

Ostotapahtuman tekevien pelaajien määrän ennustaminen  
elinaika-analyysia ja EM-algoritmia soveltamalla

Riikka Numminen

Pro gradu -tutkielma  
Toukokuu 2018

MATEMATIIKAN JA TILASTOTIETEEN LAITOS  
TURUN YLIOPISTO



TURUN YLIOPISTO

Matematiikan ja tilastotieteen laitos

NUMMINEN, RIIKKA: Ostotapahtuman tekevien pelaajien määrän ennustaminen elinaika-analyysia ja EM-algoritmia soveltamalla

Pro gradu -tutkielma, 61 s.

Sovellettu matematiikka

Toukokuu 2018

---

Tässä työssä tutkitaan menetelmää, jolla pyritään ennustamaan rajoitetulla aikavälillä kerätystä mobiilipelien pelaajadatasta ostotapahtuman tekevien pelaajien määrä. Kirjallisuuskatsauksessa esitellään elinaika-analyysin käsitteitä ja yleisimpiä malleja sekä johdettavan mallin optimiratkaisun selvittämiseen käytettävä Expectation Maximization -algoritmi. Ostotapahtuman tekevien pelaajien määrän ennustamisessa käytetään elinaika-analyysin mixture cure -mallia. Teoriaosuuden lopuksi johdetaan mallille funktiot, jotka tarvitaan EM-algoritmin käyttämiseen. Menetelmä implementoidaan R-ohjelmointikielellä.

Mallia käytetään Hipster Sheep -nimisestä pelistä [1] kerättyyn dataan ja se toimii pääasiassa odotusten mukaisesti. Joillain harvoilla otoksilla havaittiin ongelma mixture cure -mallin ja eksponenttijakauman toisistaan erottamisessa. Ongelman aiheuttavasta datasta ei pystytty erottamaan, tekevätkö kaikki pelaajat ostotapahtuman (eksponenttijakauma hitaalla ostotahdilla) vai onko datassa pelaajia, jotka eivät sitä tee, vaikka niitä seurattaisiin äärettömän pitkään (mixture cure -malli kohtuullisella ostotahdilla).

Lisäksi tutkitaan generoidulla datalla otoskoon ja ostotapahtuman tekevien pelaajien määrän vaikutusta mallin toimivuuteen. Huomataan, että malli toimii sitä paremmin, mitä enemmän pelaajia on, mitä useampi pelaaja tekee ostotapahtuman ja mitä isompi osa ostotapahtuman tekevistä pelaajista tekee sen tarkasteluhetken mennessä. Liian pienellä otoksella ja liian pienellä ostotapahtuman tekevien pelaajien sekä tarkasteluhetken mennessä tehtyjen ostotapahtumien määrällä mallin antama ratkaisu vaihtelee melko paljon otoksesta riippuen. Yhteenvedossa esitetään kehitysidea, jolla saadaan poistettua edellä mainitut ongelmat.

Avainsanat: mixture cure model, EM-algoritmi, pelidata.



# Sisältö

<b>1</b>	<b>Johdanto</b>	<b>1</b>
<b>2</b>	<b>Elinaika-analyysi</b>	<b>3</b>
2.1	Parametriset mallit . . . . .	7
2.1.1	Eksponttijakauma . . . . .	7
2.1.2	Weibull-jakauma . . . . .	9
2.1.3	Gammajakauma . . . . .	10
2.1.4	Log-normaalijakauma . . . . .	11
2.1.5	Log-logistinen jakauma . . . . .	12
2.1.6	Gompertz-jakauma . . . . .	13
2.2	Epäparametriset mallit . . . . .	15
2.3	Mixture cure -malli . . . . .	17
<b>3</b>	<b>Expectation Maximization -algoritmi</b>	<b>19</b>
3.1	Algoritmin toiminta . . . . .	19
3.2	Algoritmin konvergenssi . . . . .	22
3.2.1	Uskottavuusfunktion arvo suppenee . . . . .	22
3.2.2	Parametrien jono suppenee . . . . .	27
3.2.3	Konvergenssiaste . . . . .	28
<b>4</b>	<b>EM-algoritmin käyttö mixture cure -mallin parametrien las-</b>	
	<b>kemiseen</b>	<b>30</b>
4.1	Malli . . . . .	30
4.2	Toteutus R-ohjelmointikielellä . . . . .	35
<b>5</b>	<b>Pelidataan soveltaminen</b>	<b>39</b>
5.1	Käytetyn datan kuvailu ja muokkaus . . . . .	40
5.2	Mallin oletuksien toteutuminen . . . . .	42
5.3	Mallin käyttäminen . . . . .	45
5.4	Mallin muuttujien vaikutus laskettuihin estimaatteihin . . . . .	52
<b>6</b>	<b>Yhteenvedo</b>	<b>57</b>
	<b>Kiitokset</b>	<b>59</b>
	<b>Kirjallisuutta</b>	<b>60</b>

# 1 Johdanto

Tämän työn tavoite on tutkia elinaika-analyysissa käytetyn mixture cure -mallin toimivuutta pelialalta kerättyyn dataan. Elinaika-analyysi on tilastotieteen menetelmä, jolla pyritään ennustamaan lopputulos rajoitetulla aikavälillä kerätystä datasta. Alttiit yksilöt tekevät tarkastelun kohteena olevan tapahtuman äärellisellä seuranta-ajalla ja immuunit yksilöt eivät sitä tee, vaikka niitä seurattaisiin äärettömän pitkään. Pelidata, johon tätä sovelletaan, on netistä puhelimeen ilmaiseksi ladattavasta pelistä kerättyjä tietoja. Osa pelien ominaisuuksista on saatavilla vain jos niistä maksaa.

Tarkoituksena on ennustaa ostotapahtuman tekevien pelaajien määrä, kun käytössä on rajoitetulla aikavälillä kerättyä dataa. Uusien pelaajien saaminen edellyttää pelin mainostamista, joten mainostuskulujen minimoimiseksi halutaan mahdollisimman lyhyellä seuranta-ajalla pystyä ennustamaan pelin tuottavuus pitkällä aikavälillä. Käytettävässä mallissa on kaksi selittävää muuttujaa: ostavien pelaajien määrä ja ostotahti.

Pelialan data on muodoltaan samantapaista kuin muutenkin elinaika-analyysin sovelluskohteissa käytetty data. Data koostuu pelaajista, joista osa on tehnyt ostotapahtuman ja osa ei. Pelaajat, jotka eivät ole ostotapahtumaa tehneet, ovat sensuroituja. Sensuroiduista pelaajista ei kuitenkaan eroteta ostotapahtuman myöhemmin tekeviä pelaajia ostotapahtuman tekemiselle immuuneista pelaajista. Pelaajien jako immuuneihin ja alttiisiin yksilöihin on siis latenti muuttuja. Osa alttiista pelaajista tekee useita ostotapahtumia, mutta tässä tutkielmassa keskitytään ensimmäisen ostotapahtuman tekevien pelaajien määrään eikä kaiken kaikkiaan tehtävien ostotapahtumien määrään.

Ennen varsinaisen käsiteltävän mallin johtamista käydään läpi sen muotoilemiseen tarvittavaa elinaika-analyysin ja Expectation Maximization -algoritmin teoriaa. Elinaika-analyysistä esitellään siihen liittyviä käsitteitä ja tärkeimpiä malleja. Tämän jälkeen käydään läpi vaillinaisen datan tilanteeseen soveltuvan optimointialgoritmin toiminta. Data on vaillinaista, kun siinä on latenteja muuttujia, joiden arvoja ei tiedetä. Latenteja muuttujia sisältävän mallin optimaalisen parametrivektorin selvittämiseen käytettävä optimointimenetelmä on nimeltään EM-algoritmi.

Kirjallisuuskatsauksen jälkeen johdetaan tässä tutkielmassa käytettävälle mallille EM-algoritmin käyttämiseen tarvittavat funktiot. Tämän lisäksi saatu algoritmi toteutetaan R-ohjelmointikielellä (R), jolla havainnollistetaan algoritmin konvergenssia. Mallin toimivuutta testataan oikeasta pelistä kerätyllä datalla sekä mallin oletuksien mukaisesti generoidulla datalla. Todellisenä datana käytetään Hipster sheep -nimisestä pelistä [1] kerättyä dataa. Lopuksi vielä tutkitaan generoidulla datalla datan ominaisuuksien vaikutus-

ta siihen, miten suurimman uskottavuuden estimaatti konvergoi todelliseen arvoonsa.

## 2 Elinaika-analyysi

Aloitetaan tämän tutkielman teoriaosuus perehtymällä elinaika-analyysiin. Käydään ensin läpi tärkeimpiä peruskäsitteitä ja ominaisuuksia. Jatketaan sitten käytetyimpien parametrusten ja epäparametrusten mallien esittelemiseen ja keskitytään lopuksi vielä elinaika-analyysin osa-alueen cure-mallien ja erityisesti mixture cure -mallin perusteisiin.

Elinaika-analyysi on tilastotieteen ala, jossa tarkastellaan yksilöitä, jotka suorittavat tarkastelun kohteena olevan tapahtuman jollakin ajanhetkellä. Tällaisesta tapahtumasta käytetään nimitystä *vikaantuminen*. Aikaa, joka kuluu ennen kuin tapahtuma tehdään, kutsutaan *vikaantumisajaksi*. Elinaika-analyysissä jokainen yksilö voi tehdä tarkastelun kohteena olevan tapahtuman enintään kerran. [2]

Tarkastelun kohteena olevan tapahtuman tyyppi riippuu alasta, johon elinaika-analyysiä sovelletaan. Tästä syystä tapahtuma voi kuvata hyvinkin erilaisia asioita. Data voi olla kerätty lääketieteellisestä tutkimuksesta, jolloin tapahtuma voi olla potilaan kuolema, tai jos kerätty data on teollisuudesta, niin tapahtuma voi olla esimerkiksi jonkin laitteen tietyn osan hajoaminen. [2]

Elinaika-analyysin menetelmät riippuvat *välttöjakaumasta*, joka voidaan määrittellä välttöfunktion ja riskifunktion avulla. Olkoon satunnaismuuttuja  $T$  vikaantumisaika. Merkitään, että sen tiheysfunktio on  $f(t)$  ja kertymäfunktio on  $F(t)$ . Tällöin voidaan määrittellä *välttöfunktio*

$$S(t) = 1 - F(t) = P(T > t), \quad 0 < t < \infty, \quad (1)$$

joka määrittää todennäköisyyden, että vikaantuminen tapahtuu vasta ajanhetken  $t$  jälkeen [3]. Todennäköisyysfunktiona välttöfunktio luonnollisestikin saa arvoja väliltä  $[0,1]$ . Välttöfunktion arvo ajanhetkellä  $t = 0$  on 1, sillä kertymäfunktion arvo  $F(0) = 0$ , ja välttöfunktio on vähenevä tai vakio ajanfunktiona. Välttöfunktio määritellään usein riskifunktion avulla. Riskifunktio

$$\begin{aligned} h(t) &= \frac{f(t)}{1 - F(t)} \\ &= \lim_{\Delta t \rightarrow 0} \frac{P(t < T < t + \Delta t \mid T > t)}{\Delta t} \end{aligned} \quad (2)$$

kuvaa todennäköisyyttä, että vikaantuminen tapahtuu seuraavan ajanjakson  $\Delta t$  aikana, kun yksilö on selvinnyt ajanhetkeen  $t$  asti. [4]

Integroimalla riskifunktio muodossa (2) saadaan *kumulatiivinen riskifunktio*

$$\begin{aligned} H(t) &= \int_0^t h(u) \, du = \int_0^t \frac{f(u)}{1-F(u)} \, du = - \int_0^t \log[1-F(u)] \, du \\ &= -\log[1-F(t)] = -\log S(t), \end{aligned} \quad (3)$$

joka on riskifunktion kertymäfunktio ja täten kuvaa riskifunktion alle jäävää pinta-alaa ajanhetkeen  $t$  asti tarkasteltuna. Kaavan (3) integraalia laskettaessa tulee muistaa, että tiheysfunktio  $f(t)$  on kertymäfunktion  $F(t)$  derivaatta. Kaavasta (3) ratkaisemalla  $S(t)$  saadaan, että välttöfunktio voidaan kirjoittaa muodossa

$$S(t) = e^{-\int_0^t h(u) \, du} = e^{-H(t)}. \quad (4)$$

Yhtälö (4) osoittaa riskifunktion yhteyden välttöfunktioon. Yhtälö (2) on mahdollista kirjoittaa myös muodossa, jossa käytetään välttöfunktioita  $S(t)$  riskifunktion  $h(t)$  määrittämiseen. Tällöin riskifunktio lasketaan yhtälön

$$h(t) = -\frac{dS(t)}{dt} \frac{1}{S(t)} \quad (5)$$

avulla. [4]

Vikaantumisaajan tarkkaan määrittämiseen tarvitaan kolme ehtoa: ajan nollakohdan tulee olla yksikäsitteinen, ajankulun mitta-asteikko pitää olla sovittuna ja vikaantumisen määritelmän tulee olla täysin selvä. Ajan nollakohta pitää olla tarkkaan määritelty jokaiselle yksilölle. Yleensä se ei ole kalenterin mukaisessa ajassa sama jokaiselle yksilölle. Esimerkiksi lääketieteellisessä tutkimuksessa jokaisella potilaalla on oma nollakohtansa, joka kuvaa esimerkiksi ajanhetkeä, jolloin potilaalla todetaan tarkasteltavana oleva sairaus. Tällaisessa tapauksessa vikaantumisaika tarkoittaa aikaa sairauden toteamisesta potilaan kuolemaan kyseisestä sairaudesta johtuen. [2]

Tästä huomataan myös vikaantumisen tarkkaan määrittämisen tärkeys. Edellisessä esimerkissä olisi voitu todeta, että vikaantumisaika on aika sairauden toteamisesta potilaan kuolemaan. Tällöin vikaantumista ei olisi tarkkaan määritelty, sillä potilas voi diagnoosin saamisen jälkeen kuolla muustakin syystä kuin tutkittavana olevasta sairaudesta johtuen. [2]

Ajan etenemisen mitta-asteikkona käytetään usein reaaliaikaa, mutta on olemassa muitakin asteikkoja ajan etenemisen mittaamiseen. Esimerkiksi teollisuudessa aikaa mitataan usein kumulatiivisena käyttönä. Mitta-asteikon tulee kuitenkin olla sellainen, että vikaantumisaika saa vain positiivisia arvoja. Sen pitää olla jokaiselle yksilölle sama myös siitä syystä, että kahden samanlaisen yksilön pitäisi olla samanlaisessa tilassa, kun kumpaakin on seurattu yhtä kauan. [2]

