informaation rakentamista projektioista. Lääketieteessäkin tärkeä Radon-muunnos tarjoaa tähän mahdollisuudet, jos kappaleesta on saatavilla kuvamateriaalia useista suunnista (katso esim. Deans 1983). Pieniä partikkeleita kuvataan usein TEM (transmission electron microscope) ja SEM (scanning electron microscope) -laitteistoilla, joilla ei valitettavasti tällä hetkellä ole mahdollista saada kuvia useista suunnista. Kuitenkin mallipohjaista lähestymistapaa käyttäen voidaan myös yhdestä suunnasta otettuja projektioita hyödyntää tilastollisessa mielessä, jos kuvattuja, samaa muotomallia noudattavia partikkeleita on useita.

Koska tämä työ sivuaa vahvasti myös valonsirontaa, käydään sen ymmärtämiseksi tarvittavat käsitteet ensin läpi johdannossa, jonka jälkeen johdan-

nossa kuvaillaan lyhyesti työn tilastotieteellisen osan rakennetta.

## 1.1 Valonsironta ja työn tarkoitus

Pienten partikkelien muodon mallintaminen on tarpeellinen apuväline laajemmalle fysiikan tutkimusalueelle, valonsironnalle. Valonsironnalla tarkoitetaan valon reagoimista väliaineeseen, ja itse asiassa kaikki visuaaliset havainnot ovat valonsironnan tulosta. Kuitenkin varsinaisesta valonsirontatutkimuksesta puhutaan silloin, kun tutkittava kohde, joka reagoi valon kanssa, on joko niin pieni tai niin kaukana, ettei sen muotoa varsinaisesti nähdä. Havaitsejan ja kohteen välinen avaruuskulma on siis hyvin pieni. Kohteesta sironnutta valoa voidaan kuitenkin havaita ja tutkia.

Valonsirontatutkimus on sovellutusalueiltaan hyvin laaja. Tähtitieteessä tutkimus liittyy niin kosmisen pölyn ja komeettojen pyrstöjen tutkimiseen kuin planeettojen ja asteroidien tutkimiseen. Geofysiikassa, meteorologiassa ja kaukokartoituksessa tutkitaan valonsirontaa ilmakehän hiukkasista, pilvistä ja maan pinnalla olevista kappaleista, vaikkapa metsästä. Esimerkiksi maapallon ilmaston lämpenemisen kannalta on tärkeä tietää ilmakehän partikkeleiden kyky heijastaa lämpösäteilyä takaisin maapallolle. Tutkatekniikassa ja sotilassovelluksissa tarvitaan myös tietoa säteilyn sironnasta, vaikka säteily voi olla myös ääniaaltoja kuten sukellusveneen tunnistamiseen tarkoitetuissa tutkissa. Paperi- ja pigmenttiteollisuudessa päällysteen valonsirontaominaisuudet määrittävät pitkälle tuotteen laadun, esimerkiksi paperin valkoisuuden.

Valonsironnassa sekä kohteeseen tulevan että siitä lähtevän valon ominaisuuksia kuvaa Stokesin vektori, jossa on neljä komponenttia. Komponentit liittyvät valon intensiteettiin, lineaariseen polarisaatioon (kaksi komponenttia) ja ympyräpolarisaatioon. Sirottavan kappaleen ominaisuudet voidaan kuvata täydellisesti  $4 \times 4$  -sironnamatriisilla eli Muellerin matriisilla, jolla kerrotaan tulevan valon Stokesin vektori. Matriisin elementit riippuvat sirottavan kappaleen koosta, muodosta ja materiaalista (katso esim. Bohren

ym. 1983, luku 3).

Kappaleen absoluuttinen koko ei ole tärkeä valonsirontaongelmissa vaan koon suhde siihen tulevan valon aallonpituuteen. Koko ilmoitetaankin mieluummin ns. kokoparametrina,  $\frac{2\pi\langle r \rangle}{\lambda}$ , jossa  $\langle r \rangle$  on kappaleen keskimääräinen säde ja  $\lambda$  on valon aallonpituus. Erittäin epäsäännöllisten kappaleiden tapauksessa keskimääräinen säde ei välttämättä ole hyvä indikaattori koolle. Sen sijaan käytetään joskus sellaisen pallon sädettä, jonka tilavuus on sama kuin tarkasteltavan kappaleen tilavuus. Materiaalin valonsirontaominaisuuksia taas kuvaa sen taitekerroin eli refraktioindeksi  $m = n + in'$ , jossa reaaliosa  $n$  on valon taittumiseen ja imaginaariosa  $n'$  valon absorptioon vaikuttava osa (Muinonen 1986, luku 1). Kappaleen muoto vaikuttaa sen valonsirontaominaisuuksiin, etenkin polarisaatio on sille herkkä (Lumme ym. 1998). Siksi muodon mallintaminen on tärkeä osa valonsirontatutkimusta.

Valonsirontaongelmia ei pystytä ratkaisemaan yleisesti analyyttisessä muodossa, ja sopivan numeerisen menetelmän valinta riippuu vahvasti sekä kappaleen kokoparametrin että muodosta. Pallon tapauksessa on analyyttinen ratkaisu olemassa (Mie-teoria), ja viime aikoina on myös huomattavasti monimutkaisempi palloklusterin tapaus saatu ratkaistua (Mackowski ym. 1996). Hyvin pienille kappaleille (kokoparametri paljon alle yhden) pätee Rayleigh'n sironta, johon mm. taivaan sininen väri perustuu. Jos kokoparametri on noin sata tai suurempi, voidaan valonsirontaa approksimoida lukion fysiikasta tutulla geometrisella optiikalla, jossa valonsäteet heijastuvat peilimäisesti kappaleen pinnoilta. Kun kokoparametri on pienempi kuin sata, puhutaan pienistä partikkeleista, ja näiden valonsirontaan on vasta viime aikoina kehitetty numeerisia ratkaisuhjelmia. Pienten partikkeleiden sironta on kuitenkin hyvin laskentaintensiivinen ongelma. Kirjallisuutta valonsironta-algoritmeista tarjoavat esim. Mishchenko ym. (2000), Lumme ym. (1995) ja Muinonen (1986).

Valonsirontaongelmaa halutaan usein tutkia myös inversio-ongelmana, jolloin sironneen valon ominaisuuksista pyritään päättelemään jotain sirottavasta kappaleesta. Kun sirontaan kuitenkin vaikuttavat kaikki kolme tekijää,

muoto, koko ja refraktiaindeksi, on inversio käytännössä vaikeaa. Aineiden refraktiaindeksejä ei useinkaan tunneta tarkasti, etenkin indeksin imaginaariosaa on vaikea määrittää. Kokojakaumalle on usein joitain arvioita, mutta erityisesti muodon mallintaminen on usein hyvin paljon vain arvauksien varassa. Jos useampi kuin yksi näistä kolmesta tekijästä on epävarma, on inversio-ongelma erittäin huonosti käyttäytyvä. Kuitenkin tällaista inversiota tehdään, esimerkiksi lähteessä *Electromagnetic Scattering by Non-spherical Particles* (1999) useat artikkelit käsittelevät cirrus-pilvien partikkelien (jääkiteitä) muodon ja koon johtamista niiden valonsironnasta. Samassa lähteessä kuitenkin osoitetaan, miten oletettu väärä muoto, koko tai refraktiaindeksi muuttavat tuloksia oleellisesti (Liu ym. 1999).

Yhtenä ratkaisuna tähän inversio-ongelmaan tarjotaan tässä työssä muodon mallintamista erikseen, ei valonsirontainversiona vaan suoraan partikkeleista otetuista valokuvista. Tämä lähestymistapa tuntuu intuitiivisesti paljon suoremalta ja perustellummalta kuin valonsironnan kautta lähestyminen. Ongelmia aiheuttaa vain se, että valokuvissa näkyy kolmiulotteisten muotojen projektioita, jotka eivät aina ole edes samoin jakautuneita eri projektiosuuntiin. Lisäksi suora muodon mallintaminen ei onnistu partikkeleille, joista on mahdoton saada hyvää kuvamateriaalia, kuten kosmiselle pölylle.

## 1.2 Työn rakenne

Usein opinnäytetyössä on tapana esitellä ensin teoreettiset menetelmät, jonka jälkeen niitä lopuksi sovelletaan aineistoon. Tässä työssä on kuitenkin katsottu perustelluksi käyttää erilaista jaottelua useastakin syystä. Ensinnäkin, fysiikkaan ja tähtitieteeseen liittyvää havaintoaineistoa ei tule usein vastaan tilastotieteen laitoksella, joten monet aineistoon liittyvät termit voivat kaivata ensin pienen esittelyn. Toiseksi, sekä työssä käytetty aineisto että menetelmät jakaantuvat selvästi kolmeen osaan, joissa jokaisessa on omanlaisensa ongelma. Kolmanneksi, jotkut teoreettiset menetelmät ovat syvästi sidoksissa juuri tässä käytössä olevaan aineistoon, eikä niitä ole luontevaa



käydä läpi, ennen kuin aineisto on esitelty. Siksi varsinainen työ on jaettu kolmeen päälukuun, jossa kussakin esitellään ensin aineisto ja vasta sitten teoreettiset menetelmät.

Luvussa 2 esitellään mallinnettavat partikkelit sekä niiden teoreettinen malli. Luvussa 3 esitellään partikkeleista lasketut suureet sekä mallin estimointi. Luvussa 4 esitellään partikkeleiden valonsirontamittaukset sekä mallin verifiointi näiden pohjalta. Luvussa 5 keskustellaan tuloksista.

## **2 PROGRA<sup>2</sup>-ohjelman partikkelien mallinnus**

### **2.1 PROGRA<sup>2</sup>-ohjelma**

Johdannossa mainittiin partikkelien orientaatioon ja projektiosuuntiin liittyvä keskeinen ongelma kuvien perusteella mallintamisessa. Jos partikkelit eivät ole pallokoordinaatistossa tarkasteltuna tasaisesti orientoituneita (ks. liite A), niin eri projektiosuuntiin otetuissa kuvissa olevat partikkelien satunnaisprojektiot eivät ole samoin jakautuneita, vaan jakauma on riippuvainen projektion suunnasta. Partikkelien orientaatioon vaikuttavat sekä maan painovoima, vetäen niitä näyteastian pohjalle makaamaan, että pienten partikkelien tapauksessa mahdollisesti näyteastian ja partikkelien väliset muut vetovoimat. Partikkelien projektion muodostumista on mahdollista mallintaa esimerkiksi simuloimalla niiden orientoitumista näyteastialle, mutta tällainen lähestymistapa tuo tietenkin ei-toivottua lisätyötä ja epävarmuutta. Helpompi tilanne mallituksen ja havaintojen tulkinnan kannalta on sellainen, jossa partikkelien voi olettaa olevan (tasaisesti) satunnaisorientoituneita. Käytännössä tämä vaatii ainakin painovoiman vaikutuksen poistamista mittaustilanteesta. Rakenteilla olevalla ISS-avaruusasemalla (International Space Station) on mahdollista tutkia partikkeleita mikrogravitaatiossa (painovoiman vaikutus olematon). Tähän projektiin osallistuu myös Helsingin yliopiston tähtitieteen laitoksen planeettaryhmä.

Eräs vaihtoehtoinen toteutustapa mikrogravitaatiolle on lentokoneella teh-

tävä ns. parabolinen lento. Parabolisessa lennossa lentokone ensin nousee noin 45 asteen kulmassa ylöspäin, jonka jälkeen moottoreiden teho laskeaan minimiin. Tämän jälkeen kone ohjataan tykinammuksen rataa muistuttavalle paraboliselle radalle, jonka aikana kone saavuttaa ensin radan lakipisteen ja kääntyy tämän jälkeen laskusuuntaan päätyen taas lähes 45 asteen laskevaan kulmaan, jolloin moottorien teho palautetaan taas normaaliiksi. Parabolisen radan aikana koneen sisällä vallitsee käytännössä painoton olotila, joka kestää noin puoli minuuttia. Painottoman vaiheen aikana voidaan tehdä mikrogravitaatiokokeita. Yhden lennon aikana kone voi toistaa painottomuuden muutamia kymmeniä kertoja. Novespace-yhtiöllä<sup>1</sup>, jonka pääomistaja on Ranskan avaruusjärjestö, on käytössään lentoihin sopiva Airbus A300 Zero-G -lentokone.

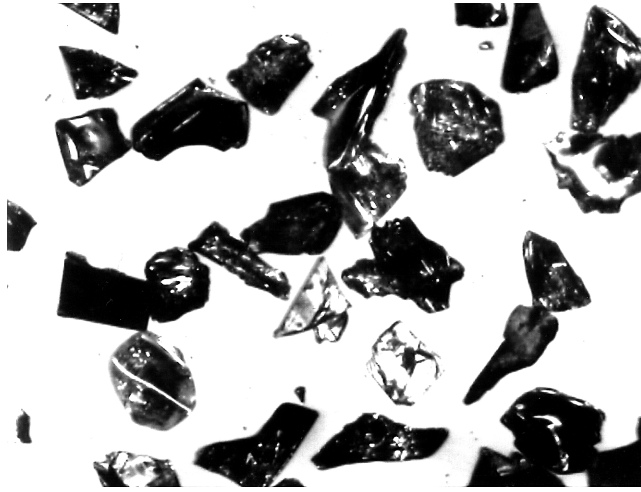
Partikkelien valonsirontamittauksista mikrogravitaatiossa vastaa ranskalainen tutkimusryhmä<sup>2</sup> PROGRA<sup>2</sup> (Propriétés Optiques des Grains Astronomiques et Atmosphériques). Ryhmällä on lennoilla mukana polarimetri, jolla kuvataan lasiastiassa olevaa partikkelinäytettä. Näytettä valaistaan laserilla (joko punaista valoa, aallonpituus 632,8 nm, tai vihreää valoa, aallonpituus 543,5 nm) ja siroava valo rekisteröidään CCD-kameralla. Näytettä kuvatesa näyteastia voidaan ensin hieman ravistaa mikrogravitaatiossa, jolloin partikkelit irtoavat astian pohjalta eivätkä niihin sen jälkeen mainittavasti vaikuta astian eikä maan vetovoimat. Tarkempaa tietoa mittauslaitteistosta löytää esimerkiksi lähteistä Worms ym. (1999, 2000).