Vikaantumisajan tarkan määrittämisen lisäksi toinen olennainen piirre elinaika-analyysissä on *sensurointi*. Sensuroinnilla tarkoitetaan sitä, että ajan nollakohtaa  $t = 0$  tai tapahtumahetkeä  $T$  ei pystytä tarkkaan määrittämään; aloitustapahtuma tai vikaantuminen on siis tapahtunut havainnointiajan ulkopuolella. Sensurointi voi olla joko *informatiivista* tai *epäinformatiivista*. Sensurointi on informatiivista, kun seurannasta poistuminen johtuu vikaantumiseen liittyvistä syistä, ja epäinformatiivista, kun se johtuu muusta kuin vikaantumiseen liittyvästä syystä. [4]

Yleisin sensurointitapa on *oikealta sensurointi*, joka tarkoittaa, että tapahtumahetkeä  $T^*$  ei tiedetä tarkkaan, mutta sen tiedetään olevan myöhemmin kuin ajanhetki  $C$ . Ajanhetki  $C$  kuvaa hetkeä, jona sensurointi tapahtuu, eli hetkeä, jolloin yksilön seuraaminen lopetetaan. Tällöin havaittu aika on  $T = \min(T^*, C)$  ja muuttuja  $\delta = \mathbb{I}(T^* \leq C)$  on *sensurointi-indeksi*, joka kuvaa, onko sensurointi tapahtunut vai ei. Muuttujan  $\delta$  määritelmässä käytetyn *indikaattorifunktion*  $\mathbb{I}(A)$  arvo on 1, kun  $A$  on tosi, ja 0, kun se on epätosi. Vastaavasti on mahdollista, että ajan nollakohtaa  $t = 0$  ei tiedetä tarkasti. Tällöin kyseessä on *vasemmalta sensurointi*. [4]

Sensurointi voidaan myös jakaa kolmeen tyyppiin: tyyppin I, tyyppin II ja satunnainen sensurointi. Tyyppin I sensuroinnissa sensurointiaika  $C$  on päätetty etukäteen ja se on sama jokaiselle yksilölle. Tyyppin II sensuroinnissa sensurointiajankohtaa ei päätetä vaan yksilöitä seurataan siihen asti, että etukäteen päätetty prosentti on vikaantunut. Satunnaisessa sensuroinnissa jokaiselle yksilölle  $i$  on oma sensurointi aikansa  $c_i$ . Sensurointiajat ovat toisistaan riippumattomia ja samoin jakautuneita. Satunnaisen sensuroinnin tapauksessa havaitaan siis vikaantumisaikojen  $t_i^*$  sijaan parit  $(t_i, \delta_i)$ , joissa

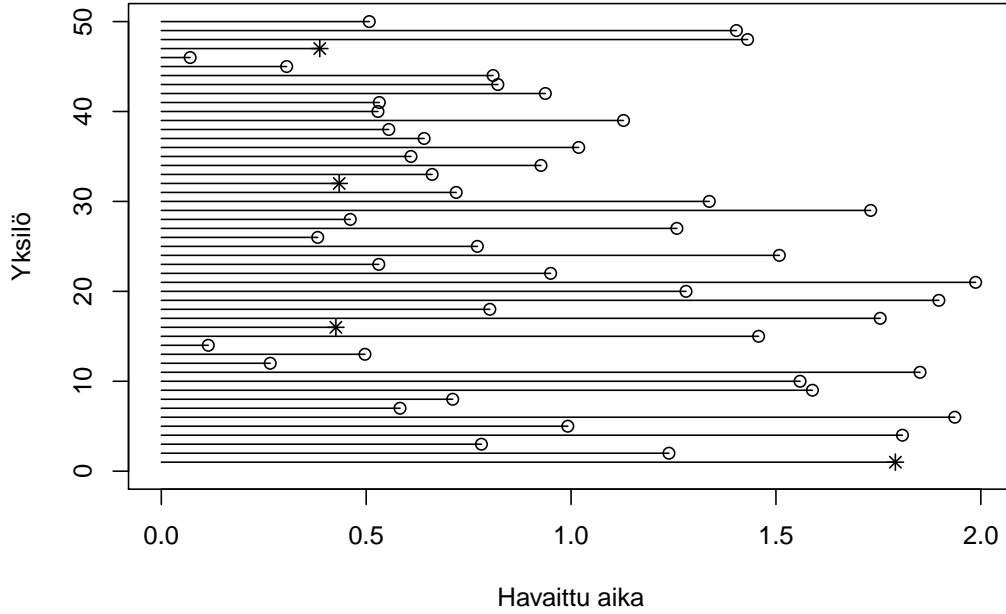
$$t_i = \min(t_i^*, c_i) \quad \text{ja} \quad (6)$$

$$\delta_i = \mathbb{I}(t_i^* \leq c_i). \quad (7)$$

[3]

Havainnollistetaan sensuroinnin vaikutusta kuvalla 1, josta nähdään, että osaa yksilöistä ei ole seurattu koko havainnointiaikaa, sillä ne ovat vikaantuneet jo aiemmin. Vikaantuneille yksilöille havaittu aika on tietysti lyhempi kuin sensurointiin asti seuratuille yksilöille

Elinaika-analyysissä havainnot noudattavat usein jotakin todennäköisyysjakaumaa, jonka muotoa kuvaillaan parametreilla  $\theta$ . Käytetyimmät jakaumat esitellään osiossa 2.1. *Uskottavuusfunktio*  $L(\theta)$  kuvaa jakauman parametrien



Kuva 1: Esimerkki sensuroidusta 50 yksilön datasta.

$\theta$  todennäköisyyttä ja elinaika-analyysissä se on muotoa

$$\begin{aligned}
 L(\theta) &= \prod_{i=1}^n L_i(\theta) \\
 &= \prod_{i=1}^n f(t_i; \theta, \mathbf{x}_i)^{\delta_i} S(t_i; \theta, \mathbf{x}_i)^{1-\delta_i},
 \end{aligned} \tag{8}$$

missä aika  $t_i$  on kaavan (6) mukainen havaittu aika, sensurointi-indeksi  $\delta_i$  on yhtälön (7) määräämä arvo ja muuttuja  $\mathbf{x}_i$  on muut havaitut arvot. Parametrien arvoja, jotka maksimoivat uskottavuusfunktion, kutsutaan *suurimman uskottavuuden estimaatiksi*. Suurimman uskottavuuden estimoinnissa määritettävät parametrien arvot suppevat todellisiin arvoihin [5]. Suurimman uskottavuuden estimaattia  $\hat{\theta}$  laskettaessa usein maksimoidaan logaritmista

uskottavuusfunktiota

$$\begin{aligned}
 l(\boldsymbol{\theta}) &= \log L(\boldsymbol{\theta}) = \log \left\{ \prod_{i=1}^n L_i(\boldsymbol{\theta}) \right\} \\
 &= \sum_{i=1}^n \log L_i(\boldsymbol{\theta}) \\
 &= \sum_{i=1}^n \{ \delta_i \log f(t_i; \boldsymbol{\theta}, \mathbf{x}_i) + (1 - \delta_i) \log S(t_i; \boldsymbol{\theta}, \mathbf{x}_i) \}
 \end{aligned} \tag{9}$$

uskottavuusfunktion (8) sijasta. [6]

## 2.1 Parametriset mallit

Vikaantumisajan  $T$  jakauma noudattaa usein jonkinlaista mallia. Tunnistettavien jakaumien pohjalta on kehitetty useita parametrisia malleja, joilla kuvataan eloonjäämisprosessia. Esitellään seuraavaksi kirjan [6] mukaisesti yleisimmin käytetyt parametriset mallit, jotka ovat eksponentti-, Weibull-, gamma-, log-normaali-, log-logistinen ja Gompertz-jakauma.

### 2.1.1 Eksponenttijakauma

Aloitetaan esittelemällä eksponenttijakauma, sillä se on parametrisista malleista yksinkertaisin. Eksponenttijakauman määrittämiseen riittää yksi parametri, jota merkitään kirjaimella  $\lambda$ . Siitä käytetään nimitystä *skaalausparametri* ja se on vakiovikaantumisnopeus. Eksponenttijakauman välttöfunktio

$$S(t; \lambda) = e^{-\lambda t}. \tag{10}$$

Välttöfunktion muodosta (10) ja riskifunktion kaavasta (5) saadaan johdettua, että riskifunktio

$$h(t; \lambda) = \lambda. \tag{11}$$

Kumulatiivinen riskifunktio  $H(t; \lambda)$  on tällöin kaavan (3) perusteella muotoa

$$H(t; \lambda) = \lambda t.$$

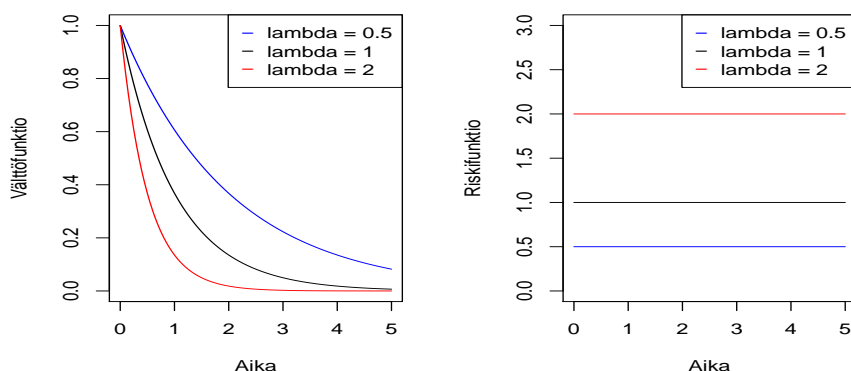
Eksponenttijakauman tiheysfunktio  $f(t)$  voidaan määritellä myös yhtälöiden (10) ja (11) avulla:

$$f(t; \lambda) = h(t; \lambda) S(t; \lambda) = \lambda e^{-\lambda t}. \tag{12}$$

Jakauman kertymäfunktio  $F(t)$  puolestaan saadaan laskettua yhtälöstä (1), jolloin saadaan, että

$$F(t; \lambda) = 1 - S(t; \lambda) = 1 - e^{-\lambda t}.$$

Kuvassa 2 havainnollistetaan parametrin  $\lambda$  vaikutusta välttöfunktion (10) muotoon. Huomataan, että välttöfunktio on sitä loivempi, mitä pienempi muuttujan  $\lambda$  arvo on. Riskifunktion (11) muotoon parametrin  $\lambda$  arvon muuttaminen ei tietenkään vaikuta riskifunktion ollessa vakio.



Kuva 2: Skaalausparametrin  $\lambda$  vaikutus eksponenttijakauman välttö- ja riskifunktioihin.

Eksponttijakauman välttöfunktion logaritmi on lineaarinen ajan suhteen

$$\log S(t; \lambda) = -\lambda t.$$

Tätä ominaisuutta voidaan hyödyntää kokeellisen datan kanssa, kun halutaan graafisesti selvittää, noudattaako se eksponenttijakaumaa.

Eksponttijakauman toinen tärkeä ominaisuus on muistittomuus. Jakauman muistittomuus tarkoittaa sitä, että

$$P(T > \Delta t) = P(T > t + \Delta t \mid T > t),$$

eli todennäköisyys selvitä vielä aika  $\Delta t$  ei muutu vaikka yksilön tiedettäisiin jo selvinneen aika  $t$ . Jokaisesta tarkasteluhetkestä katsottuna on siis sama todennäköisyys selvitä aika  $\Delta t$ . Elinaika-analyysissä eksponenttijakauman käyttö tarkoittaa, että yksilön jäljellä olevan elinajan odotusarvo on sama minkä tahansa ikäisenä. Elinaika-analyysin prosesseissa riskifunktio ei yleensä ole vakio, joten eksponenttijakauman käyttö elinaika-analyysissä on rajoitettua. Sitä kuitenkin käytetään tilanteissa, joissa tarkastellaan vikaantumisia lyhyellä aikavälillä.

### 2.1.2 Weibull-jakauma

Esitellään seuraavaksi Weibull-jakauma, joka on eksponenttijakauman tapaan yksinkertainen, mutta mukautuva. Weibull-jakauma määritellään kahdella parametrilla, joten se soveltuu useampaan tilanteeseen kuin eksponenttijakauma. Skaalausparametrin  $\lambda$  lisäksi Weibull-jakaumaa määrittää muoto parametri  $\alpha$ .

Vikaantumisajan  $T$  noudattaessa Weibull-jakaumaa (merkitään  $T \sim \text{Weib}(\lambda, \alpha)$ ) välttöfunktio on muotoa

$$S(t; \lambda, \alpha) = \exp[-(\lambda t)^\alpha]. \quad (13)$$

Vastaavasti kuin eksponenttijakauman tapauksessa, riskifunktio  $h(t)$  saadaan johdettua välttöfunktion muodosta (13) ja riskifunktion määritelmästä (5):

$$h(t; \lambda, \alpha) = \lambda \alpha (\lambda t)^{\alpha-1}. \quad (14)$$

Kumulatiivinen riskifunktio  $H(t; \lambda, \alpha)$  saadaan laskettua kaavalla (3). Weibull-jakauman tapauksessa

$$H(t; \lambda, \alpha) = (\lambda t)^\alpha. \quad (15)$$

Kuten eksponenttijakaumankin tapauksessa, myös Weibull-jakaumaa käytettäessä pystytään testaamaan, noudattaako data oletettua jakaumaa. Weibull-jakauman tapauksessa se tapahtuu ottamalla logaritmi kumulatiivisesta riskifunktiosta (15). Kumulatiivisen riskifunktion logaritmin  $\log H(t; \lambda, \alpha)$  arvot ajan  $t$  logaritmin funktiona muodostavat suoran, sillä

$$\log H(t; \lambda, \alpha) = \log \lambda + \alpha \log t.$$

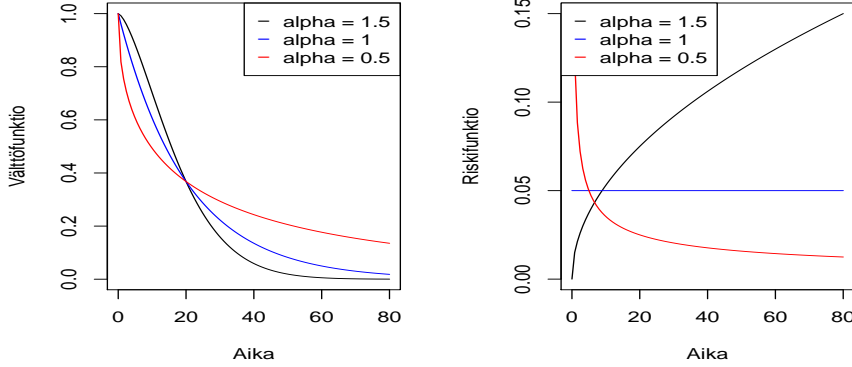
Jakauman tiheysfunktio saadaan jälleen välttöfunktion (13) ja riskifunktion (14) tulona

$$f(t; \lambda, \alpha) = h(t) S(t) = \lambda \alpha (\lambda t)^{\alpha-1} \exp[-(\lambda t)^\alpha]$$

ja kertymäfunktio saadaan välttöfunktion (13) avulla

$$F(t; \lambda, \alpha) = 1 - \exp[-(\lambda t)^\alpha].$$

Kuvassa 3 havainnollistetaan muoto parametrin  $\alpha$  vaikutusta välttö- ja riskifunktioiden muotoon. Huomataan, että eksponenttijakauma on Weibull-jakauman erikoistapaus, kun  $\alpha = 1$ . Parametrin  $\alpha$  arvoilla  $\alpha > 1$  riskifunktio on epälineaarisesti kasvava ajan  $t$  funktiona ja arvoilla  $\alpha < 1$  huomataan riskifunktion olevan epälineaarisesti vähenevä ajan  $t$  funktiona. Weibull-jakauman tapauksessa siis riskifunktio on joko monotonisesti kasvava tai vähenevä. Tämä onkin eräs Weibull-jakauman tärkeistä ominaisuuksista.



Kuva 3: Muotoparametrin  $\alpha$  vaikutus Weibull-jakauman välttö- ja riskifunktioihin.

### 2.1.3 Gammajakauma

Esitellään kolmantena parametriseena mallina gammajakauma. Sitä voidaan käyttää mallinnettaessa sellaista eloonjäämisprosessia, joka ei noudata mitään symmetristä jakaumaa. Gammajakauma määritellään kahden parametrin avulla samoin kuin Weibull-jakauma. Merkinnällä  $\Gamma(\alpha)$  tarkoitetaan gammafunktioita ja se määritellään kaavan (16) mukaisena integraalina

$$\Gamma(\alpha) = \int_0^{\infty} t^{\alpha-1} e^{-t} dt. \quad (16)$$

Gamma-jakauman tiheysfunktio on muotoa

$$f[t; T \sim \Gamma(\lambda, \alpha)] = \frac{\lambda (\lambda t)^{\alpha-1} e^{-\lambda t}}{\Gamma(\alpha)}, \quad t > 0. \quad (17)$$

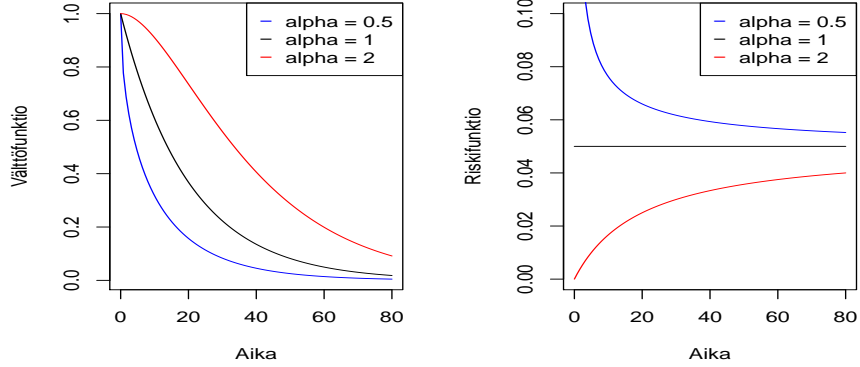
Jälleen huomataan, että jos  $\alpha = 1$ , niin gammajakauma redusoituu eksponenttijakaumaksi, joten eksponenttijakauma on myös gammajakauman erikoistapaus. Tämä nähdään myös kuvasta 4, jossa havainnollistetaan gammajakauman välttö- ja riskifunktioita erilaisilla muotoparametrin  $\alpha$  arvolla.

Keskitytään seuraavaksi erikoistapaukseen, jossa  $\lambda = 1$ . Tällöin tiheysfunktio (17) saadaan kirjoitettua yksiparametrisen gammajakauman tiheysfunktiona

$$f[t; T \sim \Gamma(\alpha)] = \frac{t^{\alpha-1} e^{-t}}{\Gamma(\alpha)}, \quad t > 0.$$

Yksiparametrisen gammajakauman kertymäfunktioita

$$F[t; T \sim \Gamma(\alpha)] = \int_0^t \frac{u^{\alpha-1} e^{-u}}{\Gamma(\alpha)} du$$



Kuva 4: Muotoparametrin  $\alpha$  vaikutus gammajakauman välttö- ja riskifunktioihin.

kutsutaan myös epätäydelliseksi gammafunktioksi. Gamma-jakauman tapauksessa välttöfunktio ja riskifunktio määritellään epätäydellisen gammafunktion avulla:

$$\begin{aligned}
 S[t; T \sim \Gamma(\alpha)] &= 1 - F[t; T \sim \Gamma(\alpha)] \\
 h[t; T \sim \Gamma(\alpha)] &= \frac{\lambda (\lambda t)^{\alpha-1} e^{-\lambda t}}{\Gamma(\alpha) \{1 - F[t; T \sim \Gamma(\alpha)]\}}.
 \end{aligned}$$

#### 2.1.4 Log-normaalijakauma

Myös log-normaalijakauma on yksi suosituista elinaika-analyysissä käytetyistä parametrisista malleista. Log-normaalijakauman käyttökelpoisuus perustuu siihen, että sen logaritmi on normaalijakautunut parametrein  $\mu$  ja  $\sigma^2$ . Eli jos elinaika-analyysissä vikaantumisaika  $T$  noudattaa log-normaalijakaumaa (merkitään  $T \sim LN(\mu, \sigma^2)$ ), niin vikaantumisaajan  $T$  logaritmi  $\log T$  noudattaa normaalijakaumaa, eli voidaan merkitä  $\log T \sim N(\mu, \sigma^2)$ .

Log-normaalijakauman tiheysfunktio on muotoa

$$f[t; T \sim LN(\mu, \sigma)] = \frac{1}{t\sigma\sqrt{2\pi}} \exp\left\{-\frac{(\log t - \mu)^2}{2\sigma^2}\right\}$$

ja sen kertymäfunktio voidaan lausua normaalijakauman kertymäfunktion  $\Phi$  avulla:

$$F[t; T \sim LN(\mu, \sigma)] = \Phi\left(\frac{\log t - \mu}{\sigma}\right) = \Phi\left(\frac{\log t}{\sigma}\right), \quad t > 0. \quad (18)$$

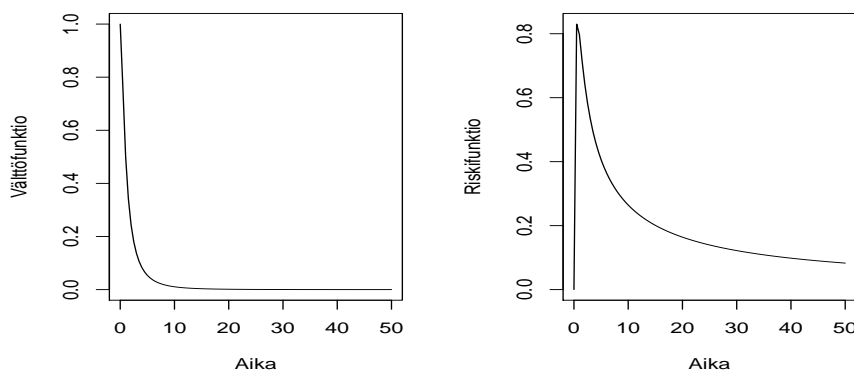
Jälleen saadaan esitettyä välttöfunktio kertymäfunktion (18) avulla:

$$S [t; T \sim LN (\mu, \sigma)] = 1 - \Phi \left( \frac{\log t}{\sigma} \right).$$

Myös riskifunktio ja kumulatiivinen riskifunktio voidaan kirjoittaa normaali-jakauman tiheys- ja kertymäfunktioiden  $\varphi$  ja  $\Phi$  avulla:

$$h [t; T \sim LN (\mu, \sigma)] = \frac{1}{t\sigma} \varphi \left( \frac{\log t}{\sigma} \right) \\ H [t; T \sim LN (\mu, \sigma)] = -\log \left[ 1 - \Phi \left( \frac{\log t}{\sigma} \right) \right].$$

Kuvassa 5 näytetään esimerkit log-normaali-jakauman välttö- ja riskifunktioiden muodoista parametrien arvoilla  $\mu = 0$  ja  $\sigma = 1$ .



Kuva 5: Havainnollistetaan log-normaali-jakauman välttö- ja riskifunktioiden muotoja parametrien arvoilla  $\mu = 0$  ja  $\sigma = 1$ .

### 2.1.5 Log-logistinen jakauma

Seuraavana on vuorossa log-logistinen jakauma, joka on muodoltaan melko samanlainen kuin log-normaali-jakauma. Log-logistisen jakauman ero log-normaali-jakaumaan on kuitenkin merkittävä. Log-normaali-jakauma sopii vain tilanteisiin, joissa ei ole tapahtunut sensurointia. Toinen log-logistisen jakauman tärkeä ominaisuus on se, että sen kertymä-, välttö- ja riskifunktiot pysytään esittämään suljetussa muodossa.

Log-logistista jakaumaa kuvataan kahdella parametrilla,  $\lambda$  ja  $\alpha$ . Vikaantumisaajan  $T$  noudattaessa log-logistista jakaumaa merkitään  $T \sim LLogist (\lambda, \alpha)$ .

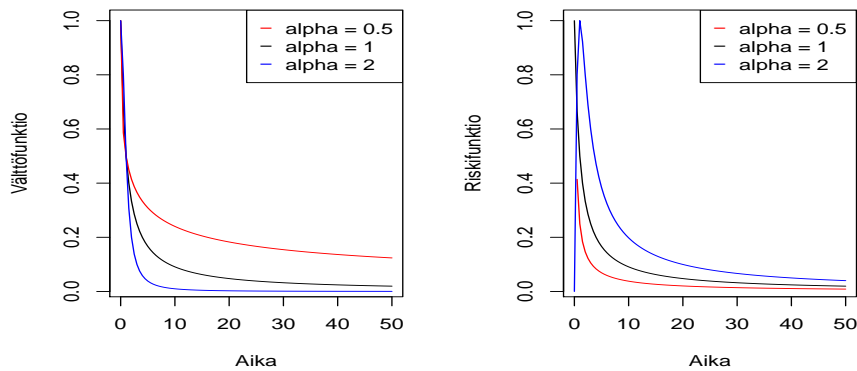
Allaolevassa listauksessa esitellään log-logistisen jakauman tiheys-, kertymä-, välttö- ja riskifunktiot, ja kuvassa 6 on esimerkit tämän jakauman välttö- ja riskifunktioista.

$$f [t; T \sim LLogist (\lambda, \alpha)] = \frac{(\alpha/\lambda) (t/\lambda)^{\alpha-1}}{[1 + (t/\lambda)^\alpha]^2}, \quad t > 0$$

$$F [t; T \sim LLogist (\lambda, \alpha)] = \frac{1}{1 + (t/\lambda)^{-\alpha}}$$

$$S [t; T \sim LLogist (\lambda, \alpha)] = \frac{1}{1 + (t/\lambda)^\alpha}$$

$$h [t; T \sim LLogist (\lambda, \alpha)] = \frac{(\alpha/\lambda) (t/\lambda)^{\alpha-1}}{1 + (t/\lambda)^\alpha}$$



Kuva 6: Esimerkki log-logistisen jakauman välttö- ja riskifunktioista.

### 2.1.6 Gompertz-jakauma

Viimeisenä parametrinenä mallina tässä tutkielmassa esitellään Gompertz-jakauma. Gompertzin laki on ehkä yleisimmin käytetty ihmisen kuolevuutta ja eloonjäämistä kuvaileva malli. Gompertz-jakaumassa riskifunktio on muotoa

$$h (t; h_0, r) = h_0 e^{rt}, \quad h_0, r > 0. \quad (19)$$

Parametria  $h_0$  kutsutaan usein iästä riippumattomaksi riskikertoimeksi ja parametri  $r$  puolestaan on iästä riippuva riskikerroin. Yhtälön (4) mukaisesti integroimalla riskifunktio (19) saadaan välttöfunktioiksi

$$S (t; h_0, r) = \exp \left[ - \int_0^t h (u; h_0, r) du \right] = \exp \left[ \frac{h_0}{r} (1 - e^{rt}) \right]. \quad (20)$$

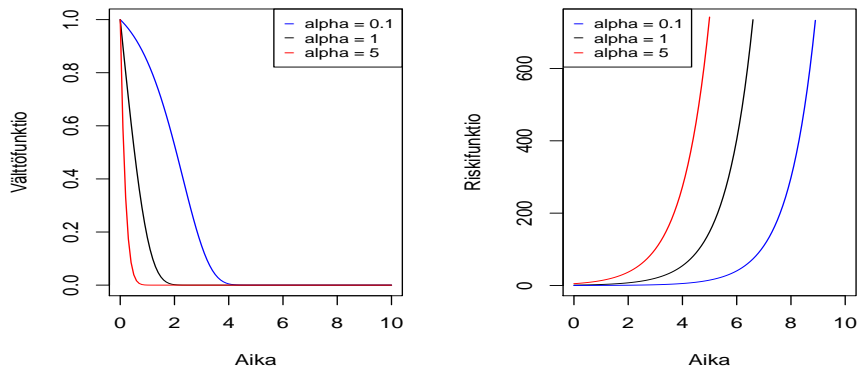
Gompertzin lain mukaisen jakauman tiheysfunktio saadaan funktioiden (19) ja (20) tulona

$$f(t; h_0, r) = h_0 e^{rt} \exp \left[ \frac{h_0}{r} (1 - e^{rt}) \right]$$

ja kertymäfunktio saadaan jälleen ratkaistua kaavasta (1), jolloin

$$F(t; h_0, r) = 1 - \exp \left[ \frac{h_0}{r} (1 - e^{rt}) \right].$$

Havainnollistetaan vielä tätäkin jakaumaa kuvaajilla. Kuvassa 7 esitetään Gompertz-jakauman välttö- ja riskifunktiot erilaisilla parametrin arvoilla.