PROGRA<sup>2</sup>-ohjelmassa on tutkittu useita eri partikkelityyppejä aggregoituneista partikkeleista yksittäisiin partikkeleihin. Vaikka ryhmämme onkin tehnyt malleja myös aggregaateille, niin tässä työssä keskitytään yksittäisiin partikkeleihin. Näistä eniten kuvamateriaalia mallintamista varten oli saatavilla boorikarbidista ( $B_4C$ , eng. boron carbide). Näyte  $B_4C$ -partikkeleista on esitelty kuvassa 1.  $B_4C$ -partikkeleita on ohjelmassa tutkittu kolmesta eri kokoluokasta, läpimitaltaan 9, 13 ja 88  $\mu m$ .

---

<sup>1</sup> Yhtiön kotisivu <http://www.novespace.fr>.

<sup>2</sup> Ryhmän kotisivu <http://www.esf.org/jcw/progra2.htm>.



Kuva 1: Esimerkki mallinnettavista  $B_4C$ -partikkeleista.

## 2.2 Teoreettinen malli PROGRA<sup>2</sup>-ohjelman partikkeleille

Mallinnettaessa PROGRA<sup>2</sup>-ohjelman  $B_4C$ -partikkeleita kannattaa sopivaa geometrista mallia etsiä lähtien partikkeleiden kuvista, esimerkkinä kuvassa 1 näkyvät partikkelit. Partikkeleista voi erottaa kulmapisteitä, joita yhdistävät melko tasaiset pinnat. Siksi sopiva approksimaatio partikkeiden muodolle voisi olla monitahokas, jonka kärkipisteinä ovat partikkelin kulmat. Kuvissa näkyvät partikkelit eivät ole täysin konvekseja, mutta poikkeamat kappaleen konvekseen verhoon ovat pieniä. Lisäksi konvekseilla kappaleilla on monia etuja yleisiin kappaleisiin verrattuna. Konvekseilla kappaleilla on joitain hyödyllisiä matemaattisia ominaisuuksia ja lisäksi konveksien kappaleiden käsittelyyn on olemassa monia algoritmeja ja valmiita tietokoneohjelmia. Siksi  $B_4C$ -partikkeliden mallintamiseen valittiin kolmiulotteisen avaruuden konvekseen monitahokas.

Annetun pistejoukon konveksin verhon muodostaminen sekä konveksin kappaleen reunapisteiden kolmiointi ovat ongelmia, jotka tulevat usein vastaan sellaisilla tieteenaloilla kuten spatiaalinen tilastotiede, tilastollinen geometria, hahmontunnistus ja tietokonegrafiikka. Tässä työssä on havaittu erityisen hyödylliseksi hollantilainen väitöskirjatyö (van de Weygaert 1991),

jossa käsitellään konveksien kappaleiden ominaisuuksia ja muodostamista, sekä QHull<sup>3</sup>-tietokoneohjelma, jolla voidaan muodostaa pistejoukon konvekssi verho ja sitä vastaava Delaunayn kolmiointi (Barber ym. 1996).

B<sub>4</sub>C-partikkelit, tai ainakin niiden projektiot kaksiulotteiseen avaruuteen (valokuvat), näyttävät olevan melko ympyrämäisiä siinä mielessä, että ne eivät ole erityisen venyneitä mihinkään tiettyyn suuntaan. Jos tätä ajatusta soveltaa teoreettisen mallin puolelle, niin voisi sanoa kappaleen ulkoreunaa eli sädettä kuvaavan satunnaisprosessin  $r(\theta, \phi)$  olevan kierto invariantti suhteessa pallokoordinaattisuuntiin  $(\theta, \phi)$ . Seuraavaksi esitellään yksityiskohtaisesti sellainen kierto invariantti satunnaisprosessi kolmiulotteisessa avaruudessa, jolla saadaan aikaan konveksien monitahokkaiden realisaatioita.

### 2.2.1 Mallin konstruointi

Tarkastellaan satunnaista joukkoa, kooltaan  $n$ , joka koostuu reaaliavaruuden  $\mathbb{R}^3$  sijaintinsa suhteen satunnaisista pisteistä  $p_1, \dots, p_n$ . Pisteiden koordinaatit on annettu pallokoordinaateissa  $p_i = (r_i, \theta_i, \phi_i)$ . Pisteiden koordinaatit ovat realisaatioita reaaliarvoisista satunnaismuuttujista  $(R, \Theta, \Phi)$ . Pallokoordinaattien suunnat  $\Theta, \Phi$  muodostavat yhdessä tasaisen jakauman pallon pinnalle (ks. tarkemmin liitteestä A).

Pisteen origosta mitattu etäisyys  $R$ , jota jatkossa kutsutaan pisteen säteeksi, noudattaa mallissa log-normaalijakaumaa. Mallissa on käytetty tavallisesta poikkeavaa parametrisointia log-normaalijakaumasta. Yleensä parametreina käytetään taustalla olevan normaalijakauman parametreja  $\mu$  ja  $\sigma^2$ , jolloin log-normaalijakauman odotusarvoksi saadaan  $\exp(\mu + \frac{\sigma^2}{2})$  ja varianssiksi  $\exp(2\mu)(\exp(2\sigma^2) - \exp(\sigma^2))$ . Koska jakauman todellinen odotusarvo ja hajonta ovat tärkeässä asemassa nyt tarkasteltavassa mallissa, on kätevää käyttää niitä myös jakauman parametreina. Merkitään näitä uusia parametreja symboleilla  $\mu_L$  (odotusarvo) ja  $\sigma_L$  (hajonta). Koska B<sub>4</sub>C-partikkelien absoluuttinen koko ei ole kiinnostuksen kohteena ja koska erikokoisten par-

<sup>3</sup> Ohjelmiston kotisivu <http://www.geom.umn.edu/software/qhull>.

tikkeliin oletetaan noudattavan muotonsa puolesta samaa mallia, asetetaan lisäksi pisteen säteen odotusarvo  $\mu_L$  arvoon yksi. Alkuperäiset log-normaalijakauman parametrit saadaan nyt siten, että

$$\mu = -\frac{1}{2} \log(1 + \sigma_L^2) \quad (1)$$

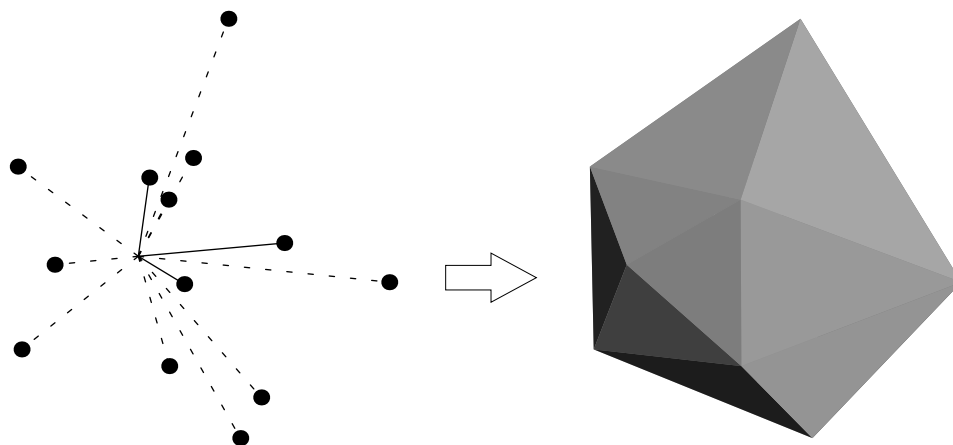
$$\sigma^2 = \log(1 + \sigma_L^2), \quad (2)$$

jossa vaihtoehtoisesta parametrisoinnista tarvitaan enää säteen hajontaa  $\sigma_L$ .

Kun satunnainen pistejoukko on arvottu, muodostetaan nämä pisteet sisältävä tilavuudeltaan pienin mahdollinen konvekssi verho. Jotkut alkuperäisistä pisteistä voivat jäädä konveksin verhon sisään eivätkä siten enää näy kulmapisteinä partikkelin ulkokuoressa. Konveksin verhon muodostava algoritmi yhdistää myös verhon pintaan jäävät pisteet toisiinsa tasomonitahokkaila, käytännössä kolmioilla. Tämä verhon pinnalle jäävien kulmapisteiden joukko ja sen kolmiointi muodostavat lopullisen realisaation mallista  $B_4C$ -partikkeleille. Mallin parametreina ovat siis alkuperäisten pisteiden määrä  $n$  sekä pisteiden säteen hajonta  $\sigma_L$ . Koska tarkoituksena on muodostaa aidosti kolmiulotteinen kappale, on syytä asettaa parametrille  $n$  määrittelyalue  $n = 4, \dots, \infty$ .

Näin muodostetun kappaleen sädettä kuvaava satunnaisprosessi  $r(\theta, \phi)$  on kierto invariantti, koska kappaleen muodostukseen käytetyn pistejoukon suuntakoordinaatit  $\Theta, \Phi$  ovat tasaisesti jakautuneet, ja pisteen säteen jakauma on riippumaton suuntakoordinaateista. Näin ollen on selvää, että myös konveksin verhonkin muodostamisen jälkeen  $r(\theta_1, \phi_1)$  ja  $r(\theta_2, \phi_2)$  ovat satunnaisuuttujina samoin jakautuneita millä tahansa valinnalla  $\theta_1, \phi_1, \theta_2, \phi_2$ . Esimerkki mallin realisaatiosta on esitetty kuvassa 2.

Realisaatioiden arpominen mallista toteutettiin Mathematica-ohjelmalla. Realisaationa ei tarvita varsinaista kolmiulotteista mallia, vaan sen kaksiulotteinen projektio, koska myös havaittu data koostuu näistä projektioista. Mallin toteutus on nähtävänä liitteessä B. Mathematican edut ohjelmankehityksessä ovat sen monipuolisuus ja valmiit komponentit, jolloin ohjelman



Kuva 2: Satunnainen konvekssi monitahokas parametrien arvoilla  $n = 13$  ja  $\sigma_L = 0,4$ . Vasemmanpuoleisessa kuvassa on origosta arvotut 13 kulmapistettä, joista kuitenkin kolme (yhtenäinen viiva) on jäänyt konveksin verhon sisäpuolelle, ja kymmenen (katkoviiva) päätynt lopulliseen kappaleeseen. Kappaleen pinta kuvassa oikealla.

saa valmiiksi erittäin nopeasti ja ohjelmakoodi pysyy lyhyenä sekä luettavana. Huonona puolena voisi pitää hitautta raskaissa tehtävissä. Tässä työssä projektioita arvottiin loppujen lopuksi 71 500 000 kappaletta, joiden arpoamiseen kului arviolta yli 13 CPU-päivää tehokkaalta PC-laitteistolta. Tarvittava aika olisi lyhentynyt murto-osaan, jos koodi olisi kirjoitettu jollain varsinaisella ohjelmointikielellä kuten FORTRAN-kielellä. Suurimpana ongelmana oli kuitenkin konveksin verhon muodostaminen, johon oli vaikeaa löytää sopivia algoritmeja valmiina, ja oman algoritmin kirjoittaminen olisi ollut melko raskas tehtävä. Onneksi aineistoa arvottiin lisää useassa erässä muun työn ohessa, joten 13 CPU-päivän kesto ei loppujen lopuksi muodostanut ongelmaa.

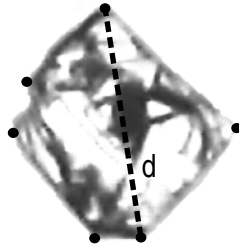
## 3 Mallin parametrien estimointi

### 3.1 B<sub>4</sub>C-partikkeleista saatu havaintoaineisto

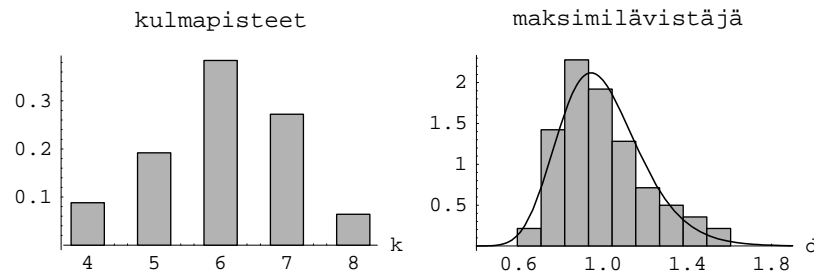
B<sub>4</sub>C-partikkelien mallintamista varten kuvista on mitattava yksittäisten partikkelien muotoa. Periaatteessa muodosta kertovan mittausaineiston muodostaa kameran tarkkuudella saatu pikselimatriisi tai -tensori koko partikkelin värisävyistä (kuvat alun perin värillisiä). Käytännössä muotoanalyysiä on kuitenkin helpompi tehdä, jos lähtökohtana on joitain järkevällä tavalla tiivistettyjä partikkelin muotoa kuvaavia tunnuslukuja pikselimatriisiin sijaan. Yleensä muotoa kuvatessa on mahdotonta löytää mitään tilastollisen päättelyn teorian mukaista tyhjentävää tunnuslukua, joka olisi samalla tulkinnallisesti järkevä ja tiivistäisi muototietoa. Käytännössä täytyy vain tyytyä kadottamaan osa datan sisältämästä muodosta kertovasta informaatiosta jotta saataisiin yksinkertainen ja helposti mitattava tunnusluku.

Partikkelit näyttäisivät kuuluvan Lumpeen ym. (1995) ehdottamaan satunnaisten monitahokkaiden luokkaan, joten partikkeleilla on selviä, teräviä kulmapisteitä. Ensimmäiseksi tunnusluvuksi on valittu partikkelin projektiossa näkyvien kulmien määrä ( $k$ ). Kulmien määrän lisäksi kiinnostava muotoa kuvaava suure on partikkelien koon vaihtelu. Muoto ja koko oletetaan usein toisistaan riippumattomiksi, kuten tässäkin, mutta koon vaihtelu antaa silti tietoa mallinnusta varten. Koon vaihtelua kuvaavaksi tunnusluvuksi on valittu partikkelin maksimilävistäjä ( $d$ ), kuitenkin niin, että maksimilävistäjän keskiarvo yli kaikkien mitattujen partikkelien skaalataan arvoon yksi.