Kuva 7: Esimerkit Gompertz-jakauman välttö- ja riskifunktioista.

## 2.2 Epäparametriset mallit

Käydään seuraavaksi lyhyesti läpi teoriaa epäparametrisista malleista. Epäparametrisissa malleissa ei nimensä mukaisesti ole jakaumaa kuvaavia parametreja, vaan estimointi perustuu suoraan havaitusta datasta laskettaviin arvoihin.

Parametrisia malleja käytettäessä täytyy pystyä määrittämään vikaantumisajan  $T$  jakauman muotoa kuvaava funktio. Epäparametristen mallien käyttö on joustavampaa, sillä niitä käytettäessä ei ole tällaisia rajoituksia. Epäparametrista mallia voidaan käyttää myös sovituksen hyvyuden graafisessa tarkastelussa. [2]

Yleisimmin käytetty epäparametrinen malli on Kaplan–Meier-malli, jota usein kutsutaan myös rajatuloestimaatiksi. Tässä osiossa esiteltävä teoria pohjautuu artikkeliin [7], mutta yhdenmukaisuuden vuoksi käytetään osion 2.1 merkintätapoja. Tärkeimpiä näistä ovat välttöfunktio  $S(t)$ , joka määritellään kuten yhtälössä (1), välttöfunktion rajatuloestimaatti  $\hat{S}(t)$  ja funktio  $n(t)$ , joka kuvaa ajanhetkeä  $t$  pidempään selviävien ja havaittavien yksilöiden määrää. Funktion  $n(t)$  arvoa laskettaessa ei oteta huomioon hetkellä  $t$  tapahtuvia vikaantumisia.

Rajatuloestimaatin laskeminen perustuu seuraaviin vaiheisiin:

1. Jaetaan aika sopivasti valittuihin aikaväleihin  $(0, u_1), (u_1, u_2), \dots, (u_{r-1}, u_r)$
2. Lasketaan jokaiselle aikavälille  $(u_{j-1}, u_j)$  osamäärä  $\hat{s}_j = \frac{S_j}{S_{j-1}}$ , joka kertoo hetkeä  $u_j$  pidempään selviävien yksilöiden osuuden hetkeä  $u_{j-1}$  pidempään selviävistä yksilöistä
3. Lasketaan aikavälin rajapistettä  $t$  pidemmälle selviävien yksilöiden välttöfunktion arvo tulona

$$S(t) = \prod_{u_j \leq t} s_j.$$

Vaiheessa 1 aikavälit  $(u_{j-1}, u_j)$  valitaan siten, että vaiheen 2 osamäärät saadaan laskettua kaavalla

$$\hat{s}_j = \frac{n_j - d_j}{n_j} = \frac{n'_j}{n_j}. \quad (21)$$

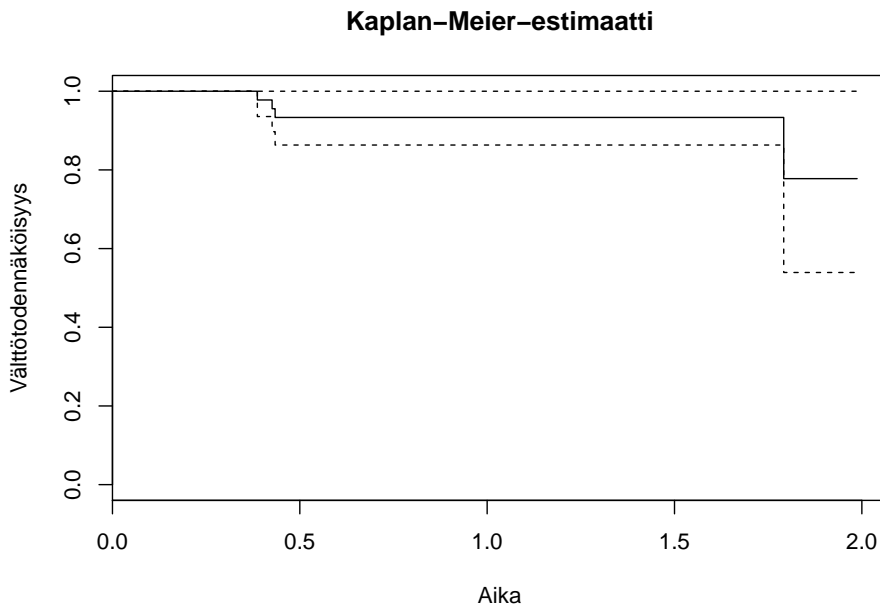
Kaavassa (21) muuttuja  $n_j$  on heti hetken  $u_{j-1}$  jälkeen havaittujen yksilöiden määrä, muuttuja  $d_j$  on aikavälillä  $(u_{j-1}, u_j)$  vikaantuvien yksilöiden määrä ja muuttuja  $n'_j$  on havaittujen yksilöiden määrä aikavälillä tapahtuvien vikaantumisten jälkeen. Vaiheessa 3 ei ole väliä sensurointien jakautumisesta

aikavälillä  $(u_{j-1}, u_j)$  kunhan kaavassa (21) esitetyt  $n_j$  ja  $n'_j$  on laskettu oikein. Vaiheessa 3 oletetaan myös, että sensurointi tapahtuu joko ennen aikavälin ensimmäistä vikaantumista tai vasta viimeisen vikaantumisen jälkeen ennen aikavälin päättymistä.

Määrä  $n_j$  on laskettu oikein, kun siihen ei ole laskettu mukaan yksilöitä, jotka vikaantuvat ajanhetkellä  $u_j$ . Siispä jokaisella ajanhetkellä  $u_j$  pätee epäyhtälöt  $n'_j \leq n(u_j) \leq n_{j+1}$ . Tämä tarkoittaa, että ajanhetkeä  $u_j$  pidempään selviävien ja havaittavien yksilöiden määrä  $n(u_j)$  on vähintään yhtäsuuri kuin aikavälillä  $(u_{j-1}, u_j)$  tapahtuvien vikaantumisten jälkeen havaittava yksilöiden määrä  $n'_j$  ja enintään yhtäsuuri kuin heti ajanhetken  $u_{j+1}$  jälkeen havaittujen yksilöiden määrä  $n_{j+1}$ . Tällöin rajatuloestimaatti on muotoa

$$\hat{S}(t) = \prod_{j=1}^k \frac{n'_j}{n_j},$$

missä  $u_k = t$  ja  $n'_j$  on määritelty samoin kuin kaavassa (21). Kaplan–Meierin rajatuloestimaattia laskettaessa data on oikealta sensuroitua ja vasemmalta katkaistua, eli yksilöistä, jotka vikaantuvat ennen havaintoajan alkamista, ei tiedetä mitään.



Kuva 8: Esimerkki Kaplan–Meier-estimaatista kuvan 1 datalle.

## 2.3 Mixture cure -malli

Esitellään vielä tämän luvun lopuksi cure-malli, minkä jälkeen pystytään määrittelemään sovellusosiossa käytettävä mixture cure -malli. Cure-malli on elinaika-analyysimalli, jota käytetään tilanteissa, joissa osa datan yksilöistä on immuuneja tarkastelun kohteena olevan tapahtuman realisoitumiselle. Tällaiset yksilöt eivät siis vikaannu, vaikka niitä seurattaisiin äärettömän pitkä aika. Matemaattisesti esitettynä tilanne, jossa cure-mallia voidaan käyttää, on sellainen, että välttöfunktion (1) raja-arvo on nolasta eroava eli

$$\lim_{t \rightarrow \infty} S(t) \neq 0.$$

Datan sopivuutta cure-mallin sovelluskohteeksi voidaan arvioida tarkastelemalla esimerkiksi osiossa 2.2 esitetyn Kaplan–Meierin rajatuloestimaatin  $\hat{S}(t)$  raja-arvoa. Kuvan 8 perusteella kuvan 1 data vaikuttaa olevan esimerkiksi tällaisesta datasta. [8]

Cure-mallin kohdalla osiossa 2 selitetty oikealta sensurointi voi tapahtua kahdesta eri syystä: yksilö on immuuni, eikä siis missään vaiheessa vikaannu, tai yksilö on altis, mutta sitä ei ole vielä seurattu riittävän pitkään. Joissain tilanteissa on mahdollista määrittää aika, johon mennessä kaikki alttiit yksilöt vikaantuvat. Aina tämä ei kuitenkaan ole mahdollista ja silloin immuunien ja alttiiden yksilöiden toisistaan erottaminen on mahdotonta. Cure-mallin kohdalla on kahdenlaisia yksilöitä, joten on luonnollista ajatella sitä kahden ryhmän sekoitemallina. Tällaisesta mallista käytetään nimitystä mixture cure -malli. [8]

Merkitään tässäkin osiossa vikaantumisaikaa muuttujalla  $T$ , ja määritellään, että muuttuja  $\zeta$  kuvaa yksilön immuuniutta. Yksilön ollessa immuuni tarkastelun kohteena olevan tapahtuman realisoitumiselle, alttiusindeksi  $\zeta = 0$ , ja muussa tapauksessa  $\zeta = 1$ . Kummallekin ryhmälle voidaan määrittellä oma välttöfunktio, joka on ehdollinen todennäköisyys, että yksilö vikaantuu vasta tarkasteluhetkeä myöhemmin. Nämä ovat muotoa

$$\begin{aligned} S_u(t) &= P(T > t | \zeta = 1) \quad \text{ja} \\ S_c(t) &= P(T > t | \zeta = 0) \equiv 1. \end{aligned}$$

Erottelussa käytetyt alaindeksit  $u$  ja  $c$  tulevat sanoista ”uncured” ja ”cured” vastaavassa järjestyksessä. Immuunien yksilöiden populaation välttöfunktio  $S_c(t) \equiv 1$ , sillä siihen kuuluvat yksilöt eivät missään vaiheessa vikaannu, joten todennäköisyys  $P(T > t) = 1$  millä tahansa  $t$ . Koko populaation marginaalinen välttöfunktio on muotoa

$$\begin{aligned} S(t) &= \pi S_u(t) + (1 - \pi) S_c(t) \\ &= \pi S_u(t) + 1 - \pi, \end{aligned} \tag{22}$$

kun yksilön vikaantumistodennäköisyys  $P(\zeta = 1) = \pi$ . [8]

Mixture cure -malli on siis parametrinen malli, jossa vikaantumisaika  $T$  noudattaa jotakin osiossa 2.1 esiteltyä jakaumaa. Tämän lisäksi mallia selittää parametri  $\pi$ , joka on alttiiden yksilöiden osuus koko populaatiosta. Tämän mallin suurimman uskottavuuden estimaatin ratkaiseminen on ongelmallista tavallisilla optimointimenetelmillä, sillä malli sisältää latentteja muuttujia.

## 3 Expectation Maximization -algoritmi

Esitellään seuraavaksi osion 2.3 mallin parametrien optimiarvojen ratkaisemiseen sopiva optimointimenetelmä. Artikkelissa [9] Dempster, Laird ja Rubin käyttivät siitä nimeä Expectation Maximization -algoritmi (EM-algoritmi). Esitellään ensin algoritmiin liittyviä käsitteitä ja sen toimintaan liittyvä teoria pääasiassa kirjaan [5] pohjautuen. Tämän jälkeen perehdytään algoritmin konvergenssin todistamiseen pääasiassa artikkelin [10] mukaisesti. Yhdenmukaisuuden vuoksi koko luvussa käytetään kirjassa [5] käytettyjä merkintöjä.

### 3.1 Algoritmin toiminta

EM-algoritmi on iteroiva algoritmi uskottavuusfunktion suurimman arvon selvittämiseksi. Sitä käytetään usein tilanteissa, joissa on esimerkiksi puuttuvaa dataa ja muiden optimointialgoritmien käyttö olisi tästä syystä hankalaa. EM-algoritmia voidaan myös käyttää tilanteissa, joissa datan epätäydellisyys ei ole selvää. Esimerkki tällaisesta tilanteesta on tapaus, jossa on latentteja muuttujia.

EM-algoritmia käytettäessä on havaittu data  $\mathbf{y}$  ja tiedetään, että on dataa, jota ei syystä tai toisesta ole voitu havaita. Puuttuvaa dataa kuvataan vektorilla  $\mathbf{z}$ . Datasta käytetään nimitystä *täydellinen data*, kun puuttuvaa dataa ei ole, ja sitä merkitään muuttujalla  $\mathbf{x}$ . Osa taulukkoon 1 kerätyistä funktioista on taulukossa vain havaitulle datalle  $\mathbf{y}$ , mutta vastaavat funktiot ovat olemassa myös täydelliselle datalle  $\mathbf{x}$ . Täydellistä dataa käytettäessä funktiot tunnistaa alaindeksistä  $c$ .

Algoritmin nimi kuvaa suoraan sen toimintaa, sillä se muodostuu kahdesta vaiheesta, joita kutsutaan nimillä E-askel (expectation) ja M-askel (maximization). Näitä vaiheita toistetaan vuorotellen, kunnes saavutetaan optimiarvo. Ensimmäisessä vaiheessa (E-askel) puuttuvan datan arvoja  $\mathbf{z}$  arvioidaan laskemalla niille ehdolliset odotusarvot tunnettujen arvojen  $\mathbf{y}$  avulla. Toisessa vaiheessa (M-askel) selvitetään uskottavuusfunktion  $L(\Psi)$  maksimoiva parametrivektori  $\hat{\Psi}$ , suurimman uskottavuuden estimaatti, E-askeleessa laskettuja arvioita hyödyntämällä. Ensin siis lasketaan arviot puuttuvalle datalle ja sen jälkeen niitä hyödynnetään täydellisen datan uskottavuusfunktion maksimoinnissa. Tämän jälkeen saatua suurimman uskottavuuden estimaattia käyttämällä lasketaan uudet arviot puuttuvalle datalle.