Tutkimuksessa oli käytössä kuusi esimerkkikuvan 1 tapaista kuvaa, joista voitiin erottaa yhteensä 125 projektiota B<sub>4</sub>C-partikkelista. Jokaisesta projektioista mitattiin kulmien määrä ja maksimilävistäjä kuten kuvassa 3. Kulmien määrä hieman epätasaisesta kappaleesta on jokseenkin subjektiivinen luku, mutta kulmapisteitä määritettäessä yritettiin saada partikkelin ensisijainen muoto vangittua välittämättä niinkään pinnan pienistä epätasaisuuksista. Koska muodon mallinnuksessa käytettiin partikkelin mallina konvek-



Kuva 3: Esimerkki partikkelista mitatuista tunnusluvuista: kulmien määrä (tässä kuusi) ja maksimilävistäjä ( $d$ ).



Kuva 4:  $B_4C$ -partikkeleista mitattujen tunnuslukujen jakaumat aineistossa. Vasemmalla kulmapisteiden ( $k$ ) jakauma, oikealla maksimilävistäjän ( $d$ ) jakauma, johon on sovitettu log-normaalijakauma.

sia muotoa, myös kulmapisteistä otettiin mukaan vain kappaleen konvekseen verhoon kuuluvat pisteet. Kuvassa 4 näkyy näiden kahden tunnusluvun jakaumat käytössä olleessa aineistossa.

### 3.2 Parametrien estimointi havaintoaineistosta

Seuraavassa käydään ensin läpi parametrien estimoinnissa tarvittavat menetelmät, jonka jälkeen menetelmiä sovelletaan havaittuun aineistoon.



### 3.2.1 Suurimman uskottavuuden menetelmä

Suurimman uskottavuuden menetelmä on yksi tilastotieteen perusmenetelmistä, kun on tarkoitus estimoida todennäköisyysmallin  $\mathbf{X} \sim f(\mathbf{x}; \boldsymbol{\theta})$  tuntematonta parametrivektoria  $\boldsymbol{\theta}$  havainnon  $\mathbf{x}$  perusteella. Päätely suurimman uskottavuuden menetelmässä perustuu mallin uskottavuusfunktioon  $L(\boldsymbol{\theta}; \mathbf{x})$ , jonka tulee olla suoraan verrannollinen havainnon  $\mathbf{x}$  todennäköisyystiheuteen, kun parametrivektorin arvo on kiinnitetty. Uskottavuusfunktio on siis muotoa  $L(\boldsymbol{\theta}; \mathbf{x}) \propto f(\mathbf{x}; \boldsymbol{\theta})$ .

Suurimman uskottavuuden estimaatti  $\hat{\boldsymbol{\theta}}$  on sellainen piste parametriavaruudessa  $\Theta$ , jolle

$$L(\hat{\boldsymbol{\theta}}; \mathbf{x}) \geq L(\boldsymbol{\theta}; \mathbf{x}) \quad \forall \boldsymbol{\theta} \in \Theta . \quad (3)$$

Kun uskottavuusfunktio on suoraan verrannollinen havainnon todennäköisyystiheuteen, voi suurimman uskottavuuden estimaatin tulkita sellaiseksi parametrin arvoksi, joka antaa suurimman todennäköisyysidentiteyden havainnolle.

Yleensä maksimointi tehdään kuitenkin log-uskottavuusfunktioon perustuen, joka on nimensä mukaisesti uskottavuusfunktion logaritmi. Laskennallisesti maksimointi johtaa yhä samaan tulokseen, koska log on aidosti monotoninen funktio. Log-uskottavuusfunktioilla on joitain etuja, kun havaittu data  $\mathbf{x}$  koostuu useammasta toisistaan riippumattomasta samoin jakautuneesta havainnosta  $(x_1, \dots, x_n)$  tai havaintovektorista. Riippumattomuuden nojalla aineiston yhteistiheysfunktio on yksittäisten tiheysfunktioiden tulo, ja sitä kautta myös uskottavuusfunktio on tulomuotoa

$$L(\boldsymbol{\theta}; \mathbf{x}) \propto \prod_{i=1}^n f(x_i; \boldsymbol{\theta}) , \quad (4)$$

jolloin log-uskottavuusfunktio sen logaritmina on summa

$$l(\boldsymbol{\theta}; \mathbf{x}) = c(\mathbf{x}) + \sum_{i=1}^n \log (f(x_i; \boldsymbol{\theta})) , \quad (5)$$

jossa  $c(\mathbf{x})$  on joku pelkästään havainnoista riippuva funktio. Kun log-uskottavuusfunktio on  $n$ :n satunnaismuuttujan  $\log (f(X_i; \boldsymbol{\theta}))$  summa, voidaan

siihen ja suurimman uskottavuuden estimaattiin soveltaa keskeisen raja-arvolauseen ja suurten lukujen lain tuloksia. Keskeisen raja-arvolauseen nojalla suurimman uskottavuuden estimaattori on asympotoottisesti normaali-jakautunut.

Suurimman uskottavuuden estimaattorin muodostaminen tehdään normaaleilla funktion maksimoinnin menetelmillä. Analyttisesti ratkaistuna estimaattori on juuri uskottavuusyhtälöille, eli yhtälöryhmälle, jossa log-uskottavuusfunktion osittaisderivaatat  $\theta$ :n suhteen asetetaan nolliksi.

### 3.2.2 Simuloitu suurimman uskottavuuden menetelmä

Uskottavuusyhtälöiden juurta ei välttämättä saada analyttisesti ratkaistua, vaan joudutaan tyytymään numeeriseen ratkaisuun. Tämä ei ole mitenkään poikkeuksellinen tilanne. Hankalampi tilanne saadaan, jos edes havaintojen todennäköisyysjakauma ei ole esitettävissä suljetussa muodossa. Tällainen tilanne voi syntyä, jos todennäköisyysmalli ei ole suoraan mikään tunnettu parametrinen jakauma, vaan havainnot synnytetään jonkin satunnaisuutta sisältävän säännön tai proseduurin avulla (käytetään jatkossa nimitystä malliproseduuri). Tästä huolimatta voidaan sanoa, että havainnoilla on todennäköisyysjakauma  $f(\mathbf{x}; \theta)$ , vaikka sitä ei voidakaan esittää suljettuna parametrinen jakaumana.

Simuloitu suurimman uskottavuuden menetelmä sopii edellä mainitun kaltaisiin tapauksiin. Vaikka todennäköisyysjakaumaa havaintoja edustavalle satunnaismuuttujalle  $\mathbf{X}$  ei tiedetä, havaitaan jakauman ominaisuuksia epäsuorasti sen tuottamien realisaatioiden  $\mathbf{x}_i$  kautta. Koska malliproseduuri on tiedossa, voidaan satunnaismuuttujan realisaatioita tuottaa simuloimalla. Kiinteillä parametrivektorin arvoilla satunnaismuuttujan simuloitujen realisaatiot noudattavat (tuntematonta) jakaumaa  $f(\mathbf{X}; \theta)$ , mutta jakaumaa voidaan approksimoida jollain sopivalla tiheysfunktioestimaattorilla. Yleisessä tapauksessa tiheysfunktioestimaattoriksi kannattaa valita epäparametrinen estimaattori, koska parametrisestä jakaumasta ei ole tietoa. Luvussa 3.2.3

tarkastellaan ydinestimaatin käyttöä tiheysfunktioestimaattorina.

Satunnaismuuttujan tiheysfunktion estimaatti antaa samalla myös estimaatin parametrien uskottavuusfunktiolle, joka on suoraan verrannollinen tiheysfunktioon vapaasti valittavalla vakiokertoimella. Tapauksessa, jossa tiheysfunktioestimaatti  $\hat{f}$  on oikean tiheysfunktion tarkentuva estimaattori, on myös näin estimoidun uskottavuusfunktion maksimoinnin tulos estimaattori oikealle suurimman uskottavuuden estimaatille.

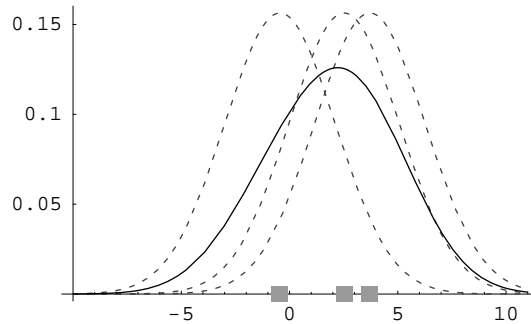
Edellä esitetyn kaltaisia lähestymistapoja estimointiin löytyy kirjallisuudesta jonkin verran. Penttinen (1984) on väitöskirjassaan tarkastellut simuloidun uskottavuuden käyttöä spatiaalisten pisteprosessien estimoinnissa. Ekonometrian puolelta löytyy simuloidun suurimman uskottavuuden menetelmän käyttöä osakekursseihin sovellettuna artikkelista Brandt ym. (2002). Artikkelissa mainitaan vastaavaa lähestymistapaa käytetyn ensi kertaa siinä ongelmakehyksessä vuonna 1995 artikkelin toisen tekijän väitöskirjassa (Santa-Clara 1995). Kirjassa Monte Carlo Statistical Methods (Robert ym. 1999) käsitellään laajemmin simulointimenetelmiä, sivuten myös suurimman uskottavuuden estimointia.

### 3.2.3 Tiheysfunktion ydinestimointi

Suosituimpia menetelmiä tiheysfunktion epäparametriseen estimointiin on nykyisin ns. ydinestimointi, jonka on ensi kerran esitellyt Parzen 1962. Ydinestimoinnissa korvataan jokainen havaintopiste  $x_i$  ydinfunktiolla  $K(x; x_i, h)$ , jonka tulee olla symmetrinen sekä ei-negatiivinen ja jonka tilavuus on yksi. Tiheysfunktion estimaatin  $\hat{f}(x)$  muodostaa ytimien summa

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K(x; x_i, h) . \quad (6)$$

Ydinfunktion valintaa voidaan tarkastella eri virhekretereillä, mutta käytännössä paljon käytetty on Gaussin ydin, eli normaalijakauman tiheysfunktio odotusarvolla  $x_i$  ja varianssilla  $h^2$ . Moniulotteisessa tapauksessa on kovarianssimatriisi muotoa  $h^2 \Sigma$ , jossa  $\Sigma$  on joko identiteettimatriisi, tai otoksen



Kuva 5: Kolme havaintoa normaalijakaumasta (harmaat neliöt) ja niihin liitetyt ytimet (katkoviiva), sekä tiheysfunktion ydinstimaatti (yhtenäinen viiva).

perusteella estimoitu korrelaatiomatriisi. Käytännössä paljon ytimen muotoa tärkeämpi kysymys on siloitusparametrin  $h$  valinta.

Siloitusparametri voidaan valita monella eri tavalla. Yksinkertaisimpia sääntöjä on Silvermanin (1986) ehdottama

$$h = \sigma \left( \frac{4}{p+2} \right)^{\frac{1}{p+4}} n^{-\frac{1}{p+4}}, \quad (7)$$

jossa  $p$  on tiheysfunktion dimensio, ja  $\sigma^2$  voidaan estimoida aineiston kovarianssimatriisin diagonaalelementeistä  $s_i$  siten, että

$$\sigma^2 = \frac{1}{p} \sum_{i=1}^p s_i. \quad (8)$$

Silvermanin ehdotus siloitusparametrille on oikean tiheysfunktion ja ydinstimaatin välisen  $L^2$ -normin mielessä optimaalinen, kun aineisto on peräisin normaalijakaumasta ja ytimenä käytetään Gaussin ydintä (esimerkki kuvassa 5). Muihin tilanteisiin sopivampia mutta monimutkaisempia menetelmiä on esitelty runsaasti esimerkiksi lähteessä Holmström (2002).

Rosenblatt (1956) on osoittanut, että kaikki ei-negatiiviset tiheysfunktioestimaatit ovat harhaisia äärellisillä otoksilla, kun otoksen tiheysfunktioperhettä ei olla rajoitettu (Webb 1999, luku 3.5; Holmström 2002, luku 3.1).

Kuitenkin, jos seuraavat ehdot ovat voimassa ytimelle ja sen siloitusparametrille (otoskoon funktiona,  $h := h(n)$ ), niin ydineestimaatti on tiheysfunktion asymptoottisesti harhaton sekä asymptoottisesti tarkentuva estimaatti:

$$\int_{-\infty}^{\infty} K(x)dx = 1 \quad (9a)$$

$$\sup_x K(x) < \infty \quad (9b)$$

$$\lim_{x \rightarrow \infty} xK(x) = 0 \quad (9c)$$

$$\lim_{n \rightarrow \infty} h(n) = 0 \quad (9d)$$

$$\lim_{n \rightarrow \infty} nh(n) = \infty . \quad (9e)$$

Nämä pätevät Gaussin ytimelle ja kaavan (7) mukaiselle silotusparametrille. Tuloksen todistus on liitteessä C.

### 3.2.4 Malliparametrien estimointi $B_4C$ -partikkeleille

#### Johdattelua

Luvussa 3.2 edellä esiteltyjä menetelmiä on nyt tarkoitus soveltaa  $B_4C$ -partikkeleille ja niitä kuvaavalle monitahokasmallille. Mallin realisaatioiden muotoon vaikuttavat vapaat parametrit olivat  $n$  ja  $\sigma_L$  (luku 2.2.1). Havaittu aineisto taas puolestaan koostuu projektioissa näkyvien kulmien määrästä  $k$ , sekä projektiokappaleen maksimilävistäjästä  $d$  (luku 3.1).

Kiinnitetyillä parametrien arvoilla mallin tuottamat realisaatiot havaittavista muuttujista  $k$  ja  $d$  noudattavat todennäköisyysjakaumaa, jonka analyttinen määrittäminen on kuitenkin vaikea, jollei mahdoton tehtävä. Ensinnäkin, säteistä ja suunnista muodostuvat satunnaispisteet eivät muodosta kappaletta sellaisenaan, vaan vain konveksin verhon ulkoreunalle päätyvät pisteet vaikuttavat kappaleen ulkomuotoon. Verhoon päätyvät pisteet noudattavat nyt alkuperäisen jakauman ehdollista muotoa ehtona sijoittuminen konveksin verhon reunalle. Toiseksi, näistä pisteistä muodostetaan vielä satunnaisprojektiio alempaan ulottuvuuteen, jolloin jakauma muuntuu lisää. Kolmanneksi, tästä jakaumasta lasketaan vielä muunnoksena tunnusluvut  $k$  ja  $d$ , joista eteenkin  $d$  maksimina olisi vaikea johtaa, vaikka edelliset vaiheet

selvitettäisiinkin. Johtopäätöksenä on, ettei tunnuslukujen jakauman selvittäminen analyttisesti ole realistinen tavoite. Koska malliproseduuri on kuitenkin tiedossa, voidaan tästä jakaumasta saada teoriassa rajattomasti näytteitä simuloimalla, joten luvussa 3.2.2 esitetty simuloitu suurimman uskottavuuden menetelmä sopii mallin parametrien estimointiin.