Täydellisen datan uskottavuusfunktioita  $L_c(\Psi)$  kuvaa täydellistä dataa  $\mathbf{x}$  vastaavan satunnaismuuttujan  $X$  todennäköisyysjakauman tiheysfunktio  $g_c(\mathbf{x}; \Psi)$ , jossa vektori  $\Psi = (\Psi_1, \dots, \Psi_d)^T$  koostuu jakauman parametrien arvoista. Laskennallisista syistä käytetään usein logaritmista uskottavuusfunktioita  $l_c(\Psi) = \log L_c(\Psi) = \log g_c(\mathbf{x}; \Psi)$ . Logaritmfunktio on monotoninen,

joten logaritminen uskottavuusfunktio saa suurimman arvonsa samoilla parametrien  $\Psi$  arvoilla kuin uskottavuusfunktio. Havaitun datan  $\mathbf{y}$  tiheysfunktio  $g(\mathbf{y}; \Psi)$  on muotoa  $g(\mathbf{y}; \Psi) = \prod_{j=1}^n f(\mathbf{w}_j; \Psi)$ , kun havaittu data  $\mathbf{y} = (\mathbf{w}_1^T, \dots, \mathbf{w}_n^T)^T$  ja  $n$ -kokoisen satunnaisotannan havainnot noudattavat satunnaisvektorin  $W$  todennäköisyysjakaumaa  $f(\mathbf{w}; \Psi)$ .

Taulukko 1: EM-algoritmin ominaisuuksista kerrottaessa käytetyt merkinnät.

$\mathbf{x}$	Täydellistä dataa kuvaava vektori
$\mathbf{y}$	Havaittua dataa kuvaava vektori
$\mathbf{z}$	Puuttuvaa dataa kuvaava vektori
$\mathbf{X}$	Täydellistä dataa vastaava satunnaismuuttuja
$\mathbf{Y}$	Havaittua dataa vastaava satunnaismuuttuja
$g(\mathbf{y}; \Psi)$	Satunnaismuuttujan $\mathbf{Y}$ todennäköisyysjakauman tiheysfunktio
$\Psi = (\Psi_1, \dots, \Psi_d)^T$	Satunnaismuuttujan $\mathbf{Y}$ tiheysfunktion parametrivektori
$\Omega$	Parametrivektorin määrittelyjoukko
$L(\Psi) = g(\mathbf{y}; \Psi)$	Uskottavuusfunktio
$l(\Psi)$	Logaritminen uskottavuusfunktio
$k(\mathbf{x} \mathbf{y}; \Psi)$	Satunnaismuuttujan $\mathbf{X}$ ehdollinen tiheysfunktio
$\Psi^{(k)}$	Parametrivektori $k$ iteraatiokierroksen jälkeen
$Q(\Psi; \Psi^{(k)})$	Täydellisen datan uskottavuusfunktion ehdollinen odotusarvo
$\hat{\Psi}$	Parametrivektorin $\Psi$ suurimman uskottavuuden estimaatti

Varsinaisesti EM-algoritmia käytettäessä on siis kaksi otosavaruutta  $\mathcal{X}$  ja  $\mathcal{Y}$  sekä joukosta pisteeseen -kuvaus otosavaruudesta  $\mathcal{X}$  otosavaruuteen  $\mathcal{Y}$ . Pystytään havaitsemaan vain epätäydellisen datan vektori  $\mathbf{y} = \mathbf{y}(\mathbf{x})$  joukossa  $\mathcal{Y}$  sen sijasta, että havaittaisiin täydellisen datan vektori  $\mathbf{x}$  avaruudessa  $\mathcal{X}$ . Havaitun datan jakauman tiheysfunktio voidaan tällöin kirjoittaa muodossa

$$g(\mathbf{y}; \Psi) = \int_{\mathcal{X}(\mathbf{y})} g_c(\mathbf{x}; \Psi) d\mathbf{x},$$

missä  $\mathcal{X}(\mathbf{y})$  on otosvaruuden  $\mathcal{X}$  osajoukko, jonka määrää yhtälö  $\mathbf{y} = \mathbf{y}(\mathbf{x})$ .

EM-algoritmi pyrkii selvittämään epätäydellisen datan logaritmissen uskottavuusfunktion  $\log L(\Psi)$  maksimin  $\hat{\Psi}$  iteratiivisesti hyödyntämällä täydellisen datan logaritmissa uskottavuusfunktiota  $\log L_c(\Psi)$ . Sitä ei kuitenkaan tiedetä, joten käytetään sen ehdollista odotusarvoa, kun tiedetään havaittu data  $\mathbf{y}$  ja parametrien senhetkiset arvot  $\Psi^{(k)}$ . Tällöin EM-algoritmin vaiheet ovat seuraavaa muotoa:

**E-askel:** Lasketaan täydellisen datan logaritmissen uskottavuusfunktion ehdollinen odotusarvo

$$Q(\Psi; \Psi^{(k)}) = E_{\Psi^{(k)}} \{\log L_c(\Psi) | \mathbf{y}\}. \quad (23)$$

**M-askel:** Etsitään  $\Psi^{(k+1)} \in \Omega$ , joka maksimoi funktion  $Q(\Psi; \Psi^{(k)})$ , eli jolle

$$Q(\Psi^{(k+1)}; \Psi^{(k)}) \geq Q(\Psi; \Psi^{(k)}) \quad (24)$$

kaikilla  $\Psi \in \Omega$ .

Yleistetyssä EM-algoritmissa (GEM-algoritmi) maksimointivaiheessa valitaan parametrivektori  $\Psi^{(k+1)}$ , jolla funktion  $Q$  arvo on aiemman parametrivektorin  $\Psi^{(k)}$  tuottamaa arvoa suurempi, mutta ei edellytetä sen maksimoivan funktiota  $Q$ . Eli M-askel on muotoa: etsitään  $\Psi^{(k+1)} \in \Omega$ , joka toteuttaa yhtälön

$$Q(\Psi^{(k+1)}; \Psi^{(k)}) \geq Q(\Psi^{(k)}; \Psi^{(k)}).$$

EM-algoritmeilla on monia hyviä ominaisuuksia verrattuna muihin algoritmeihin. Osa näistä ominaisuuksista osoitetaan osiossa 3.2. Mainittakoon ensimmäisenä ominaisuutena, että EM-algoritmi on parantava algoritmi, sillä uskottavuusfunktion arvo kasvaa jokaisella iteraatiokierroksella (ks. lause 3). Toinen tärkeä ominaisuus on, että algoritmin konvergenssi pystytään osoittamaan melko yleisten ehtojen ollessa voimassa. Osiossa 3.2 osoitetaan algoritmin konvergenssi lokaaliin optimiin.

EM-algoritmi on usein helppo implementoida, sillä siinä käytetään täydellistä dataa ja laskuissa ei tarvitse laskea uskottavuusfunktion arvoa eikä sen derivaattoja. EM-algoritmin suorittamiseen ei tarvita paljoa muistia, mikä on hyvä asia. EM-algoritmia käytettäessä iteraatiokierrosten määrä voi olla suuri, mutta yhden kierroksen suorittaminen on melko nopeaa, joten algoritmi ei suuresta iteraatiokierrosten määrästä huolimatta ole muita optimointimenetelmiä hitaampi.

EM-algoritmia käytettäessä tarvittavien funktioiden johtaminen on usein helpompi kuin muita menetelmiä käytettäessä, sillä EM-algoritmissa mak-

simoidaan vain täydellisen datan logaritmisin uskottavuusfunktion ehdollinen odotusarvo. E-askeleessakaan se ei yleensä ole kovinkaan monimutkasta. EM-algoritmin konvergenssia on helppo tarkkailla seuraamalla uskottavuusfunktion arvon kasvua. Tällä menetelmällä pystytään myös laskemaan jonkinlaiset uskottavat arviot puuttuvalle datalle.

Esitellään osion lopuksi vielä määritelmässä 1 toinen tapa merkitä EM-algoritmin M-askeleessa olevaa epäyhtälöä (24) ja käytetään sitä suljetun kuvauksen määritelmässä 2.

**Määritelmä 1.** Funktion  $Q(\Psi; \Psi^{(k)})$  maksimoivien pisteiden joukko on

$$M(\Psi^{(k)}) = \arg \max_{\Psi \in \Omega} Q(\Psi; \Psi^{(k)}).$$

**Määritelmä 2.** Pisteeltä joukolle -kuvaus  $M(\Psi^{(k)})$  on *suljettu* pisteessä  $\Psi = \Psi_0$ , jos ehdoista

$$\begin{aligned} \Psi^{(k)} &\rightarrow \Psi_0, & \Psi^{(k)} &\in \Omega \\ \phi^{(k)} &\rightarrow \phi_0, & \phi^{(k)} &\in M(\Psi^{(k)}) \end{aligned}$$

seuraa, että  $\phi_0 \in M(\Psi_0)$ . Kuvaus on *suljettu joukossa*  $\Omega$ , jos se on suljettu jokaisessa pisteessä  $\Psi \in \Omega$ . [11]

## 3.2 Algoritmin konvergenssi

Nyt kun EM-algoritmin toimintaperiaate ja siihen liittyvät käsitteet on käyty läpi, voidaan keskittyä algoritmin konvergenssin osoittamiseen. Tässä osiossa esitettävät lauseet ja niiden todistukset ovat pääasiassa lähteestä [10]. Osoitetaan, että sekä uskottavuusfunktion  $L(\Psi)$  että parametrivektorin  $\Psi$  arvot konvergoivat. Aloitetaan uskottavuusfunktion konvergenssista, jonka avulla sitten osoitetaan parametrien konvergenssi. Osion lopussa käydään vielä läpi EM-algoritmin konvergenssiasteen laskeminen.

### 3.2.1 Uskottavuusfunktion arvo suppenee

Aloitetaan uskottavuusfunktion konvergenssin osoittaminen todistamalla, että millä tahansa EM-algoritmin iteraatiokierroksella saatava uskottavuusfunktion arvo  $L(\Psi^{(k)})$  on vähintään yhtäsuuri kuin edellisellä iteraatiokierroksella saatu uskottavuusfunktion arvo  $L(\Psi^{(k-1)})$ . EM-algoritmin monotonisuus on alun perin esitetty artikkelissa [9], mutta lauseen 3 todistuksessa käytetään kirjassa [5] esitetyn todistuksen mukaisia merkintöjä.

**Lause 3.** (EM-algoritmin monotonisuus) Uskottavuusfunktion arvo ei pienene, kun käytetään EM-algoritmia uuden parametrivektorin  $\Psi$  laskemiseen. Eli jokaisella  $k = 0, 1, 2, \dots$  pätee

$$L(\Psi^{(k+1)}) \geq L(\Psi^{(k)}). \quad (25)$$

*Todistus.* Määritellään ehdollisen tiheysfunktion  $k(\mathbf{x}|\mathbf{y}; \Psi)$  olevan muotoa

$$k(\mathbf{x}|\mathbf{y}; \Psi) = \frac{g_c(\mathbf{x}; \Psi)}{g(\mathbf{y}; \Psi)}, \quad (26)$$

kun tiedetään havaittu data  $\mathbf{y}$  ja tämänhetkiset parametrien arvot  $\Psi$ . Havaitun datan uskottavuusfunktio saadaan muotoon

$$\begin{aligned} \log L(\Psi) &= \log g(\mathbf{y}; \Psi) \\ &= \log g_c(\mathbf{x}; \Psi) - \log k(\mathbf{x}|\mathbf{y}; \Psi) \\ &= \log L_c(\Psi) - \log k(\mathbf{x}|\mathbf{y}; \Psi) \end{aligned} \quad (27)$$

ottamalla yhtälöstä (26) puolittain logaritmi ja sijoittamalla saatu lauseke logaritmissen uskottavuusfunktion taulukossa 1 esitettyyn määritelmään. Ottamalla odotusarvot yhtälön (27) kummaltakin puolelta, saadaan uskottavuusfunktio muotoon

$$\begin{aligned} \log L(\Psi) &= E_{\Psi^{(k)}} \{\log L_c(\Psi) | \mathbf{y}\} - E_{\Psi^{(k)}} \{\log k(\mathbf{X}|\mathbf{y}; \Psi) | \mathbf{y}\} \\ &= Q(\Psi; \Psi^{(k)}) - \mathcal{H}(\Psi; \Psi^{(k)}), \end{aligned} \quad (28)$$

missä funktiolla  $\mathcal{H}$  on merkitty täydellisen datan ehdollisen tiheysfunktion logaritmin odotusarvoa ehdoilla  $\mathbf{y}$  ja  $\Psi^{(k)}$ . Nyt yhtälön (28) perusteella

$$\begin{aligned} \log L(\Psi^{(k+1)}) - \log L(\Psi^{(k)}) &= \{Q(\Psi^{(k+1)}; \Psi^{(k)}) - Q(\Psi^{(k)}; \Psi^{(k)})\} \\ &\quad - \{\mathcal{H}(\Psi^{(k+1)}; \Psi^{(k)}) - \mathcal{H}(\Psi^{(k)}; \Psi^{(k)})\}. \end{aligned} \quad (29)$$

Yhtälön (29) oikeanpuolen ensimmäisen termin tiedetään olevan ei-negatiivinen kaavan (24) perusteella, joten epäyhtälö (25) toteutuu, kun

$$\mathcal{H}(\Psi^{(k+1)}; \Psi^{(k)}) - \mathcal{H}(\Psi^{(k)}; \Psi^{(k)}) \leq 0. \quad (30)$$

Epäyhtälö (30) on tosi, sillä jokaisella  $\Psi$

$$\begin{aligned}
& \mathcal{H}(\Psi; \Psi^{(k)}) - \mathcal{H}(\Psi^{(k)}; \Psi^{(k)}) \\
&= E_{\Psi^{(k)}} \{ \log k(\mathbf{X}|\mathbf{y}; \Psi) | \mathbf{y} \} - E_{\Psi^{(k)}} \{ \log k(\mathbf{X}|\mathbf{y}; \Psi^{(k)}) | \mathbf{y} \} \\
&= E_{\Psi^{(k)}} \{ \log k(\mathbf{X}|\mathbf{y}; \Psi) - \log k(\mathbf{X}|\mathbf{y}; \Psi^{(k)}) | \mathbf{y} \} \\
&= E_{\Psi^{(k)}} \left\{ \log \left[ \frac{k(\mathbf{X}|\mathbf{y}; \Psi)}{k(\mathbf{X}|\mathbf{y}; \Psi^{(k)})} \right] | \mathbf{y} \right\} \\
&\leq \log \left[ E_{\Psi^{(k)}} \left\{ \frac{k(\mathbf{X}|\mathbf{y}; \Psi)}{k(\mathbf{X}|\mathbf{y}; \Psi^{(k)})} \right\} | \mathbf{y} \right] \\
&= \log \int_{\mathcal{X}(\mathbf{y})} k(\mathbf{x}|\mathbf{y}; \Psi^{(k)}) d\mathbf{x} \\
&= 0.
\end{aligned}$$

Siispä uskottavuusfunktion arvo uusilla parametrien arvoilla  $\Psi^{(k+1)}$  on vähintään yhtäsuuri kuin edellisellä iteraatiokierroksella lasketuilla parametrien arvoilla  $\Psi^{(k)}$  ja täten uskottavuusfunktion arvojen jono  $\{L(\Psi^{(k)})\}$  on kasvava.  $\square$

EM-algoritmin monotonisuuden lisäksi vaaditaan, että uskottavuusfunktion arvojen jono  $\{L(\Psi^{(k)})\}$  on rajoitettu, jotta sen tiedetään suppenevan arvoon  $L^*$ . Arvo  $L^*$  on uskottavuusfunktion ääriarvokohta, mikäli  $L^* = L(\Psi^*)$  kriittisessä pisteessä  $\Psi^*$ . Piste  $\Psi^*$  on *kriittinen piste*, jos se on yhtälön  $\frac{\partial L(\Psi)}{\partial \Psi} = \mathbf{0}$  ratkaisu.

Oletuksista **OL1**, **OL2** ja **OL3** seuraa suoraan, että uskottavuusfunktion arvojen jono  $\{L(\Psi^{(k)})\}$  on ylhäältä rajoitettu millä tahansa parametrivektorin alkuarvauksella  $\Psi_0 \in \Omega$ .

**OL1** Parametrien määrittelyjoukko  $\Omega$  on  $r$ -ulotteisen euklidisen avaruuden  $\mathbb{R}^r$  osajoukko.

**OL2** Parametrivektorin alkuarvauksen määrittelemä avaruuden  $\Omega$  osajoukko  $\Omega_{\Psi_0} = \{\Psi \in \Omega : L(\Psi) \geq L(\Psi_0)\}$  on kompakti millä tahansa  $L(\Psi_0) > -\infty$ .

**OL3** Funktio  $L$  on jatkuva määrittelyjoukossa  $\Omega$  ja differentioituva määrittelyjoukon  $\Omega$  sisäosassa.

Nyt tiedetään oletukset, joiden tulee olla voimassa, jotta uskottavuusfunktion arvojen jono suppenee. Lemmassa 4 esitetään lisää oletuksia, joiden avulla uskottavuusfunktion arvojen jonon  $\{L(\Psi^{(k)})\}$  kriittiseen pisteeseen

suppeneminen EM-algoritmin tuottamissa pisteissä  $\Psi^{(k)}$  todistetaan lauseessa 6. Merkitään lokaalien maksimien joukkoa kirjaimella  $\mathcal{M}$  ja kriittisten pisteiden joukkoa kirjaimella  $\mathcal{S}$ .

**Lemma 4.** Olkoon  $\{\Psi^{(k)}\}$  jono GEM-algoritmillä laskettuja arvoja  $\Psi^{(k+1)} \in \mathbf{M}(\Psi^{(k)})$ . Oletetaan, että kuvaus  $\mathbf{M}(\Psi^{(k)})$  on suljettu joukon  $\Omega$  kriittisten pisteiden joukon  $\mathcal{S}$  komplementissa  $\overline{\mathcal{S}}$ , ja että

$$L(\Psi^{(k+1)}) > L(\Psi^{(k)}), \text{ kaikilla } \Psi^{(k)} \notin \mathcal{S}.$$

Tällöin jono  $\{\Psi^{(k)}\}$  suppenee aina kriittiseen pisteeseen  $\Psi^* \in \mathcal{S}$  ja  $L(\Psi^{(k)})$  suppenee monotonisesti pisteeseen  $L^* = L(\Psi^*)$ .

Lemman 5 tulos on riittävä ehto sille, että EM-algoritmissa käytetty kuvaus  $\mathbf{M}$  on suljettu. Lauseen 6 todistuksessa käytetään tätä ehtoa ja lemmaa 4, kun osoitetaan, että uskottavuusfunktion arvo konvergoi ääriarvokohtaan  $L(\Psi^*)$  jollakin kriittisellä pisteellä  $\Psi^*$ . Lauseen 6 todistuksessa käytetty merkintä  $D^{1,0}$  tarkoittaa ensimmäistä derivaattaa ensimmäisen parametrivektorin suhteen.

**Lemma 5.** (Jatkuvuusehto) EM-algoritmin kuvaus  $\mathbf{M}$  on suljettu, jos funktio  $Q(\Psi'; \Psi)$  on jatkuva sekä pisteessä  $\Psi'$  että pisteessä  $\Psi$ .

**Lause 6.** Oletetaan, että funktio  $Q$  täyttää lemmassa 5 esitetyn jatkuvuusehdon. Tällöin kaikki EM-algoritmillä lasketun jonon  $\{\Psi^{(k)}\}$  rajapisteet ovat uskottavuusfunktion  $L(\Psi)$  kriittisiä pisteitä ja uskottavuusfunktion arvot  $L(\Psi^{(k)})$  suppenevat monotonisesti arvoon  $L^* = L(\Psi^*)$  jollakin kriittisellä pisteellä  $\Psi^*$ .

*Todistus.* Lemman 5 jatkuvuusehto on riittävä osoittamaan lemmassa 4 ensimmäisen oletuksen voimassaolo, joten riittää osoittaa lemmassa 4 toinen oletus kaikille pisteille  $\Psi^{(k)} \notin \mathcal{S}$ . Valitaan parametrivektori  $\Psi^{(k)}$  siten, että se on joukon  $\Omega$  sisäpiste. Piste  $\Psi^{(k)}$  maksimoi funktion  $\mathcal{H}(\Psi; \Psi^{(k)})$  joukossa  $\Omega$  yhtälön (30) mukaan, joten sen derivaatta häviää kyseisissä pisteissä, eli  $D^{1,0}\mathcal{H}(\Psi^{(k)}; \Psi^{(k)}) = 0$ . Siispä uskottavuusfunktion derivaatta  $DL(\Psi^{(k)}) = D^{1,0}Q(\Psi^{(k)}; \Psi^{(k)})$ . Uskottavuusfunktion derivaatta  $DL(\Psi^{(k)}) \neq 0$  missä tahansa kriittisten pisteiden joukkoon kuulumattomassa pisteessä  $\Psi^{(k)} \notin \mathcal{S}$ . Tästä seuraa, että piste  $\Psi^{(k)}$  ei ole funktion  $Q(\Psi^{(k)}; \Psi^{(k)})$  lokaali maksimi joukossa  $\Omega$ . M-askeleen määritelmässä olevasta yhtälöstä (24) saadaan, että

$$Q(\Psi^{(k+1)}; \Psi^{(k)}) > Q(\Psi^{(k)}; \Psi^{(k)}).$$

Tämä yhdessä yhtälön (30) kanssa osoittaa, että  $L(\Psi^{(k+1)}) > L(\Psi^{(k)})$  kaikilla pisteillä  $\Psi^{(k)} \notin \mathcal{S}$ . Uskottavuusfunktion konvergenssi seuraa tästä.  $\square$

Nyt on osoitettu, että uskottavuusfunktion arvo suppenee jossakin kriittisessä pisteessä, mutta kriittisen pisteen laatua (minimi, maksimi vai satulapiste) ei vielä ole osoitettu. Lemmassa 7 esitetään kirjan [5] mukaisesti, että EM-algoritmi voi myös konvergoida uskottavuusfunktion satulapisteeseen eikä täten aina saavuta välttämättä edes lokaalia maksimia.

**Lemma 7.** Epäyhtälöstä (30) saadaan, että

$$\frac{\partial \mathcal{H}(\Psi; \Psi^{(k)})}{\partial \Psi} \Big|_{\Psi = \Psi^{(k)}} = 0, \quad (31)$$

eli että funktion  $\mathcal{H}$  derivaatta pisteessä  $\Psi^{(k)}$  häviää. Olkoon  $\Psi_0$  mielivaltainen parametrivektori  $\Psi$ . Valitaan  $\Psi^{(k)} = \Psi_0$  ja käytetään yhtälöä (31) yhtälön (28) derivaatan lausekkeeseen. Tällöin saadaan, että

$$\frac{\partial \log L(\Psi)}{\partial \Psi} \Big|_{\Psi = \Psi_0} = \frac{\partial Q(\Psi; \Psi_0)}{\partial \Psi} \Big|_{\Psi = \Psi_0}. \quad (32)$$

Uskottavuusfunktion  $L$  ja täydellisen datan uskottavuusfunktion ehdollisen odotusarvon  $Q$  derivaatat ovat siis yhtäsuuret pisteessä  $\Psi_0$ . Oletetaan seuraavaksi, että  $\Psi = \Psi^*$  ja että  $\Psi^*$  on funktion  $L(\Psi)$  kriittinen piste. Tällöin saadaan kriittisen pisteen määritelmästä ja yhtälöstä (32), että

$$\frac{\partial \log L(\Psi)}{\partial \Psi} \Big|_{\Psi = \Psi^*} = \frac{\partial Q(\Psi; \Psi_0)}{\partial \Psi} \Big|_{\Psi = \Psi^*} = 0.$$

Saadaan siis, että EM-algoritmissa maksimoitavan funktion  $Q$  kriittiset pisteet ovat myös logaritmisin uskottavuusfunktion kriittisiä pisteitä. Jos piste  $\Psi^*$  on funktion  $Q(\Psi; \Psi^*)$  globaali maksimi joukossa  $\Omega$ , niin EM-algoritmi voi konvergoida uskottavuusfunktion  $L(\Psi)$  satulapisteeseen, sillä kriittisen pisteen  $\Psi^*$  laadusta uskottavuusfunktiolle  $L$  ei ole takeita.

Esitetään vielä lauseessa 8, että uskottavuusfunktio konvergoi lokaaliin maksimiin. Konvergenssia globaaliin maksimiin ei pystytä takaamaan EM-algoritmile kuten ei muillekaan optimointialgoritmeille.

**Lause 8.** Oletetaan, että funktio  $Q$  täyttää lemmän 5 jatkuvuusehdon ja että millä tahansa  $\Psi \in \mathcal{S} \setminus \mathcal{M}$  pätee

$$\sup_{\Psi' \in \Omega} Q(\Psi'; \Psi) > Q(\Psi; \Psi).$$

Tällöin kaikki jonon  $\{\Psi^{(k)}\}$  rajapisteet ovat funktion  $L$  lokaaleja maksimeja ja  $L(\Psi^{(k)})$  konvergoi monotonisesti arvoon  $L^* = L(\Psi^*)$  jollakin lokaalilla maksimilla  $\Psi^*$ .

### 3.2.2 Parametrien jono suppenee

Osoitetaan seuraavaksi, että parametrivektoreiden jono  $\{\Psi^{(k)}\}$  suppenee arvoon  $\Psi^*$ , mikäli uskottavuusfunktion  $L(\Psi)$  arvo suppenee arvoon  $L^* = L(\Psi^*)$ . Merkitään

$$\mathcal{S}(a) = \{\Psi \in \mathcal{S} : L(\Psi) = a\}$$

ja

$$\mathcal{M}(a) = \{\Psi \in \mathcal{M} : L(\Psi) = a\},$$

jotka kuvaavat niitä kriittisten pisteiden joukkoon  $\mathcal{S}$  tai lokaalien maksimien joukkoon  $\mathcal{M}$  kuuluvia pisteitä  $\Psi$ , joissa uskottavuusfunktio saa arvon  $a$ . Lemman 4 oletuksien nojalla  $L(\Psi) \rightarrow L^*$  ja kaikki jonon  $\{\Psi^{(k)}\}$  rajapisteteet kuuluvat joukkoon  $\mathcal{S}(L^*)$ .

Aloitetaan parametrien jonon suppenemisen todistaminen lauseen 9 mukaisella tilanteella, jossa uskottavuusfunktion  $L(\Psi)$  optimi  $L^*$  saavutetaan vain yhdessä kriittisessä pisteessä  $\Psi^*$ .

**Lause 9.** Olkoon jono  $\{\Psi^{(k)}\}$  lemmän 4 oletuksien mukaisella GEM-algoritmilla laskettuja parametrivektorin arvoja. Jos  $\mathcal{S}(L^*) = \{\Psi^*\}$ , missä  $L^*$  on lemmassa 4 esiintyvä funktion  $L(\Psi^{(k)})$  rajapiste, niin silloin  $\Psi^{(k)} \rightarrow \Psi^*$ .

Lauseen 9 ehtoa  $\mathcal{S}(L^*) = \{\Psi^*\}$  voidaan relaksoida olettamalla, että  $\|\Psi^{(k+1)} - \Psi^{(k)}\| \rightarrow 0$ , kun  $k \rightarrow \infty$ . Näin voidaan olettaa, sillä se on välttämätön ehto parametrien jonon suppenemiselle. Lauseessa 10 tätä ehtoa hyödynnetään, kun osoitetaan parametrien jonon  $\{\Psi^{(k)}\}$  suppenevan johonkin kriittisten pisteiden joukon  $\mathcal{S}(L^*)$  pisteeseen.

**Lause 10.** Olkoon jono  $\{\Psi^{(k)}\}$  lemmän 4 oletuksien mukaisella GEM-algoritmilla laskettuja parametrivektorin arvoja. Jos  $\|\Psi^{(k+1)} - \Psi^{(k)}\| \rightarrow 0$ , kun  $k \rightarrow \infty$ , niin kaikki jonon  $\{\Psi^{(k)}\}$  rajapisteteet kuuluvat yhdistettyyn ja kompaktiin joukkoon  $\mathcal{S}(L^*)$ , jossa  $L^*$  on lemmassa 4 esitetty ääriarvokohta funktiolle  $L(\Psi^{(k)})$ . Erityisesti jos  $\mathcal{S}(L^*)$  on diskreetti, niin parametrivektori  $\Psi^{(k)}$  konvergoi johonkin pisteeseen  $\Psi^* \in \mathcal{S}(L^*)$ .