### Ydinestimointi

Ei ole syytä olettaa että havaitut muuttujat  $k$  ja  $d$  olisivat toisistaan riippumattomia, joten muuttujien jakauman muodostaminen kiinteillä parametriarvoilla edellyttää täyden kaksidimensionaalisen jakauman estimointia. Koska muuttuja  $k$  on kuitenkin diskreetti, voidaan jakauman estimointia yksinkertaistaa jakamalla estimointi  $k$ :n arvojen todennäköisyyksien estimointiin, ja  $d$ :n jatkuvan (yksiulotteisen) jakauman estimointiin kiinteällä  $k$ :n arvolla. Tällöin jakauman dekompositio on

$$f(k, d; n, \sigma_L) = P(K = k; n, \sigma_L) f(d|K = k; n, \sigma_L) . \quad (10)$$

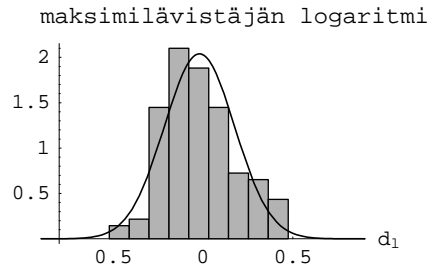
Koska mahdollisia  $k$ :n arvoja on vähän ( $k = 3, \dots, n$ ), voidaan todennäköisyydet  $P(K = k; n, \sigma_L)$  estimoida simuloidusta aineistosta suoraan  $k$ :n frekvensseistä. Jos  $P(K = k; n, \sigma_L) = 0$  jollain  $k$ , niin tätä vastaava yhteisjakauma  $f(k, d; n, \sigma_L) = 0$  kaikilla  $d$ . Siksi kutakin parametriparia  $n, \sigma_L$  vastaavaan jakaumaan tarvitsee  $k$ :n todennäköisyyksien lisäksi estimoida ehdolliset yksiulotteiset jakaumat  $f(d|K = k; n, \sigma_L)$  vain niillä  $k$ :n arvoilla, joilla  $P(K = k; n, \sigma_L) > 0$ . Kiinnitettyä parametriparia kohden tarvitsee siis estimoida yksiulotteisia ehdollisia  $d$ :n jakaumia noin 5-20. Lukumäärä riippuen enimmäkseen parametrin  $n$  arvosta. Nämä jakaumaestimaatit muodostetaan ydinestimoinnilla.

Itse asiassa  $k$ :n frekvenssejäkään ei tarvitse laskea erikseen, koska yhteistheysfunktion ydinestimaatti  $\hat{f}$ , kun havaintoja, joilla  $K = k$ , on  $n_k$  kappaletta, on

$$\begin{aligned} \hat{f}(k, d; n, \sigma_L) &= \hat{P}(K = k; n, \sigma_L) \hat{f}(d|K = k; n, \sigma_L) \\ &= \frac{n_k}{n} \frac{1}{n_k} \sum_i^{n_k} K(d; d_i, h) = \frac{1}{n} \sum_i^{n_k} K(d; d_i, h) . \end{aligned} \quad (11)$$

Tarkennusta vaille jää nyt ainoastaan silotusparametrin  $h := h(n_k)$  kaavassa (7) oleva hajonta  $\sigma$ . Sen arvioimiseen on kaksi vaihtoehtoa: joko siitä havaintojoukosta  $d_i$  joilla  $K = k$ , tai yli kaikkien kaksikulotteiseen jakaumaan kuuluvien havaintojen  $d_i$ . Tässä päädyttiin valitsemaan jälkimmäinen vaihtoehto, koska joillain arvoilla  $k$  voi havaintoja  $d_i$  olla melko vähän, jolloin hajontaestimaattori ei välttämättä olisi luotettava siinä joukossa. Lisäksi on syytä olettaa, että hajonta on suurinpiirtein vakio kaikilla  $k$ . Tällä tavalla menetellen voi olla pieni mahdollisuus saada yliestimoitu hajonta ja sitä kautta liian voimakas silotus, mutta kun havaintoja yhteensä on hyvin runsaasti, niin silotusparametrinkaan vaikutus ei todennäköisesti ole suuri.

Muuttuja  $d$ , jonka ehdollisia jakaumia on tarkoitus estimoida, on projektio-partikkelin kulmien maksimilävistäjä, ja näin ollen ei-negatiivinen muuttuja. Muuttujan ehdolliset jakaumat muistuttavat muodoltaan log-normaalijakaumaa; jakaumat ovat kutakuinkin yksihuippuisia, ja nollaa lähestyttäessä todennäköisyystiheys lähestyy myös nollaa (kuva 4). Ydinestimointi ei perusmuodossaan anna hyviä estimaatteja rajoitettujen (määrittelyjoukko ei ole koko  $\mathbb{R}^n$ ) muuttujien jakaumille, vaan määrittelyjoukon reunoilla täytyy käyttää joitain erityismenettelyjä. Joissain tapauksissa erityismenettelyt voidaan kiertää muuntamalla alkuperäistä muuttujaa sopivasti, jolloin saadaan koko  $\mathbb{R}^n$ :ssä määritelty muuttuja. Tässä tapauksessa logaritmuunnos sopii hyvin, koska se muuntaa muuttujan  $d$  jakauman rajoittamattomaksi ja myös likimain normaaliseksi (kuva 6). Termiä 'likimain normaalin' käytetään tässä melko vapaasti, koska jakaumaa ei tarvitse olettaa normaaliksi, vaan jakauman likimaista muotoa käytetään hyväksi ainoastaan ei-parametrisen tiheysfunktioestimaatin tyyppin valintaan. Uuden muuttujan  $d_l = \log(d)$  tiheysfunktioille saadaan hyvä ydinestimaatti käyttämällä luvussa 3.2.3 esiteltyä Gaussin ydintä, ja valitsemalla silotusparametriksi Silvermanin ehdottama normaalijakaumalle optimaalinen  $h$  kaavalla (7). Kun aineistoa voi ja on syytä simuloida riittävästi, ei silotusparametrin valinta muodostu enää kriittiseksi hyvän estimaatin muodostamisen kannalta. Kuvassa 7 on esimerkki tiheysfunktioestimaatista havaituille muuttujille  $k$ ,  $d$  eräällä parametriparin



Kuva 6: Logaritmimuunnetun muuttujan  $d_l$  jakauma havaitussa aineistossa sekä jakaumaan sovitettu normaalijakauma.

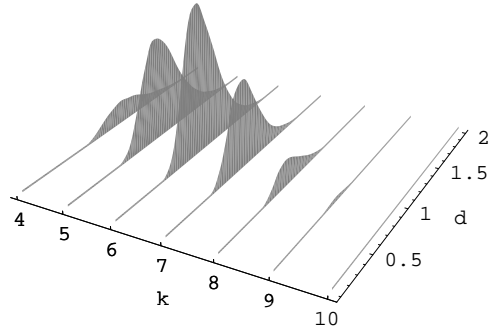
arvolla.

### Tietokonetoteutus

Simuloitu suurimman uskottavuuden estimointi toteutettiin itse kirjoitetulla FORTRAN 90 -kielisellä ohjelmalla, jonka lähdekoodi on nähtävissä liitteessä D. Estimointia varten simuloitiin lisää aineistoa asteittain liittäen uutta aineistoa vanhaan alkaen 5 000 havainnosta yhtä parametriparia  $n, \sigma_L$  kohti, päätyen 100 000 havaintoon parametriparia kohti. Aidossa havaitussa datassa oli 125 kappaletta havaintoja.

Yksi päämäärähän on optimoida uskottavuusfunktioita, eli löytää pinnan maksimikohta. Tässä työssä ei ole kuitenkaan käytetty mitään varsinaista optimointialgoritmia suurimman uskottavuuden estimaatin löytämiseksi. Tähän oli useita syitä. Ensinnäkin, simuloinnin tuloksena saatu piste suurimman uskottavuuden pinnalla parametrien funktiona on satunnaismuuttuja, jonka arvo muuttuu sitä mukaa kun uutta aineistoa arvotaan. Vaikka estimaatti onkin tarkentuva, voi heilahtelua olla kuitenkin paljon vielä suurillakin simulointimäärillä. Toisaalta kaikki nopeat optimointialgoritmit tarvitsevat vähintäänkin arvion kaksiulotteisen pinnan ensimmäisistä tai jopa toisista osittaisderivaatoista tietyissä pisteissä. Derivaatat pitäisi approksimoida numeerisesti, mikä on tunnetusti epävarma tehtävä jopa ei-satunnaiselle funktiolle, saati sitten sellaiselle, jossa voi olla reilustikin satunnaisvirhettä. Jos taas käytetään jotain optimointialgoritmia jossa ei vaadita

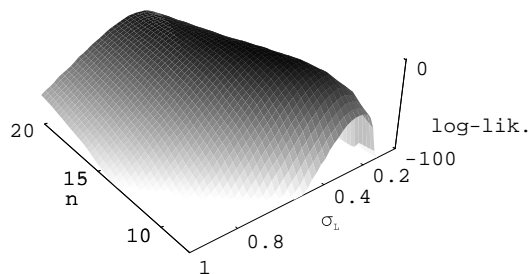




Kuva 7: Kaksiulotteisen tiheysfunktion ydineestimaatti havaituille muuttujille  $k$ ,  $d$ , kun mallin parametrit ovat  $n = 13$ ,  $\sigma_L = 0,4$ . Varsinainen ydineestimointi on tehty muunnetulle muuttujalle  $d_l$ , mutta kuvassa on palattu alkuperäisiin muuttuja-arvoihin.

derivaattojen muodostamista, kuten simuloitua jäähtytystä tai geneettisiä algoritmeja, tulee funktion arvo laskettua helposti hyvin monessa pisteessä, mikä on raskasta. Lisäksi on vaikea etukäteen tietää millä simuloitujen havaintojen määrällä uskottavuus ei enää heilahtelee liikaa.

Toiseksikin, suurimman uskottavuuden estimaattihan ei piste-estimaattina ole ainoa kiinnostuksen kohde, vaan järkevä estimointi antaa myös jonkinlaisen arvion estimoinnin luotettavuudelle. Luottamusväliestimointiin tarvitaan aina myös uskottavuuspinnan arvoja suurimman uskottavuuden pisteen ulkopuolelta. Siksi käytännöllinen ja vieläpä toteutukseltaan yksinkertainen lähestymistapa onkin laskea uskottavuuspinnan arvot jossain järkevässä parametrien arvojen kaksiulotteisessa hilassa, jolloin välipisteet interpoloimalla saadaan kuva koko uskottavuuspinnasta, kuten kuvassa 8. Hilassa käytettiin parametrille  $n$  arvoja  $8, 9, \dots, 20$  ja parametrille  $\sigma_L$  kahta eri tarkkuutta: yleisesti arvot olivat  $0,025$  välein lähtien arvosta  $0,025$  ja päätyen arvoon  $1$ , mutta kiinnostavampi väli  $(0,25; 0,5)$  käytiin läpi  $0,01:n$  suuruisin askelein.



Kuva 8: Parametrien  $n$  ja  $\sigma_L$  simuloitu uskottavuuspinta. Kutakin pinnan laskettua pistettä varten on simuloitu 100 000 havaintoa monitahokasmallista ja laskettujen pisteiden väliset pisteet on arvioitu lineaarisella interpoloinnilla. Kuvasta näkyy, että parametrit eivät ole toisistaan riippumattomia ja että parametrin  $\sigma_L$  vaikutus näkyy log-uskottavuudessa paljon herkemmin.

## Luottamusväleistä

Tapoja muodostaa uskottavuusfunktioon perustuvia luottamusvälejä on kolmea tyyppiä, Waldin testisuureeseen, uskottavuusosamäärään sekä piste-määrään perustuva. Näistä vaihtoehdoista sekä Waldin testisuure että pistemäärästä tarvitsevat periaatteessa Fisherin informaatiomatriisin  $i$  arvoa pisteessä  $\hat{\theta}$ , joka käytännössä korvataan havaitulla informaatiolla  $j(\hat{\theta})$ . Kuten edellä on jo todettu, derivaatat arvioituvat tällaisessa tapauksessa epäluotettavasti, erityisesti havaitun informaation vaatima Hessen matriisi. Ainoa käyttökelpoinen luottamusväliestimaattori on siten uskottavuusosamäärään eli Neyman-Pearsonin testisuureeseen perustuva. Uskottavuusosamäärään testisuure  $t$  vektoriparametrille on muotoa

$$t = 2(l(\hat{\theta}; \mathbf{x}) - l(\theta_0; \mathbf{x})) \overset{\sim}{\sim} \chi_d^2, \quad (12)$$

ja testisuureen sekä luottamusvälin välisen suhteen vuoksi  $\alpha$  suuruiseen luottamusalueeseen kuuluu tällöin ne parametriavaruuden  $\Theta$  pisteet, joilla

$$\{\theta \mid 2(l(\hat{\theta}; \mathbf{x}) - l(\theta; \mathbf{x})) \leq \chi_d^2(\alpha)\}. \quad (13)$$

Usein tarkastellaan log-uskottavuusfunktion sijaan sen skaalattua versiota  $l(\theta; \mathbf{x}) - l(\hat{\theta}; \mathbf{x})$ , jonka maksimiarvo 0 saavutetaan suurimman uskottavuuden estimaatilla. Jos asymptoottinen  $\chi_d^2$ -jakaumatulos pätee, saadaan 95%

luottamusväliksi skaalatulle log-uskottavuusfunktiolle yhtälön (13) nojalla se alue, jossa  $l(n, \sigma_L; \mathbf{x}) \leq -2,996$ , kun parametrivektorin dimensio  $d$  on kaksi ( $2,996 = \chi_2^2(0,95) / 2$ ). Periaatteessa pitäisi tarkastaa toteutuuko asymp-toottisuus riittävällä tarkkuudella, jotta luottamusvälin voisi tulkita 95% väliksi (Ekholm 1997). Log-uskottavuusfunktiota pitäisi tarkastella tuolla 95% alueella verraten sitä vastaavaan normaali- eli kvadraattiseen approk-simaatioon  $\frac{1}{2} l''(\hat{\theta}; \mathbf{x}) (\theta - \hat{\theta})^2$  (yksiulotteiselle parametrille). Tässä tarvitaan kuitenkin taas sellaista mitä ei voi saada, eli log-uskottavuusfunktion toista derivaattaa. Approksimaatio antaa kuitenkin ainakin suuntaa-antavan luot-tamusvälin.