*Todistus.* Oletuksen **OL2** mukaan parametrien arvojen jono  $\{\Psi^{(k)}\}$  on rajoitettu. Rajoitetun jonon  $\{\Psi^{(k)}\}$  rajapisteteiden joukko on yhtenäinen ja kompakti, jos  $\|\Psi^{(k+1)} - \Psi^{(k)}\| \rightarrow 0$ , kun  $k \rightarrow \infty$ . Lemman 4 nojalla jonon  $\{\Psi^{(k)}\}$  rajapisteteet kuuluvat joukkoon  $\mathcal{S}(L^*)$ . Tästä seuraa haluttu tulos.  $\square$

Yleistetään seuraavaksi lauseessa 11 lause 10 tilanteeseen, jossa kriittisten pisteiden joukkoon  $\mathcal{S}$  kuuluvien parametrivektoreiden sijasta tarkastellaan koko määrittelyjoukkoa  $\Omega$ .

**Lause 11.** Olkoon jono  $\{\Psi^{(k)}\}$  lemmän 4 oletuksien mukaisella GEM-algoritmilla laskettuja parametrivektorin arvoja. Määritellään, että

$$\mathcal{L}(L) = \{\Psi \in \Omega : L(\Psi) = L\}.$$

Tehdään myös lisäoletus, että  $D^{1,0}Q(\Psi^{(k+1)}; \Psi^{(k)}) = 0$  ja oletetaan sen olevan jatkuva sekä pisteessä  $\Psi^{(k+1)}$  että pisteessä  $\Psi^{(k)}$ . Tällöin  $\Psi^{(k)}$  suppenee kriittiseen pisteeseen  $\Psi^*$ , jossa  $L(\Psi^*) = L^*$ . Arvo  $L^*$  on uskottavuusfunktion  $L(\Psi^{(k)})$  ääriarvo, jos joko  $\mathcal{L}(L^*) = \{\Psi^*\}$  tai  $\|\Psi^{(k+1)} - \Psi^{(k)}\| \rightarrow 0$ , kun  $k \rightarrow \infty$ , ja  $\mathcal{L}(L^*)$  on diskreetti.

Esitellään seuraavaksi unimodaalin funktion määritelmä, jotta osiossa 4 voidaan hyödyntää seurausta 13, joka osoittaa parametrivektorin konvergoivan unimodaalin uskottavuusfunktion globaaliin maksimiin.

**Määritelmä 12.** (Unimodaali funktio) Uskottavuusfunktio  $L(\Psi)$  on unimodaalinen maksimoinnin suhteen, jos ehdosta  $L(\Psi') > L(\Psi)$  seuraa, että  $\Delta\Psi = \Psi' - \Psi$  on parantava suunta pisteestä  $\Psi$  kaikille pistepareille  $\Psi$  ja  $\Psi'$ . Maksimointitehtävässä suunta  $\Delta\Psi$  on parantava, jos  $\nabla L(\Psi)^T \Delta\Psi > 0$ . [12]

**Seuraus 13.** Oletetaan, että uskottavuusfunktio  $L(\Psi)$  on unimodaalinen joukossa  $\Omega$ ,  $\Psi^*$  on sen ainut kriittinen piste ja  $D^{1,0}Q(\Psi'; \Psi)$  on jatkuva pisteissä  $\Psi'$  ja  $\Psi$ . Tällöin jokaiselle EM-algoritmilla lasketulle jonolle  $\{\Psi^{(k)}\}$  pätee, että parametrivektori  $\Psi^{(k)}$  suppenee uskottavuusfunktion  $L(\Psi)$  maksimoivaan pisteeseen  $\Psi^*$ .

### 3.2.3 Konvergenssiaste

Käydään vielä läpi kirjan [5] mukainen tapa määrittää ja laskea käytännössä EM-algoritmin konvergenssiaste.

Määritelmässä 1 esitelty kuvaus  $M$  on kuvaus joukosta  $\Omega$  itseensä siten, että

$$\Psi^{(k+1)} = M(\Psi^{(k)})$$

kaikilla  $k = 0, 1, 2, \dots$ . Jos parametrivektori  $\Psi^{(k)}$  konvergoi johonkin pisteeseen  $\Psi^*$  ja kuvaus  $M(\Psi)$  on jatkuva, niin piste  $\Psi^*$  on kuvauksen  $M$  kiintopiste eli  $\Psi^* = M(\Psi^*)$ .

Kuvaukselle  $\Psi^{(k+1)} = M(\Psi^{(k)})$  pisteessä  $\Psi^*$  lasketusta Taylorin sarjasta saadaan pisteen  $\Psi^*$  ympäristössä, että

$$\Psi^{(k+1)} - \Psi^* \approx J(\Psi^*)(\Psi^{(k)} - \Psi^*). \quad (33)$$

Yhtälössä (33)  $J(\Psi)$  on  $d \times d$  Jacobin matriisi kuvaukselle  $\mathbf{M}(\Psi) = (M_1(\Psi), \dots, M_d(\Psi))^T$ , jossa alkio

$$J_{ij}(\Psi) = \frac{\partial M_i(\Psi)}{\partial \Psi_j}.$$

Jacobin matriisista  $J(\Psi^*)$  puhutaan usein konvergenssiasteen matriisina.

Parametrivektorin  $\Psi$  todellisen havaitun konvergenssiasteen mitta on globaali konvergenssiaste, joka määritellään raja-arvona

$$r = \lim_{k \rightarrow \infty} \frac{\|\Psi^{(k+1)} - \Psi^*\|}{\|\Psi^{(k)} - \Psi^*\|}, \quad (34)$$

missä  $\|\cdot\|$  on mikä tahansa normi  $d$ -dimensioisessa euklidisessa avaruudessa  $\mathbb{R}^d$ . Merkitään, että Jacobin matriisin  $J(\Psi^*)$  suurin ominaisarvo on  $\lambda_{\max}$ . Tiettyjen säännöllisyysehtojen ollessa voimassa, tiedetään konvergenssiasteen olevan  $r = \lambda_{\max}$ .

Konvergenssiaste voidaan laskea myös komponentteittain. Tällöin komponentin  $i$  konvergenssiaste on

$$r_i = \lim_{k \rightarrow \infty} \frac{\|\Psi_i^{(k+1)} - \Psi_i^*\|}{\|\Psi_i^{(k)} - \Psi_i^*\|}. \quad (35)$$

Tiettyjen ehtojen ollessa voimassa konvergenssiaste  $r$  on komponenttien konvergenssiasteiden  $r_i$  maksimi:

$$r = \max_{1 \leq i \leq d} r_i.$$

Algoritmi tietysti konvergoi jos ja vain jos jokainen komponentti konvergoi. Siispä on luonnollista, että algoritmin konvergenssiaste määräytyy hitaimmin konvergoivan komponentin mukaan, sillä konvergenssia ei saavuteta ennen kuin jokainen komponentti on saavuttanut optimiarvonsa.

Käytännössä yleensä kuitenkin käytetään yhtälöiden (34) ja (35) tilalla muotoja

$$\begin{aligned} r &= \lim_{k \rightarrow \infty} \frac{\|\Psi^{(k+1)} - \Psi^{(k)}\|}{\|\Psi^{(k)} - \Psi^{(k-1)}\|} \quad \text{ja} \\ r_i &= \lim_{k \rightarrow \infty} \frac{\|\Psi_i^{(k+1)} - \Psi_i^{(k)}\|}{\|\Psi_i^{(k)} - \Psi_i^{(k-1)}\|}. \end{aligned}$$

## 4 EM-algoritmin käyttö mixture cure -mallin parametrien laskemiseen

Kuten osiossa 3 todettiin, EM-algoritmi on mixture cure -mallin parametrien laskemiseen sopiva menetelmä, sillä mixture cure -mallissa on latentteja muuttujia. Johdetaan nyt EM-algoritmin käyttämiseen tarvittavat funktiot mixture cure -mallin mukaiselle datalle, implementoidaan saatu algoritmi R-ohjelmointikielellä ja havainnollistetaan sekä dataa että algoritmin konvergenssia kuvilla. Mallin mukaisessa tilanteessa data on sensuroitua, joten vikaantumisaika  $T$  tiedetään tarkasti vain osalle yksilöistä. Mallissa on oletuksena myös, että osa yksilöistä on immuuneja tarkastelun kohteena olevalle tapahtumalle, mutta immuunien ja alttiiden yksilöiden osuuksia ei tiedetä.

### 4.1 Malli

Sovelluksen kohteena olevassa mallissa havaittuina muuttujina ovat aika  $T$ , joka on joko vikaantumisaika  $T^*$  tai sensurointiaika  $C$  yhtälön (6) mukaisesti, ja yhtälön (7) mukaisesti määritetty sensurointi-indeksi  $\delta$ . Voidaan siis osiossa 3 käytettyjen merkintöjen mukaisesti kirjoittaa, että  $\mathbf{Y} = (T, \delta)$ . Osiossa 2.3 esitelty alttiusindeksi  $\zeta$  on tässä tapauksessa vaillinaista dataa.

Muuttuja  $\zeta_i$  kuvaa siis yksilön alttiutta tehdä tarkastelun kohteena oleva tapahtuma ja se on näyte satunnaismuuttujasta  $\zeta \sim \text{Bern}(\pi)$ . Täydelliselle datalle muuttujan  $\zeta_i$  arvot tiedetään, mutta tässä tapauksessa se on latentti muuttuja, joten sen arvoja ei tiedetä. Yksilön vikaantumisaika  $t_i^*$  on realisaatio satunnaismuuttujasta  $T^* \sim \text{Exp}(\lambda)$  ja näitä merkintöjä käyttäen, havaittu aika  $t_i = \min(t_i^*, c_i)$  ja sensurointi-indeksi  $\delta_i = \mathbb{I}(t_i^* \leq c_i)$ , kun yksilön  $i$  sensurointiaika  $c_i$  on joko erikseen määritetty tai se on näyte jotakin satunnaisjakaumaa noudattavasta satunnaismuuttujasta  $C$ . Sensurointiaika on näyte satunnaismuuttujasta  $C$ , kun kyseessä on osiossa 2 mainittu satunnainen sensurointi.

Oletetaan, että vikaantumisaika  $T$  noudattaa eksponentijakaumaa, joten alttiin populaation välttöfunktio on kaavan (10) mukaan muotoa

$$S(t|\zeta = 1) = P(T > t|\zeta = 1) = e^{-\lambda t}$$

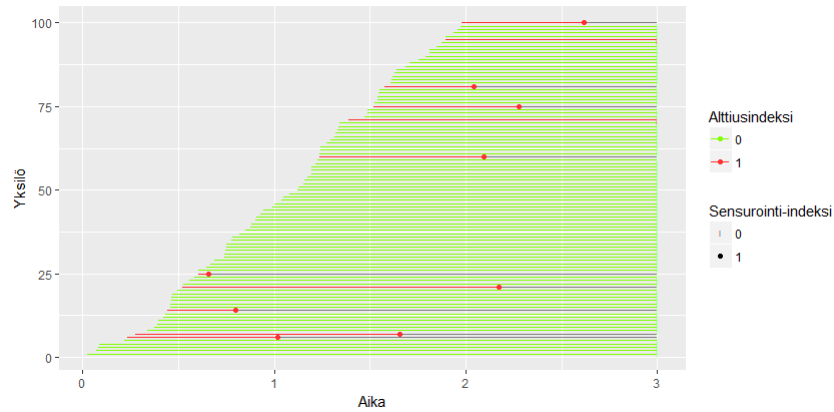
ja tiheysfunktio on kaavan (12) mukaisesti

$$f(t|\zeta = 1) = P(T = t|\zeta = 1) = \lambda e^{-\lambda t}.$$

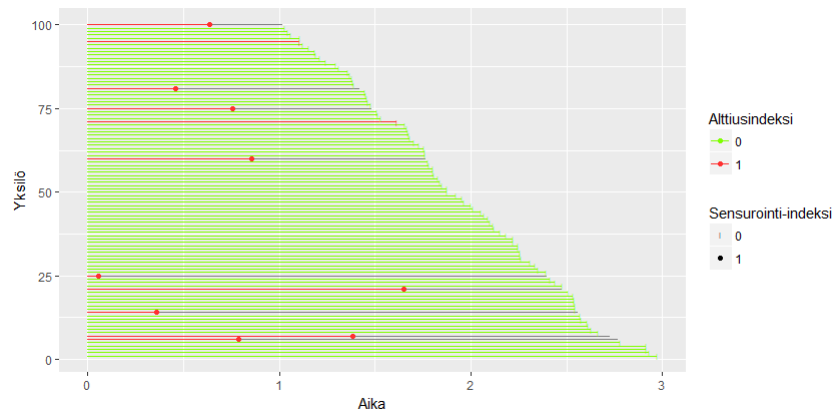
Oletetaan, että alttiiden yksilöiden osuus koko populaatiosta on  $P(\zeta = 1) = \pi$ . Tällöin immuunien yksilöiden osuus  $P(\zeta = 0) = 1 - \pi$ , sillä kyseessä on kahden ryhmän sekoitemalli. Saadaan, että mallissa on kaksi

EM-algoritmilla estimoitavaa parametria  $\theta = (\pi, \lambda)$ , jotka ovat tapahtuman tekevien yksilöiden osuus ja vakiovikaantumistahti.

Kuvassa 9 havainnollistetaan, miltä tällainen data näyttää, kun aloitusajat on piirretty kalenteriajan mukaisesti ja sensurointi-aika on kaikille sama tyyppin I sensuroinnin mukaisesti. Samaa dataa voidaan havainnollistaa myös kuvan 10 mukaisesti. Siinä havainnollistetaan seuranta-aikaa. Kuvasta 10 on helpompi vertailla yksilöiden vikaantumisaikoja. Data on generoitu siten, että aloittaminen tapahtuu yhden aikayksikön aikana ja maksimiseuranta-aika on kolme aikayksikköä. Generoidussa datassa on 100 yksilöä ja generoinnissa käytetty alttiiden yksilöiden osuus on 10 %.



Kuva 9: Yksilöiden seuranta-ajat kalenterin mukaisesta aloitusajankohdasta tyyppin I mukaiseen sensurointiajankohtaan.



Kuva 10: Yksilöiden seuranta-aikojen pituudet tyyppin I sensuroinnin tapauksessa.

Mixture cure -mallin uskottavuusfunktio ja muut EM-algoritmin käyttöä varten tarvittavat funktiot on johdettu artikkelissa [13] tilanteelle, jossa vikaantumisaika  $T$  noudattaa Weibull-jakaumaa. Tämän tutkielman sovelluksen tilanne on Weibull-jakauman erikoistapaus, joten esitetään nyt juuri tämän tapauksen funktiot, jotka on itse johdettu.