Tulkinnaltaan uskottavuusosamäärän luottamusväli vastaa Bayes-päätelyn puolelta luottoväliä, kun parametrivektorin prioritiheysfunktio on epäinfor-matiivinen eli tasainen (esimerkiksi Tanner 1993). Itse asiassa koko (simuloi-tu) suurimman uskottavuuden estimointi vastaa tällaista Bayes-päätelyä. Luottovälin muodostamisessa suosittu metodi on ns. suurimman posteriori-tiheyden alue. Suurimman posterioritiheyden alue on sellainen parametria-varuuden väli  $[\theta_1, \theta_2]$ , jossa parametrin posterioritiheysfunktio täyttää ehdot

$$\int_{\theta_1}^{\theta_2} f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) d\theta = \alpha \quad (14a)$$

$$f_{\Theta|\mathbf{X}}(\theta|\mathbf{x}) \geq f_{\Theta|\mathbf{X}}(\theta^*|\mathbf{x}) \quad , \text{ kun } \theta \in [\theta_1, \theta_2] \text{ ja } \theta^* \notin [\theta_1, \theta_2] . \quad (14b)$$

Koska epäinformatiivisella prioritiheysfunktiolla parametrivektorin postero-ritiheys on suoraan verrannollinen havaintojen todennäköisyyteen ehdolla parametrivektori, on se myös suoraan verrannollinen parametrivektorin us-kottavuusfunktioon. Todennäköisyys  $\alpha$  on pinta-alan osuus kokonaispinta-alasta, joka tiheysfunktiolle on yksi mutta uskottavuusfunktiolle jotain muu-ta. Siksi suurimman posterioritiheyden alue  $(n, \sigma_L) \in \Theta$  saadaan tässä esti-moinnissa muotoon

$$(14a) \Rightarrow \frac{\iint_{\Theta} L(n, \sigma_L; \mathbf{x}) d\sigma_L dn}{\int_4^{\infty} \int_0^{\infty} L(n, \sigma_L; \mathbf{x}) d\sigma_L dn} = \alpha , \quad (15)$$

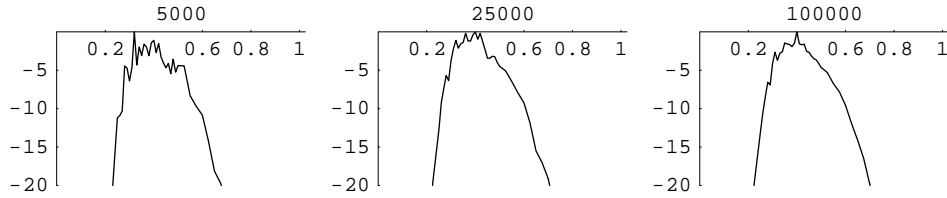
ja ehto (14b) uuteen muotoon suoraan merkintöjä muuttamalla. Huomaa, että muuttujaa  $n$  tarkastellaan tässä poikkeuksellisesti jatkuvana käyttämäl-lä lineaarisella interpolaatiolla saatua jatkuvaa uskottavuuspintaa (kuten

kuvassa 8). Näin saadut luottamusrajat muuttujalle ovat ei-kokonaislukuina hieman oudot, mutta toisaalta luottamusalue on tällöin kooltaan tasan 95%, jolloin eri menetelmillä johdettujen luottamusalueiden vertaileminen on helpompaa. Myös alueen numeerinen määrittäminen on helpompaa. Huomaa myös, että kaavan (15) nimittäjässä olevan skaalaustekijän parametria  $n$  koskeva integrointi alkaa arvosta 4, kuten mallin määrittelyssä on sovittu.

Suurimman posteritiheyden alue voidaan ratkaista numeerisesti. Kaavan (15) nimittäjänä oleva integraali on skaalaustekijä, joka kertoo koko uskottavuuspinnan rajaaman tilavuden. Kuten aikaisemmin on todettu, on meillä käytössä pinnan arvoja vain tietyn hilan alueelta, mutta uskottavuusfunktion arvo tippuu noltaan hyvissä ajoin ennen hilan rajoja, joten tilavuus saadaan arvioitua tarkasti laskemalla se vain tässä hilassa. Osoittajan integraalissa oleva luottamusalue  $\Theta$  saadaan haettua muuntamalla integraali muotoon

$$\iint_{\Theta} L(n, \sigma_L; \mathbf{x}) d\sigma_L dn = \int_4^{\infty} \int_0^{\infty} \mathbf{I}_{\{(n, \sigma_L) | L(n, \sigma_L; \mathbf{x}) \geq c\}} L(n, \sigma_L; \mathbf{x}) d\sigma_L dn, \quad (16)$$

ja etsimällä numeerisesti vakio  $c$  jolla ehto (15) toteutuu. Tällöin  $\Theta = \{(n, \sigma_L) | L(n, \sigma_L; \mathbf{x}) \geq c\}$ . Näin laskettuna saatiin rajaksi  $c$  arvo 0,005087, joka logaritmoituna tarkoittaa rajaa -5,281 log-uskottavuusfunktiolle. Arvo on aika paljon uskottavuusosamäärän rajaa -2,996 alempi. Tämä johtuu siitä, että uskottavuusfunktio eli parametrien jakauma ei vastaa kovin hyvin normaaliapproksimaatiota vaan on hieman vino, ja lisäksi jakauman hännät ovat painavammat. Bayesin suurimman posteritiheyden alue vaikuttaa tässä estimoinnissa näistä kahdesta menetelmästä realistisemmalta ja luotettavammalta.



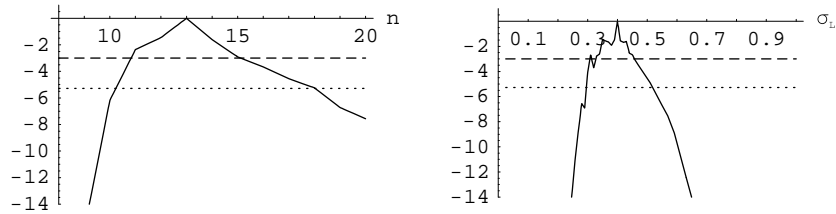
Kuva 9: Parametrin  $\sigma_L$  profiili-log-uskottavuus, kun simuloitua aineistoa on ollut käytössä ensin 5000, sitten 25 000, ja viimein 100 000 havaintoa.

### Estimoinnin tulokset

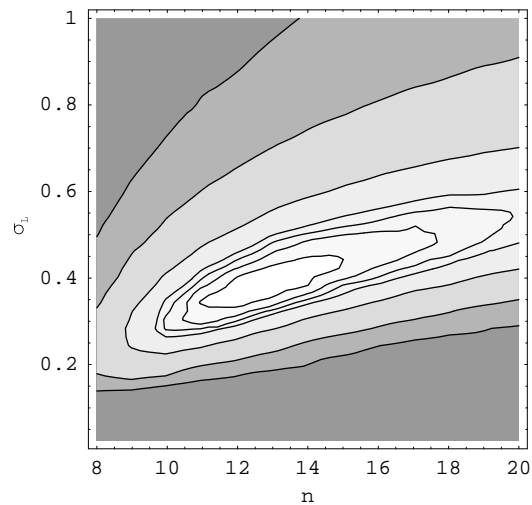
Jos tarkastellaan parametrien profiiliuskottavuusfunktioita<sup>4</sup>, niin parametri  $\sigma_L$  näyttää tarvitsevan enemmän simuloitua aineistoa kuin  $n$ , jotta log-uskottavuus saavuttaisi edes summittaisesti oikean muodon. Kuvassa 9 näkyy funktion heilahtelun tasoittuminen, kun simuloitua aineistoa, jonka perusteella tiheysfunktioestimaatti muodostetaan, lisätään. Lopullisessa tiheysfunktioestimaatissa käytettiin 100 000 havaintoa. Vaikka uskottavuusfunktio vielä heilahtelee jonkin verran, on havaintojen simulointiin ja uskottavuuspinnan muodostamiseen tarvittava tehtävä jo tarpeeksi raskas, joten 100 000 havaintoa riittää.

Uskottavuuspinnan maksimi pysyy tasaisesti samassa kohdassa jo 15 000 simuloidusta havainnosta lähtien, joten piste-estimaatti on varsin luotettava. Suurimman uskottavuuden estimaatiksi saadaan parametrille  $n$  arvo 13, ja parametrille  $\sigma_L$  arvo 0,4. Luottamusväli sen sijaan heilahtelee hiukan enemmän, mutta on kuitenkin jo riittävän tarkasti saatavilla. Kuvassa 10 näkyy parametrien yksiulotteisten 95% luottamusvälien rajat profiili-log-uskottavuusfunktioista: parametrille  $n$  välit ovat  $[10,8; 15,1]$  (uskottavuusosamäärä) ja  $[10,2; 18,0]$  (suur. post.tiheys). Parametrille  $\sigma_L$  välit ovat  $[0,31; 0,46]$  (uskottavuusosamäärä) ja  $[0,30; 0,52]$  (suur. post.tiheys). Kuvassa 11 näkyy parametrien yhteisvaikutus kolmiulotteisen log-uskottavuusfunktion tasa-arvokäyräkuvassa.

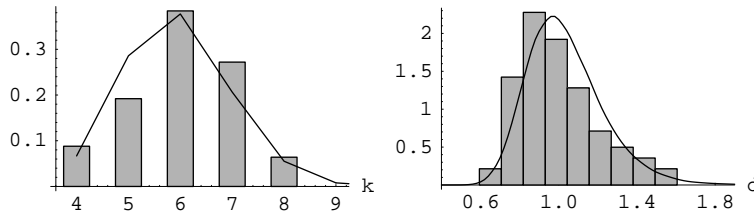
<sup>4</sup> Toisen muuttujan suhteen maksimoitua log-uskottavuusfunktiota, esim. parametrille  $\sigma_L$  se on  $l_{prof} = \max_n l(n, \sigma_L; \mathbf{x})$ .



Kuva 10: Parametrien  $n$  (vasemmalla) ja  $\sigma_L$  (oikealla) profiili-log-uskottavuusfunktiot ja luottamusvälit määräävät rajat, joista ylempi katkoviiva merkitsee uskottavuusosamäärän luottamusväliä ja alempi suurimman posterioritiheyden aluetta, kun havaintoja on simuloitu 100 000 kappaletta.



Kuva 11: Log-uskottavuusfunktion tasa-arvokäyriä. Sisin käyrä rajaa uskottavuusosamäärän luottamusalueen, seuraava suurimman posterioritiheyden alueen. Tämän jälkeen käyrät ovat korkeudella  $-7.5$ ,  $-10$ ,  $-20$ ,  $-50$ ,  $-100$ ,  $-1000$  ja  $-2500$ .



Kuva 12: Tunnuslukujen  $k$  (vasemmalla) ja  $d$  (oikealla) todennäköisyysjakaumat sekä aidossa aineistossa (pylväät) että suurimman uskottavuuden estimaatin mukaisilla parametreilla simuloidussa aineistossa (viivat).

Suurimman uskottavuuden estimaatilla  $n = 13$ ,  $\sigma_L = 0,4$  pitäisi saada aikaan monitahokasmallista partikkeleita, joilla on ainakin mitattujen tunnuslukujen puolesta samat ominaisuudet kuin  $B_4C$ -partikkeleilla. Yksi esimerkki näillä parametrien arvoilla luodusta partikkelista esitettiin jo kuvassa 2. Edellä mainituilla parametrien arvoilla voidaan tarkastella sekä tunnuslukujen havaittua jakaumaa että tämän jakauman ydineestimaattia simuloidun aineiston kautta, ja näin varmentaa mallin oikeellisuus. Kuvassa 5 oli jakauman ydineestimaatti, mutta kuvassa 12 esitetään lisäksi tunnuslukujen yksilotteiset reunajakaumat sekä havaitulle että simuloidulle aineistolle. Yhteensopivuus on silmämääräisesti aivan tyydyttävä, joten malli on järkevä, ja parametrien estimointi on sujunut hyvin.

## 4 Mallin verifiointi polarisaatiomittausten avulla

Johdannossa todettiin, että erityisesti valon polarisaatio on herkkä sirottavan kappaleen muodolle. Tehtäessä valonsirontainversiota käytetään usein juuri havaittua polarisaatiota, johon sitten sovitetaan valonsirontasimulaation tuottamaa polarisaatiota, optimoiden näin kappaleen kokoa, muotoa ja refraktiaindeksiä. Tässä muodon mallintamista on lähestytty suoremmin, mutta polarisaatiomittauksia kannattaa silti käyttää hyväksi. Vaikka

partikkelien muoto on saatu mallinnettua, ei  $B_4C$ -partikkelien refraktioideksiä ole pystytty määrittämään laboratoriomittauksissa. Refraktioideksin määrittäminen voidaan tehdä nyt paljon luotettavammin valonsirontainversiona kun muoto ja koko on kiinnitetty, ja näin ollen inversiossa on enää yksi tuntematon tekijä.

Toinen hyvä puoli polarisaation tutkimisessa on, että se antaa mallin sovituspseuduurista riippumattoman menetelmän tarkastella mallin järkevyyttä verrattuna mitattuihin havaintoihin oikeista partikkeleista. Jos malli on hyvä, saadaan valonsirontainversiosta refraktioideksille joku järkevän kokoluokan luku, ja lisäksi mallin tuottama polarisaatio vastaisi suurin piirtein mitattua.

#### 4.1 Polarisaatiomittaukset

Kuten luvussa 2.1 mainittiinkin, on PROGRA<sup>2</sup>-ryhmä mitannut  $B_4C$ -partikkelien valonsirontaa painottomilla lennoilla. Mittauksissa käytettyä laseria voidaan liikuttaa kaareissa näytteen ympärillä, jolloin valonsirontaominaisuudet saadaan mitattua vaihekulman  $\alpha$  funktiona. Vaihekulma  $\alpha$  tarkoittaa valonlähteen, näytteen ja havaitsijan välistä kulmaa. Valon polarisaatiota mitattaessa on tieto vaihekulmasta välttämätön, koska polarisaatiolla on yleisesti useita vaihekulmasta riippuvia ominaisuuksia. Yksi kiinnostava ominaisuus on lähellä takaisinsirontakulmaa ( $\alpha \rightarrow 0^\circ$ ) usein tapahtuva polarisaation kääntäminen positiivisesta negatiiviseksi. Tämä ominaisuus on lähes universaali ilmeten monilla erikokoisilla ja -tyyppisillä kappaleilla joiden pintaa peittää pölykerros.

Lennoilla käytetyssä laitteistossa on  $B_4C$ -partikkeleita kuvattaessa käytetty punaista laseria, aallonpituudeltaan 632,8 nm. Vaihekulman  $\alpha$  arvolla  $180^\circ$  polarisaatio lähestyy nollaa, samoin takaisinsirontakulmalla  $\alpha = 0^\circ$ , ja takaisinsirontasuunnassa siroavan valon havaintolaitteisto tulee laserin tielle, joten aivan näillä vaihekulman ääripäillä ei ole havaintoja. Mittauksissa käytetyt vaihekulmat eivät ole täsmälleen samoja eri kokoluokkien (9, 13



ja  $88 \mu\text{m}$ ) näytteille. Sekä pienin että suurin mitattu vaihekulma löytyy  $88 \mu\text{m}$  näytteelle:  $15^\circ$  ja  $160^\circ$ . Muut mitatut kulmat ovat tuolla välillä noin 5 - 20 asteen välein. Mitattu polarisaatio eri kokoluokissa esitellään tarkemmin vasta myöhemmin kuvassa 16 yhdessä valonsirontasimulaatioista saadun polarisaatiokäyrän kanssa.

## 4.2 Valonsirontasimulaatiot

### 4.2.1 Simuloinnin toteutus

Mallinnettavien partikkelien kokoparametrit ovat noin 45 ( $9 \mu\text{m}$ ), 65 ( $13 \mu\text{m}$ ) ja 440 ( $88 \mu\text{m}$ ). Tämän kokoluokan partikkelien valonsirontaa voidaan approksimoida geometrisella optiikalla eli säteenseurantakoodeilla (eng. ray tracing, RT) (pienimmän kokoluokan tapauksessa aletaan jo lähestyä rajaa, jossa RT ei ole enää sopiva approksimaatio). RT-koodeissa tutkittavaa kappaletta kohti lähetetään 'valonsäteitä', jotka sitten kohtaavat kappaleen pintaelementin. Pintaelementissä osa säteen intensiteetistä heijastuu peilimäisesti pois päin pintaelementistä, kun taas osa taittuu kappaleen sisälle. Konveksin kappaleen tapauksessa ulos heijastunut osa ei enää kohtaa kappaletta, kun sen sijaan sisään taittunut kohtaa taas jossain vaiheessa kappaleen pintaelementin sisältäpäin. Kappaleen sisälle taittuvan säteen intensiteetti vähenee refraktioidexin imaginaariosan määräämää vauhtia sisällä kuljetun matkan funktiona, joten sisällä kulkevaa sädettä tarvitsee seurata vain tiettyyn rajaan saakka. Kaikkien ulos heijastuvien säteiden Stokesin vektorin komponentit otetaan talteen heijastuskulman funktiona, ja näin saadaan simuloitua kappaleen valonsirontaominaisuudet.

Tässä työssä käytettiin RT-koodia<sup>5</sup> joka on esitelty tarkemmin artikkelissa Macke ym. (1996). Ohjelma ajettiin yksittäiselle partikkelille siten, että sille arvottiin 100 satunnaisorientaatiota, joiden suhteen tulos keskiarvoistettiin. Jokaisella orientaatiolla partikkeliä kohti ammuttiin 1 000 sädettä, joiden

---

<sup>5</sup> Saatavilla osoitteesta <http://www.ifm.uni-kiel.de/fb/fb1/me/research/Projekte/RemSens/SourceCodes/source.html>.

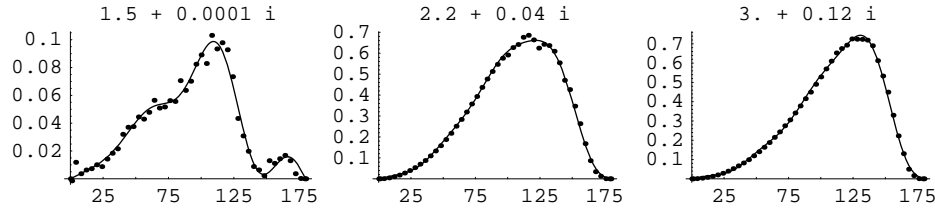
heijastumiset kerättiin talteen siten, että vaihekulman vaihteluväli 0 - 180 astetta oli jaettu neljän asteen mittaisiin alueisiin.

#### 4.2.2 Refraktioindeksin määrittäminen

RT-koodi tarvitsee kappaleen polarisaatiokäyrän muodostamiseen vielä tiedon kappaleen refraktioindeksistä. Koska indeksiä ei tiedetä, optimoidaan polarisaatiota refraktioindeksin funktiona niin, että tulos on mahdollisimman lähellä mitattua polarisaatiota kaikissa kolmessa kokoluokassa. Kuten simuloidun uskottavuudenkin toteutuksessa, niin myös tässä optimoinnissa on ehkä järkevintä jättää varsinaiset optimointialgoritmit rauhaan, ja toteuttaa optimointi raa'alla laskentavoimalla käymällä läpi polarisaatioita refraktioindeksin reaali- ja imaginaariosan muodostamassa kaksiulotteisessa hilassa.

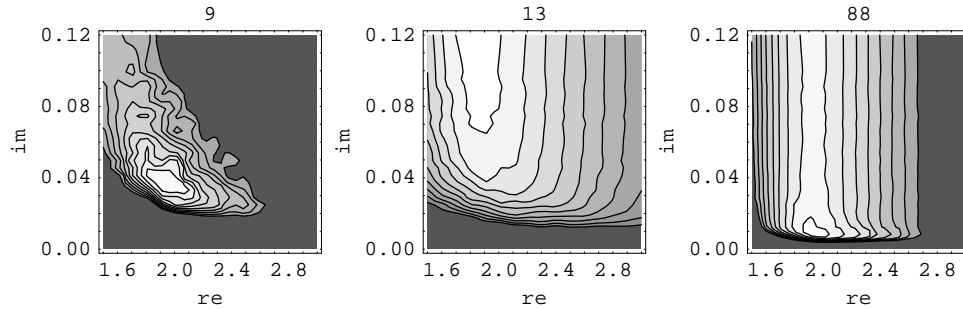
Optimointia varten täytyy ensin miettiä sopiva etäisyysmitta havaitun ja simuloidun polarisaation välillä. Yleisesti käytettynä menetelmänä pienimmän neliösumman etäisyys on hyvä valinta. Havaittu polarisaatio on mitattu vain tietyillä diskreeteillä kulmilla, joten neliölliset etäisyydet on syytä laskea juuri näiden kulmien kohdalla. Myös RT-koodi antaa neljän asteen välein diskreetoitua aineistoa. Olisi periaatteessa mahdollista saada RT-simulaatiosta ulos polarisaatio juuri samoille vaihekulman arvoille kuin havaitussa aineistossa, mutta houkuttelevampi ajatus on sovittaa diskreetin simuloidun aineiston tilalle jatkuva polarisaatiokäyrä, jolloin sovitteesta saadaan laskettua polarisaatio mille tahansa kulmalle. Koska simulaation tuloksena saadut diskreetit pisteet polarisaatiokäyrällä näyttävät sisältävän vain vähän satunnaisvaihtelua, saadaan sopivalla käyränsovituksella käytettyä hyväksi koko neljän asteen välein laskettu simulaatioaineisto. Sovitetulla käyrällä on yksittäisiä pisteitä pienempi satunnaisvaihtelu.

Sopivaa sovitekäyrää mietittäessä on syytä kiinnittää huomiota polarisaation tunnettuihin ominaisuuksiin: se lähestyy nollaa kun vaihekulma  $\alpha \rightarrow 0$



Kuva 13: RT-simulaation tulokset (pisteet) ja niihin pienimmän neliösumman menetelmällä sovitettu kahdeksan termin sinisarja vaihekulman  $\alpha$  funktiona. Kaikki käyrät ovat  $9 \mu\text{m}$  kokoisille partikkeleille, mutta kolmella eri refraktioindeksillä. Kuvasta näkyy, että sinisarja tasoittaa sopivasti satunnaisvaihtelua vasemmanpuoleisimmalla refraktioindeksillä, eikä heilahtele liikaa tasaisillakaan aineistoilla kuten oikeanpuoleisimmassa kuvassa.

tai  $\alpha \rightarrow 180^\circ$ . Nämä ominaisuudet toteuttaa esimerkiksi sinisarja  $s(\alpha) = \sum_{k=1}^{\infty} \beta_k \sin(k\alpha)$ , jonka voi katkaista johonkin sopivaan kohtaan  $K$ . Koska sarjan funktiot muodostavat ortogonaalisen kannan välillä  $[0, \pi]$ , voi sarjan kertoimet  $\beta_k$  helposti sovittaa simulaatioaineistoon pienimmän neliösumman menetelmällä. Kuvassa 13 näkyy muutama esimerkki RT-simulaation tuloksista ja niihin sovitetusta sinisarjasta, kun sarja on katkaistu kohtaan  $K = 8$ . Neliölliset erotukset havaintojen ja RT-simulaation tulosten välillä on laskettu käyttäen erilaisia refraktioindeksejä: indeksin reaaliosa on vaihdellut välillä  $[1,5; 3]$   $0,1:n$  suuruisilla askeleilla, kun taas imaginaariosaan on käytetty ensin arvoa  $0,0001$ , sen jälkeen väli  $[0,005; 0,1]$  on käyty tasavälisin  $0,005:n$  mittaisiin askeleihin, ja viimeiseksi on käytetty arvoa  $0,12$ . Näiden pisteiden muodostamassa hilassa on siis  $352$  refraktioindeksiä, joilla polarisaatiokäyrä täytyy laskea. Tätä varten arvottiin  $500$  kappaletta monita-hokasmallin partikkeleita käyttäen estimoinnin antamia arvoja mallin parametreille. Kullakin refraktioindeksillä valittiin näistä  $500$  partikkelista  $30$  suuruinen joukko ilman takaisinpanoa, joille RT-koodi ajettiin. Varsinainen tulos yhdelle refraktioindeksille on näiden  $30$  partikkelin polarisaatioista keskiarvoistettu polarisaatio. Kun jokaisella refraktioindeksillä käytetään



Kuva 14: Rmse-erot havaitun aineiston polarisaatioon refraktioidexsin reaali (re) ja imaginaariosan (im) funktiona kaikille kolmelle kokoluokalle. Tasa-arvokäyrät on piirretty tasavälisesti suhteessa kunkin kokoluokan rms-arvojen minimiin, ensimmäinen käyrä on 10 prosenttia rmse-minimiä korkeammalla, seuraavat siitä 20 prosenttiyksikön korotuksin aina 190 prosenttiin rmse-minimistä.

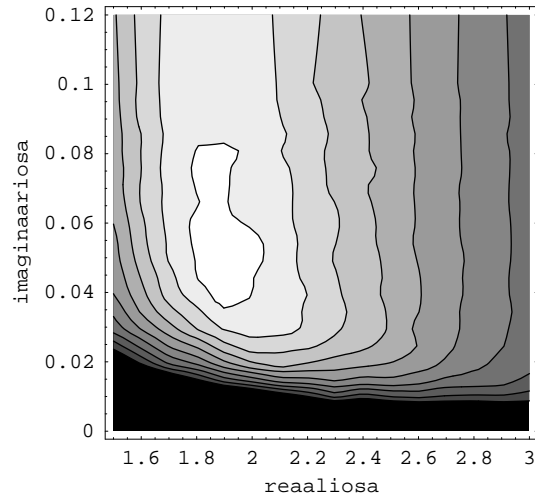
eri partikkelikokoelmaa, saadaan tulokseen mukaan myös partikkelien muodon satunnaisvaihtelun aiheuttama vaikutus mikä on läsnä myös aidoissa mittauksissa. Nämä 352 kappaletta kolmenkymmenen partikkelin RT-ajoa laskettiin jokaiselle kolmesta kokoluokasta. Kuvassa 14 on näiden laskujen tulokset: havaitun ja simuloidun aineiston välisistä neliösummista lasketut rmse-arvot (root mean square error) kaikille kokoluokille erikseen.

Kuvasta 14 nähdään koon vaikutus polarisaation herkkyyteen refraktioidexsin imaginaariosan suhteen. Isommilla 13 ja 88  $\mu\text{m}$  partikkeleilla sisään taittunut säde ehtii joka tapauksessa kulkea partikkelin sisällä sen verran pitkän matkan, että jo suhteellisen pienillä imaginaariosan arvoilla sen intensiteetti heikkenee olemattomiin ennenkuin se pääsee ulos partikkelista. Siten polarisaatio ei enää juurikaan muutu kun imaginaariosaa nostetaan. Sen sijaan pienin 9  $\mu\text{m}$ :n kokoluokka on vielä herkkä myös imaginaariosan muutoksille. Refraktioidexsin reaaliosta eri kokoluokat antavat kuitenkin melko yhtenevää informaatiota.

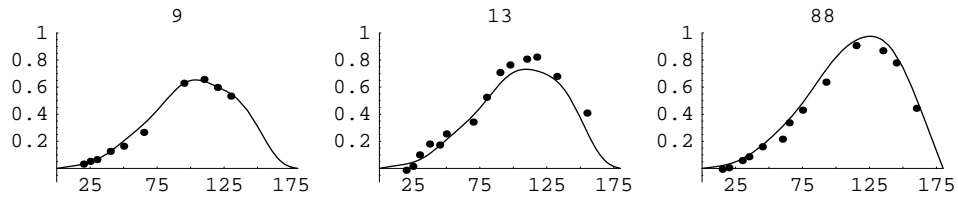
Koska on syytä olettaa, että  $\text{B}_4\text{C}$ -materiaalin refraktioidexsi on vakio kai-

kissa kokoluokissa, pitää kolmen kokoluokan tulokset yhdistää, ja etsiä kaikille luokille sopiva yhteisen minimin antava refraktiaindeksi. Eri kokoluokille on eri määrä havaintoja eivätkä havaintovaihekulmatkaan ole samoja, joten yhdistämiseen on ainakin kaksi kilpailevaa vaihtoehtoa. Kaikki neliölliset erotukset voitaisiin laskea yhteen ja jakaa kaikkien havaintojen määrällä ja ottaa neliöjuuri, jolloin jokainen havainto vaikuttaisi yhtä suurella painolla, mutta suuremmat kokoluokat ( $13 \mu m$ , 14 havaintoa ja  $88 \mu m$ , 13 havaintoa) saisivat suuremman painon kuin pienin kokoluokka ( $9 \mu m$ , 10 havaintoa). Toinen vaihtoehto on jakaa ensin jokaisen kokoluokan neliösummat havaintojen määrällä ja laske vasta nämä kolme summaa yhteen ja ottaa neliöjuuri, jolloin jokainen kokoluokka saa yhtä suuren painon. Tässä työssä on päädytty kannattamaan jälkimmäistä lähestymistapaa, koska pienin kokoluokka on myös herkin refraktion imaginaariosalle, ja tätä herkkyyttä menettäisiin hieman jos jokaiselle havainnolle annettaisiin yhtä suuri paino, kun pienimmän kokoluokan havaintoja sattuu olemaan vähiten. Minimi löytyy onneksi kummallakin tavalla samasta kohtaa, refraktiaindeksillä  $1,9 + 0,04i$ .

Kuvassa 15 näkyy havaintojen ja mallin yhteensopivuus yhtäaikaaisesti kaikilla kokoluokilla. Mitään varsinaista luottamusväliä minimille on vaikea antaa, mutta kuvasta kannattaa seurata vaikkapa ensimmäisen tasa-arvokärän rajaamaa aluetta jossa yhteensopivuus on vielä melko hyvä. Tämä alue kattaa reaaliosan puolesta noin välin  $[1,8; 2,05]$  ja imaginaariosan puolesta välin  $[0,035; 0,08]$ . Kuvassa 16 näkyy  $B_4C$ -partikkeleille sovitetulla refraktiaindeksillä ja monitahokasmallilla saatu polarisaatiokäyrä yhdessä havaintojen kanssa, ja yhteensopivuus on varsin tyydyttävä. Malli on siis saatu käyttäytymään sekä muotonsa että valonsirontaominaisuuksiensa puolesta kuten aidot  $B_4C$ -partikkelit.



Kuva 15: Kaikkien kolmen kokoluokan yhdistetty yhteensopivuus havaittuun aineistoon refraktioindeksin reaali- ja imaginaariosan funktiona. Minimi löytyy refraktioindeksillä  $1,9 + 0,04i$ . Ensimmäinen käyrä on taas 10 prosentin korotus yhdistettyyn rmse-minimiin, ja seuraavat siitä 20 prosenttiyksikön välein aina 190 prosenttiin.



Kuva 16: Sovitetulle mallille refraktioindeksillä  $1,9 + 0,04i$  simuloitu polarisaatiokäyrä eri kokoluokissa yhdessä PROGRA<sup>2</sup>-projektin oikeista B<sub>4</sub>C-partikkeleista mittaaman polarisaation (pisteet) kanssa.

## 5 Päätelmät

Simuloitu suurimman uskottavuuden menetelmä on selvästi erittäin lupaava ja käyttökelpoinen mitä erilaisimpiin parametrien estimointitehtäviin. Sen käyttö ei ole mitenkään rajoittautunut juuri tässä työssä esiteltyyn ongelmaan, vaan sitä voi pitää laajenuksena tavalliseen suurimman uskottavuuden menetelmään. Kaikki normaalit suurimman uskottavuuden estimoinnit voidaan suorittaa simuloimalla ja saada approksimatiivisia tuloksia, mutta myös monet tavallisen uskottavuuspäätelyn ulottumattomissa olevat ongelmat ratkeavat simuloimalla.

Simuloitu uskottavuuspäätely tarjoaa myös joitain mahdollisuuksia verrata eri estimaattorien suorituskykyä ja ylipäätensä mallin identifioituvuutta. Kahta kilpailevaa harhatonta estimaattoria vertaillaan yleensä varianssien avulla. Varianssit näkyvät kuitenkin myös log-uskottavuusfunktion kaarevuudessa suurimman uskottavuuden arvion kohdalla. Simuloidusta log-uskottavuusfunktioista voidaan siten helposti todeta, kummalla kahdesta (tai useammasta) kilpailevasta estimaattorista on pienin varianssi eli suurin kaarevuus. Estimaattorin asemassa tässä voi olla esimerkiksi kaksi vaihtoehtoa, hieman eri tavalla mitattavaa tunnuslukua.

Myös tunnusluvun eli mitattavan suureen järkevyyttä näkyy log-uskottavuusfunktiossa. Jos malli ei ole identifioituva eli tunnusluku ei kerro mitään mallin parametreista, on funktio silloin vakioarvoinen. Tässä tapauksessa myös simuloitu uskottavuus olisi hyvin lähellä vakiofunktioita, eikä maksimoinnissa löydetäisi luotettavaa maksimia. Tämän työn tapauksessa log-uskottavuudesta voidaan nähdä, että valitut tunnusluvut ovat järkeviä ja antavat informaatiota mallin parametreista, koska uskottavuudella on selkeä yksi globaali maksimi.

Tässä työssä esitellyllä menetelmällä on hyviä sovelluskohteita pienten partikkelien valonsirontatutkimuksessa. Työn pohjalta on jo yhteistyössä ranskalaisen PROGRA<sup>2</sup>-ryhmän kanssa valmisteltu artikkeli, joka on hyväksytty julkaistavaksi arvovaltaiseen alan lehteen (JQSRT, Journal of Quantitative

Spectroscopy and Radiative Transfer). Eteenkin refraktiaindeksin arvioinnissa tällä menetelmällä voi olla paljon sovelluskohteita.

Juuri B<sub>4</sub>C-partikkelien mallinnus on jo nyt antanut hyödyllisiä tuloksia myös käytännössä. Alunperin polarisaatiokäyrä 13  $\mu\text{m}$  kokoluokassa *PROGRA*<sup>2</sup>-ryhmän mittauksissa oli hieman epäjohdonmukainen verrattuna kokoluokkiin 9 ja 88  $\mu\text{m}$ . Kun tässä esitelty teoreettinen malli antoi polarisaatiosta hyvin erilaiset tulokset, ranskalaiset lennättivät 13  $\mu\text{m}$ :n kokoluokan uudelleen. Uudet tulokset olivat jo paremmin ymmärrettäviä, mutta vieläkin teoreettinen malli näytti, että 9 ja 88  $\mu\text{m}$ :n kokoluokat sopivat hyvin yhteen, kun taas 13  $\mu\text{m}$ :n kokoluokka käyttäyi omalaatuisesti. Uusi tarkempi tarkastelu osoitti, että ranskalaisten käyttämässä näytteessä oli mukana muutamia paljon 13  $\mu\text{m}$ :n kokoluokkaa suurempia partikkeleita, jotka aiheuttivat vääristymää polarisaatioon. Kun tämä korjattiin, saatiin vihdoin hyvä yhteensopivuus sekä teoreettisen mallin että havaintojen kanssa. Vaikka havainnoissa oli jo etukäteen havaittu ongelmia, ei luultavasti ainakaan toista korjauskierrosta olisi osattu tehdä ilman teoreettisen mallin tukea.



## Viitteet

- Barber, C. B. – Dobkin, D. P. – Huhdanpää H. (1996): The Quickhull algorithm for convex hulls. *ACM Transactions on Mathematical Software* 22:469–483.
- Brandt, M. – Santa-Clara, P. (2002): Simulated likelihood estimation of diffusions with an application to exchange rate dynamics in incomplete markets. *Journal of Financial Economics* 63:161–210.
- Bohren, C. – Huffman, D. (1983): *Absorption and Scattering of Light by Small Particles*. John Wiley & Sons, USA.
- Deans, S. (1983): *The Radon Transform and Some of Its Applications*. John Wiley & Sons, USA.
- Ekholm, A. (1997): Johdatus uskottavuuspäätelyyn. Opetusmoniste. Tilastotieteen laitos, Helsingin yliopisto.
- Electromagnetic Scattering by Nonspherical Particles (1999): *Journal of Geophysical Research* 104.
- Holmström, L. (2002): Funktioiden estimointi. Luentomoniste. Rolf Nevalinna -instituutti.
- Liu, Y. – Arnott, P. – Hallet, J. (1999): Particle size distribution retrieval from multispectral optical depth: Influences of particle nonsphericity and refractive index. *Journal of Geophysical Research* 104:31753–31762.
- Lumme, K. – Rahola, J. – Muinonen, K. – Volten H. (1995): Scattering by Rough Particles and Stochastic Aggregates. In: 4th International Congress, Optical Particle Sizing 583–592. NürnbergMesse GmbH, Nürnberg.
- Lumme, K. – Rahola, J. (1998): Comparison of light scattering by stochastically rough spheres, best-fit spheroids and spheres. *Journal of Quantitative Spectroscopy & Radiative Transfer* 60:439–450.

- Macke, A. – Mueller, J. – Raschke, E. (1996): Single scattering properties of atmospheric ice crystals. *Journal of the Atmospheric Sciences* 53:2813–2825.
- Mackowski, D. – Mishchenko, M. (1996): Calculation of the  $T$  matrix and the scattering matrix for ensembles of spheres. *Journal of Optical Society of America* 13:2266–2278.
- Mishchenko, M. – Hovenier, J. – Travis, L. (ed) (2000): *Lighth Scattering by Nonspherical Particles: Theory, Measurements and Applications*. Academic Press, USA.
- Muinsonen, K. (1986): Klassinen sähkömagneettinen sironta epäsäännöllisistä hiukkasista. Pro gradu -tutkielma. Teoreettisen fysiikan laitos, Helsingin yliopisto.
- Parzen, E. (1962): On Estimation of a Probability Density Function and Mode. *Annals of Mathematical Statistics* 33:1065–1076.
- Penttinen, A. (1984): Modelling interactions in spatial point patterns: parameter estimation by the maximum likelihood method. Ph.D. Dissertation. Jyväskylä studies in computer science, economics and statistics 7.
- Robert, C. P. – Casella, G. (1999): *Monte Carlo Statistical Methods*. Springer-Verlag, New York.
- Rosenblatt, M. (1956): Remarks on some Nonparametric Estimates of a Density function. *Annals of Mathematical Statistics* 27:832-835.
- Santa-Clara, P. (1995): Simulated likelihood estimation of diffusions with an application to the short term interest rate. Teoksessa: *Essays on the Theory and Estimation of Term Structure Models*. INSEAD.
- Silverman, B. W. (1989): *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, London.

- Tanner, M. (1993): *Tools for Statistical Inference: Methods for the Exploration of Posterior Distributions and Likelihood Functions*. 2nd ed., Springer-Verlag, New York.
- Webb, A. (1999): *Statistical Pattern Recognition*. Arnold, London.
- van de Weygaert, R. (1991): Voids and the geometry of large scale structure. Ph.D. Dissertation. University of Leiden, Netherlands.
- Worms, J-C. – Renard, J-B. – Hadamcik, E. – Levasseur-Regourd, A-C. (1999): Results of the PROGRA<sup>2</sup> Experiment: An Experimental Study in Microgravity of Scattered Polarized Light by Dust Particles with Large Size Parameter. *Icarus* 142:281–297.
- Worms, J-C. – Renard, J-B. – Hadamcik, E. – Brun-Huret, N. – Levasseur-Regourd, A-C. (2000): Light scattering by dust particles with the PROGRA<sup>2</sup> instrument - comparative measurements between clouds under microgravity and layers on the ground. *Planetary and Space Science* 48:493–505.

## A Tasainen jakauma pallokoordinaateissa

Tasainen jakauma pallokoordinaateissa tarkoittaa sitä, että origosta lähtevät pallokoordinaattisuunnat  $(\theta, \phi)$  ovat jakautuneet tasaisesti pallon pinnalle. Tasaisessa kaksiulotteisessa jakaumassa tietyn alueen todennäköisyys on aina suoraan verrannollinen alueen pinta-alaan, mutta pallokoordinaateissa tämä pinta-ala lasketaan pallon pinnalta. Lähestytään satunnaismuuttujien  $(\Theta, \Phi)$  yhteisjakaumaa  $f_{\Theta, \Phi}$  kertymäfunktion ja pallon pinnalla olevan alueen alasta lähtien.

Lähdetään riippumattomuudesta. Jos ja vain jos muuttujat  $\Theta$  ja  $\Phi$  ovat keskenään riippumattomia, niin yhteiskertymäfunktio separoituu, eli  $F_{\Theta, \Phi}(\theta, \phi) = F_{\Theta}(\theta) F_{\Phi}(\phi)$ . Toisaalta taas tasainen jakauma merkitsi sitä, että kertymäfunktio on suoraan verrannollinen kertyneen alueen pinta-alaan. Navalta lähtevän pallosegmentin osan  $\Theta \in (0, \theta)$ ,  $\Phi \in (0, \phi)$  ala  $\omega$  pallon pinnalla on

$$\omega = 2 \sin^2\left(\frac{\theta}{2}\right) \phi, \quad (17)$$

joten kertymäfunktio separoituu halutulla tavalla:

$$F_{\Theta, \Phi}(\theta, \phi) \propto \omega \Rightarrow F_{\Theta}(\theta) \propto \sin^2\left(\frac{\theta}{2}\right) \text{ ja } F_{\Phi}(\phi) \propto \phi. \quad (18)$$

Samalla nähdään, että muuttujan  $\Phi$  jakauma on tasainen jakauma, ja muuttujan  $\Theta$  jakauma saadaan derivoimalla kertymäfunktio

$$F_{\Theta}(\theta) \propto \sin^2\left(\frac{\theta}{2}\right) \Rightarrow f_{\Theta}(\theta) \propto \frac{\sin \theta}{2}. \quad (19)$$

Satunnaissuuntien arpominen tasaisesta jakaumasta pallokoordinaateissa on nyt helppoa. Riippumattomuuden ansiosta voidaan  $\phi$  arpoa suoraan tasaisesta jakaumasta väliltä  $[0, 2\pi]$ . Kulman  $\theta$  arpomiseen tarvitaan ensin jakaumamuunnos  $Y = \frac{1}{2}(1 - \cos(\Theta))$ , jotta päästään tasaiseen jakaumaan. Todistetaan seuraavaksi muunnos. Muunnos on määrittelyalueellaan  $[0, \pi]$  monotonisesti kasvava funktio, joten uuden muuttujan  $Y$  jakaumaksi saadaan

$$f_Y(y) = f_{\Theta}(\theta(y)) \frac{d\theta(y)}{dy} \propto \frac{\sin(\arccos(2y-1))}{2} \frac{1}{\sqrt{(1-y)y}} = 1. \quad (20)$$

Kun uuden muuttujan  $Y$  määrittelyalue on väli  $[0, 1]$ , on  $Y$  siis tasajakautunut. Käänteismuunnoksella  $\theta := \theta(y) = \arccos(2y - 1)$  saadaan tasajakaumasta oikeaa jakaumaa noudattava kulma  $\theta$  arvottua.

## B Lähdekoodi monitahokasmallin realisaatioiden arpomiseen

Seuraava lähdekoodi toimii Mathematica-ohjelman versiossa 4.

```
Needs[ "Statistics`ContinuousDistributions`"];
Needs[ "Calculus`VectorAnalysis`"];
Needs[ "DiscreteMath`Combinatorica`"];
Needs[ "DiscreteMath`ComputationalGeometry`"];

(* Johdetaan thetan kertymäfunktion käänteisfunktio
   satunnaislukujen arvontaa varten. *)
f[theta_] := Sin[theta]/2;
F[theta_] = Integrate[ f[x], {x, 0, theta}];
ans = Solve[ x == F[y], y];
Finv[x_] = (y /. ans[[2, 1]]);

(* Reparametroitu log-normaalijakauma odotusarvolla 1 *)
nMu[sigma_] := -1/2 Log[sigma^2 + 1];
nSigma[sigma_] := Sqrt[ Log[sigma^2 + 1]];
rf[sigma_] := LogNormalDistribution[nMu[sigma], nSigma[sigma]]

(* Yhden monitahokkaan luova ohjelma, argumenttina
   parametrit n ja sigma, tuloksena tunnusluvut k ja d,
   eli projektion kulmien määrä ja maksimilävistäjä *)
proj[ n_, sigma_] :=
Module[ {säteet, suunnat, data, p, k, d},
  (* arvottaa säteet ja suunnat *)
  säteet = RandomArray[ rf[ sigma], n];
  suunnat =
  Table[ {Finv[ Random[]], Random[Real, {0, 2 Pi}]}, {n}];
  (* koordinaattimuunnos *)
  data =
  Table[ CoordinatesToCartesian[ {säteet[[i]], suunnat[[i, 1]],
    suunnat[[i, 2]]}, Spherical], {i, n}];
  (* projektio ja konvekssi verho *)
  p = data[[All, {1, 2}]];
  p = Part[ p, ConvexHull[ p]];
  (* tunnusluvut *)
  k = Length[ p];
  p = Map[ Function[ x, Apply[ Subtract, x]], KSubsets[ p, 2]]^2;
  d = Sqrt[ Max[ Map[ Function[ x, Apply[ Plus, x]], p]]];
  {n, sigma, k, d}
]
```

## C Rosenblatin ehtojen todistus ydineestimaatille

Todistetaan väite Gaussin ytimen ja kaavan (7) mukaisen silotusparametrin muodostaman ydineestimaatin asymptoottinen harhattomuus ja konsistenssi Rosenblatin ehdoilla.

Väite (9a):

$$\int_{-\infty}^{\infty} K(x)dx = 1 \quad (21)$$

Todistus:

Ydin on normaalijakauman tiheysfunktio, ja siten pinta-alaltaan yksi.

Väite (9b):

$$\sup_x K(x) < \infty \quad (22)$$

Todistus:

Normaalijakauman tiheysfunktio on äärellinen, kunhan varianssi  $\sigma^2 > 0$ . Havaintoja täytyy siis olla useampi kuin yksi, jotta varianssi ei olisi nolla. Todennäköisyys että varianssi on nolla, kun havaintoja on kaksi tai enemmän on nolla, ja tapauksena muutenkin patologinen. Väite on siis voimassa todennäköisyydellä yksi kun otos on kooltaan kaksi tai enemmän, kuten tiheysfunktion estimoinnissa kuuluukin.

Väite (9c):

$$\lim_{x \rightarrow \infty} xK(x) = 0 \quad (23)$$

Todistus:

Käytetään l'Hospitalin sääntöä  $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = \lim_{x \rightarrow \infty} \frac{f'(x)}{g'(x)}$ , ja valitaan  $f(x) = x$  ja  $g(x) = K(x)^{-1}$ . Nyt  $f'(x) = 1$  ja

$$g'(x) = (x - \mu) \frac{\sqrt{2\pi}}{h} e^{\frac{1}{2}(\frac{x-\mu}{h})^2}, \quad (24)$$

mikä lähestyy ääretöntä, kun  $x$  kasvaa rajatta. Siten derivaattojen suhde lähestyy nollaa, kun  $x$  kasvaa rajatta.

Väite (9d):

$$\lim_{n \rightarrow \infty} h(n) = 0 \quad (25)$$

Todistus:

Silotusparametrin kaavassa  $h(n) = \sigma \left( \frac{4}{p+2} \right)^{\frac{1}{p+4}} n^{-\frac{1}{p+4}}$  hajonta  $\sigma$  on äärellinen ja  $n^{-\frac{1}{p+4}}$  lähestyy nollaa, joten  $h(n)$  lähestyy nollaa.

Väite (9e):

$$\lim_{n \rightarrow \infty} nh(n) = \infty \quad (26)$$

Todistus:

Samassa silotusparametrin kaavassa hajonta on taas äärellinen, mutta  $n$ :stä riippuva osa  $n^{\frac{p+3}{p+4}}$  lähestyy ääretöntä, joten  $nh(n)$  lähestyy ääretöntä.



## D Simuloidun suurimman uskotavuuden estimoinnin toteutus

Ohjelman lähdekoodi noudattaa FORTRAN 90 -standardia.

```
PROGRAM likelihood
! Maximum likelihood - ohjelma, kevät 2002, Antti Penttilä
! Laskee ML-pinnan, kun annetaan otosdata ja simuloitu mallidata.

IMPLICIT NONE
INTEGER, PARAMETER :: statN = 2, fprec = SELECTED_REAL_KIND(12)
INTEGER, DIMENSION(:, :, :), ALLOCATABLE :: histo
INTEGER, DIMENSION(:), ALLOCATABLE :: nVal, kVal
INTEGER :: simuN, sampleN, nN, sN, kN, i, j, k, stat, ti1, ti2, ti3, ii, &
& jj, kk
REAL (KIND=fprec), PARAMETER :: M_PI = 3.141592653589793_fprec
REAL (KIND=fprec), DIMENSION(:, :, :), ALLOCATABLE :: pSample
REAL (KIND=fprec), DIMENSION(:, :), ALLOCATABLE :: sampleD, MLD, sigmaD, &
& meanD, simuD
REAL (KIND=fprec), DIMENSION(:), ALLOCATABLE :: sVal
REAL (KIND=fprec) :: tf, h, globmax
CHARACTER (LEN=60), PARAMETER :: outfileP = 'pSample.dat', &
& outfileML = 'ML.dat'
CHARACTER (LEN=200) :: simufile, samplefile
LOGICAL :: status

! Tiedostojen nimet komentoriviltä
CALL GETARG (1, simufile)
CALL GETARG (2, samplefile)

! Luetaan taulukot ja varataan muistia
OPEN (8, FILE=simufile, ACTION='READ', STATUS='OLD')
INQUIRE(8, OPENED=status)
IF (.NOT. status) THEN
  WRITE (*,*) 'Tiedostoa ', simufile, ' ei saada avatuksi'
  STOP
END IF
READ (8, *) simuN, nN, sN, kN
ALLOCATE (simuD(simuN,2), nVal(nN), sVal(sN), kVal(kN), STAT=stat)
IF (stat /= 0) THEN
  WRITE (*,*) 'Muistin varaus epäonnistui, taulukko'
  WRITE (*,*) simufile
  STOP
END IF
READ (8,*) nVal(:)
READ (8,*) sVal(:)
READ (8,*) kVal(:)
```

```

OPEN (7, FILE=samplefile, ACTION='READ', STATUS='OLD')
INQUIRE(7, OPENED=status)
IF (.NOT. status) THEN
  WRITE (*,*) 'Tiedostoa ', samplefile, ' ei saada avatuksi'
  STOP
END IF
READ (7, *) sampleN
ALLOCATE (sampleD(sampleN,statN), STAT=stat)
IF (stat /= 0) THEN
  WRITE (*,*) 'Muistin varaus epäonnistui, taulukko'
  WRITE (*,*) samplefile
  STOP
END IF

! Logaritmuunnos muuttujalle d
DO i=1,sampleN
  READ(7, *) sampleD(i,:)
  sampleD(i,2) = LOG (sampleD(i,2))
END DO
CLOSE (7)

! Lukumäärät kulmien määrälle k, sekä keskiarvot ja hajonnat
ALLOCATE (histo(nN,sN,kN), meanD(nN,sN), sigmaD(nN,sN), STAT=stat)
IF (stat /= 0) THEN
  WRITE (*,*) 'Muistin varaus epäonnistui, taulukko'
  WRITE (*,*) 'histo ja sigmaD'
  STOP
END IF
histo(:, :, :) = 0
sigmaD(:, :) = 0.0_fprec
meanD(:, :) = 0.0_fprec
! Jokaisen otospisteen todennäköisyys
ALLOCATE (pSample(nN,sN,sampleN), STAT=stat)
IF (stat /= 0) THEN
  WRITE (*,*) 'Muistin varaus epäonnistui, taulukko'
  WRITE (*,*) 'pSample'
  STOP
END IF
pSample(:, :, :) = 0.0_fprec
! ML-pinta
ALLOCATE (MLD(nN,sN), STAT=stat)
IF (stat /= 0) THEN
  WRITE (*,*) 'Muistin varaus epäonnistui, taulukko'
  WRITE (*,*) 'MLD'
  STOP
END IF

globmax = 0.0_fprec
globmax = -1/globmax ! pienin mahdollinen luku

```

```

! Päälooppi yli parametrien n ja sigma_L arvojen
DO i=1,nN
DO j=1,sN

! Luetaan pala simuloitua aineistoa
DO k=1,simuN
READ(8, *) ti1, tf, simuD(k,:)
IF (ti1 /= nVal(i) .OR. tf < sVal(j)-0.00001_fprec .OR. &
& tf > sVal(j)+0.00001_fprec) THEN
WRITE (*,*) 'Problem on line', i*j*k
WRITE (*,*) 'should be', nVal(i), sVal(j)
WRITE (*,*) 'is', ti1, tf
CLOSE (8)
STOP
END IF
END DO

! Skaalataan lävistäjä arvoon 1 ja logaritmoidaan
simuD(:,2) = simuD(:,2) * simuN / SUM (simuD(:,2))
simuD(:,2) = LOG (simuD(:,2))

! k:n arvojen todennäköisyydet, d_l:n keskiarvo ja hajonta
DO k=1, simuN
ti1 = MINLOC (ABS (kVal-simuD(k,1)),1)
histo(i,j,ti1) = histo(i,j,ti1) + 1
meanD(i,j) = meanD(i,j) + simuD(k,2)
sigmaD(i,j) = sigmaD(i,j) + simuD(k,2)**2
END DO
meanD(i,j) = meanD(i,j) / simuN
sigmaD(i,j) = SQRT (sigmaD(i,j) / simuN - meanD(i,j)**2)

! Kernel-estimaatti
DO k=1, sampleN
DO ii=1, simuN
IF (simuD(ii,1) == sampleD(k,1)) THEN
ti1 = MINLOC (ABS (kVal-sampleD(k,1)),1)
h = sigmaD(i,j) * (4.0_fprec/3.0_fprec)**(1.0_fprec/5.0_fprec) &
& * histo(i,j,ti1)**(-1.0_fprec/5.0_fprec)
tf = EXP (-0.5_fprec * ((sampleD(k,2) - simuD(ii,2)) / h)**2) / &
& (sqrt(2*M_PI) * h * simuN)
pSample(i,j,k) = pSample(i,j,k) + tf
END IF
END DO
END DO

! ML-pinta
MLD(i,j) = SUM (LOG (pSample(i,j,:)))
IF (MLD(i,j) > globmax) THEN

```

```

        globmax = MLD(i,j)
    END IF

    END DO
END DO
CLOSE (8)

MLD(:, :) = MLD(:, :) - globmax ! Skaalataan ML-pinta

! Tulokset tiedostoon
OPEN (7, FILE=outfileP, ACTION='WRITE', STATUS='REPLACE')
INQUIRE (7, OPENED=status)
IF (.NOT. status) THEN
    WRITE (*,*) 'Tiedostoa ', outfileP, ' ei saada avatuksi'
    STOP
END IF
WRITE (7, '(3I5)') nN, sN, sampleN
DO i=1, nN
    WRITE (7, '(I5)', ADVANCE='NO') nVal(i)
END DO
WRITE (7, *) ''
DO i=1, sN
    WRITE (7, '(F6.3,X)', ADVANCE='NO') sVal(i)
END DO
WRITE (7, *) ''
DO i=1, sampleN
    WRITE (7, '(F5.1,F6.3)') sampleD(i,:)
END DO
DO i=1, nN
    DO j=1, sN
        DO k=1, sampleN
            WRITE (7, '(E30.16E3)', ADVANCE='NO') pSample(i,j,k)
        END DO
    WRITE (7, *) ''
    END DO
END DO
CLOSE (7)

OPEN (7, FILE=outfileML, ACTION='WRITE', STATUS='REPLACE')
INQUIRE (7, OPENED=status)
IF (.NOT. status) THEN
    WRITE (*,*) 'Tiedostoa ', outfileML, ' ei saada avatuksi'
    STOP
END IF
WRITE (7, '(3I5)') nN, sN, sampleN
DO i=1, nN
    WRITE (7, '(I5)', ADVANCE='NO') nVal(i)
END DO
WRITE (7, *) ''

```

```
DO i=1, sN
  WRITE (7, '(F6.3,X)', ADVANCE='NO') sVal(i)
END DO
WRITE (7, *) ''
DO i=1, sampleN
  WRITE (7, '(F5.1,F6.3)') sampleD(i,:)
END DO
DO i=1, nN
  DO j=1, sN
    WRITE (7, '(ES24.15)') MLD(i,j)
  END DO
END DO
CLOSE (7)

END PROGRAM likelihood
```