Taulukossa 2 esitetään todennäköisyydet, että yksilö vikaantuu tai selviää ajanhetkellä  $t$  ehdolla, että se on immuuni tai altis. Näiden ehdollisten todennäköisyyksien avulla saadaan marginaalinen välttöfunktio (22) kirjoitettua eksplisiittisesti. Se on tällöin muotoa

$$\begin{aligned} S(t) &= P(T > t) = P(T > t | \zeta = 0) P(\zeta = 0) + P(T > t | \zeta = 1) P(\zeta = 1) \\ &= 1 \cdot (1 - \pi) + e^{-\lambda t} \cdot \pi \\ &= 1 - \pi + \pi e^{-\lambda t}. \end{aligned} \tag{36}$$

Vastaavalla tavalla saadaan laskettua marginaalinen tiheysfunktio

$$\begin{aligned} f(t) &= P(T = t) = P(T = t | \zeta = 0) P(\zeta = 0) + P(T = t | \zeta = 1) P(\zeta = 1) \\ &= 0 \cdot (1 - \pi) + \lambda e^{-\lambda t} \cdot \pi \\ &= \pi \lambda e^{-\lambda t}. \end{aligned}$$

Taulukko 2: Välttö- ja tiheysfunktioiden laskemiseen tarvittavat ehdolliset todennäköisyydet.

$P(T \zeta)$	$T = t$	$T > t$
$\zeta = 0$	0	1
$\zeta = 1$	$\lambda e^{-\lambda t}$	$e^{-\lambda t}$

Epätäydellisen data uskottavuusfunktio on kaavan (8) mukaisesti muotoa

$$L(\boldsymbol{\theta} | \mathbf{t}, \boldsymbol{\delta}) = \prod_{i=1}^n (\pi \lambda e^{-\lambda t_i})^{\delta_i} (1 - \pi + \pi e^{-\lambda t_i})^{1-\delta_i} \tag{37}$$

ja sen logaritminen uskottavuusfunktio joko ottamalla logaritmi yhtälöstä (37) tai kaavan (9) perusteella on

$$l(\boldsymbol{\theta} | \mathbf{t}, \boldsymbol{\delta}) = \sum_{i=1}^n \left\{ \delta_i (\log \pi + \log \lambda - \lambda t_i) + (1 - \delta_i) \log (1 - \pi + \pi e^{-\lambda t_i}) \right\}.$$

EM-algoritmia käytettäessä tarvitaan kuitenkin epätäydellisen datan uskottavuusfunktion sijasta täydellisen datan uskottavuusfunktio. Se on mixture cure -mallin tapauksessa muotoa

$$L_c(\boldsymbol{\theta}|\mathbf{t}, \boldsymbol{\delta}, \boldsymbol{\zeta}) = \prod_{i=1}^n (\pi \lambda e^{-\lambda t_i})^{\delta_i} \left[ (1 - \pi)^{\mathbb{I}(\zeta_i=0)} (\pi e^{-\lambda t_i})^{\mathbb{I}(\zeta_i=1)} \right]^{1-\delta_i}. \quad (38)$$

Otetaan yhtälöstä (38) logaritmi, jolloin saadaan täydellisen datan logaritminen uskottavuusfunktio lopulta muotoon

$$l_c(\boldsymbol{\theta}|\mathbf{t}, \boldsymbol{\delta}, \boldsymbol{\zeta}) = \sum_{i=1}^n \{ \delta_i [\log \pi + \log \lambda - \lambda t_i] + (1 - \delta_i) [\mathbb{I}(\zeta_i = 0) \log(1 - \pi) + \mathbb{I}(\zeta_i = 1) (\log \pi - \lambda t_i)] \}. \quad (39)$$

EM-algoritmin E-askeleessa olevan funktion (23) muodostamisessa tarvitaan taulukossa 3 esitetyt ehdolliset todennäköisyydet. Taulukon 3 arvot on laskettu seuraavilla kaavoilla

$$P(\zeta_i = j | T_i > t_i, \boldsymbol{\theta}^{(k-1)}) = \frac{P(T_i > t_i | \zeta_i = j, \boldsymbol{\theta}^{(k-1)}) P(\zeta_i = j)}{P(T_i > t_i | \zeta_i = 0, \boldsymbol{\theta}^{(k-1)}) P(\zeta_i = 0) + P(T_i > t_i | \zeta_i = 1, \boldsymbol{\theta}^{(k-1)}) P(\zeta_i = 1)} \text{ ja} \quad (40)$$

$$P(\zeta_i = j | T_i = t_i, \boldsymbol{\theta}^{(k-1)}) = \frac{P(T_i = t_i | \zeta_i = j, \boldsymbol{\theta}^{(k-1)}) P(\zeta_i = j)}{P(T_i = t_i | \zeta_i = 0, \boldsymbol{\theta}^{(k-1)}) P(\zeta_i = 0) + P(T_i = t_i | \zeta_i = 1, \boldsymbol{\theta}^{(k-1)}) P(\zeta_i = 1)}, \quad (41)$$

joissa  $j \in \{0, 1\}$ . Kaavoilla (40) ja (41) lasketaan ehdolliset todennäköisyydet, että yksilö on immuuni tai altis, jos se on selvinnyt tai vikaantunut ajanhetkellä  $t$ .

Taulukko 3: Funktion  $Q$  muodostamiseen tarvittavat ehdolliset todennäköisyydet.

$P(\zeta T)$	$\zeta = 0$	$\zeta = 1$
$T > t$	$\frac{1-\pi}{1-\pi+\pi e^{-\lambda t}}$	$\frac{\pi e^{-\lambda t}}{1-\pi+\pi e^{-\lambda t}}$
$T = t$	0	1

Nyt yhtälöt (39), (40) ja (41) yhdistämällä saadaan EM-algoritilla mak-

simoitava funktio (23) muotoon

$$\begin{aligned}
Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(k-1)}) &= E_{\zeta|T, \boldsymbol{\theta}^{(k-1)}} [l(\boldsymbol{\theta}|t, \delta, \zeta)] \\
&= \sum_{i=1}^n \sum_{j=0}^1 P(\zeta_i = j | T_i > t_i, \boldsymbol{\theta}^{(k-1)})^{1-\delta_i} P(\zeta_i = j | T_i = t_i, \boldsymbol{\theta}^{(k-1)})^{\delta_i} l_c(\boldsymbol{\theta}|t_i, \delta_i, \zeta_i) \\
&= \sum_{i:\delta_i=0} \left[ \frac{1-\pi^{(k-1)}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} \log(1-\pi) + \frac{\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} (\log \pi - \lambda t_i) \right] \\
&+ \sum_{i:\delta_i=1} [\log \pi + \log \lambda - \lambda t_i],
\end{aligned} \tag{42}$$

kun viimeisintä muotoa edeltävät välivaiheet on jätetty kirjoittamatta. Derivoimalla yhtälö (42) erikseen kummankin parametrin,  $\pi$  ja  $\lambda$ , suhteen ja ratkaisemalla saatujen derivaattafunktioiden nollakohdat, saadaan parametreille laskettua estimaatit

$$\pi^{(k)} = \frac{1}{n} \left[ \sum_{i:\delta_i=0} \frac{\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} + \sum_{i:\delta_i=1} 1 \right] \quad \text{ja} \tag{43}$$

$$\lambda^{(k)} = \frac{\sum_{i:\delta_i=1} 1}{\sum_{i:\delta_i=0} \frac{\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}}{1-\pi^{(k-1)}+\pi^{(k-1)}e^{-\lambda^{(k-1)}t_i}} + \sum_{i:\delta_i=1} t_i}. \tag{44}$$

Osiossa 3 osoitettiin lemmassa 7, että EM-algoritmilla lasketut funktion  $Q$  maksimoivat parametrien arvot ovat myös täydellisen datan uskottavuusfunktion kriittisiä pisteitä, joten kaavoilla (43) ja (44) lasketut arvot maksimoivat myös uskottavuusfunktion (39).

Täydellisen datan logaritminen uskottavuusfunktio (39) saadaan jaettua kahteen osaan siten, että toinen osa sisältää kaikki parametria  $\pi$  sisältävät termit ja toinen sisältää kaikki parametria  $\lambda$  sisältävät termit. Voidaan siis merkitä, että

$$l_c(\pi, \lambda | \mathbf{t}, \boldsymbol{\delta}, \boldsymbol{\zeta}) = l_\pi(\pi | \mathbf{t}, \boldsymbol{\delta}, \boldsymbol{\zeta}) + l_\lambda(\lambda | \mathbf{t}, \boldsymbol{\delta}, \boldsymbol{\zeta}). \tag{45}$$

Yhtälössä (45) oleva funktio  $l_\pi$  on muotoa

$$\begin{aligned}
l_\pi(\pi | \mathbf{t}, \boldsymbol{\delta}, \boldsymbol{\zeta}) &= \\
&\sum_{i=1}^n \{ \delta_i \log \pi + (1 - \delta_i) [\mathbb{I}(\zeta_i = 0) \log(1 - \pi) + \mathbb{I}(\zeta_i = 1) \log \pi] \}
\end{aligned} \tag{46}$$

ja funktio  $l_\lambda$  on muotoa

$$l_\lambda(\lambda|\mathbf{t}, \boldsymbol{\delta}, \boldsymbol{\zeta}) = \sum_{i=1}^n \{\delta_i [\log \lambda - \lambda t_i] - (1 - \delta_i) \mathbb{I}(\zeta_i = 1) \lambda t_i\}. \quad (47)$$

Sekä funktio (46) että funktio (47) koostuu konkaavien termien summasta, sillä logaritmifunktio on konkaavi samoin kuin lineaariset termit ovat. Konkaavien funktioiden summana myös funktio (45) on konkaavi. Toisaalta, konkaavit funktion ovat unimodaalisia maksimoinnin suhteen [12], joten seurauksen 13 mukaan tiedetään siis, että EM-algoritmi konvergoi täydellisen datan logaritmisen uskottavuusfunktion globaaliin maksimiin.

Täydelliselle datalle pystytään suoraan laskemaan paramertien  $\pi$  ja  $\lambda$  optimaaliset arvot. Derivoidaan yhtälö (46) muuttujan  $\pi$  suhteen ja yhtälö (47) muuttujan  $\lambda$  suhteen. Täydellisen datan logaritmiselle uskottavuusfunktiolle optimaaliset arvot saadaan selvittämällä laskettujen derivaattojen nollakohdat. Datan ollessa täydellistä, tiedetään, mitkä yksilöistä ovat immuuneja ja mitkä alttiita. Siitä syystä parametrien  $\pi$  ja  $\lambda$  optimiarvot saadaan ratkaistua sijoittamalla havaitut muuttujien  $t$ ,  $\delta$  ja  $\zeta$  arvot yhtälöihin

$$\pi = \frac{1}{n} \sum_{i=1}^n [\delta_i + (1 - \delta_i) \mathbb{I}(\zeta_i = 1)] = \frac{1}{n} \sum_{i=1}^n \zeta_i \quad \text{ja} \quad (48)$$

$$\lambda = \frac{\sum_{i=1}^n \delta_i}{\sum_{i=1}^n t_i [\delta_i + (1 - \delta_i) \mathbb{I}(\zeta_i = 1)]} = \frac{\sum_{i=1}^n \delta_i}{\sum_{i=1}^n t_i \zeta_i}. \quad (49)$$

Yhtälöitä (48) ja (49) ei siis voida käyttää puuttuvalle datalle parametrien  $\pi$  ja  $\lambda$  laskemiseen, sillä puuttuvalle datalle ei tiedetä muuttujan  $\zeta$  arvoa. Tästä syystä käytetään EM-algoritmia parametrien optimaalisten arvojen selvittämiseen.

## 4.2 Toteutus R-ohjelmointikielellä

Funktioiden johtamisen jälkeen toteutettiin saman asian R-ohjelmointikielellä RStudio-ohjelmointiympäristöä käyttämällä. Luetteloinnissa 1 esitetään funktio `cureEM`, jolla saadaan laskettua mixture cure -mallin parametrien optimiarvot EM-algoritmeilla. Funktio `cureEM` ottaa kaksi parametria: ensimmäinen näistä, `params`, on parametrien  $\pi$  ja  $\lambda$  alkuarvot sisältävä vektori ja toinen, `D`, on havaittu data. Havaittu data koostuu muuttujista `time`, `censored` ja `weight`. Muuttuja `time` on havaittu aika  $t_i$  ja muuttuja `censored` on sensurointi-indeksi  $\delta_i$ . Muuttuja `weight` kuvaa todennäköisyyttä, että yksilö on altis tarkastelun kohteena olevalle tapahtumalle, ja se lasketaan jokaisella iteraatiokierröksellä uudelleen yhtälöitä (40) ja (41) käyttäen arvolla  $j = 1$ .

```

1 #EM- algoritmi mixture cure -mallille
2
3 cureEM <- function(params, D){
4
5     lambda_prev <- params[1]
6     pi_prev <- params[2]
7     n <- length(D$time)
8     time_censored <- D$time[D$censored==0]
9     time_not_censored <- D$time[D$censored==1]
10
11 #Lopetuskriteeri
12 eps <- 10^(-5)
13 continue <- TRUE
14
15 while(continue){
16
17     #Kertoimien laskeminen
18     D$weight[D$censored == 0] <- pi_prev*exp(-lambda_prev*time_censored) /
19       (1-pi_prev + pi_prev*exp(-lambda_prev*time_censored))
20     D$weight[D$censored == 1] <- 1
21
22     #Parametrien maksimointi
23     pi <- sum(D$weight)/n
24
25     lambda <- sum(D$weight[D$censored == 1])/
26       (sum(D$weight[D$censored == 0]*time_censored) +
27         sum(D$weight[D$censored == 1]*time_not_censored))
28
29     #Lopetuskriteerin toteutumisen tarkistus ja parametrien arvojen paivitys
30     if(sqrt((pi_prev - pi)^2 + (lambda_prev - lambda)^2) < eps){
31         continue <- FALSE
32     }else{
33         lambda_prev <- lambda
34         pi_prev <- pi
35     }
36
37 }
38
39 result <- c(lambda, pi)
40 return(result)
41 }

```

Luettelointi 1: Mixture cure -mallin parametrien laskeminen EM-algoritmilla.

Funktion `cureEM` toiminnan oikeellisuus varmistettiin generoimalla testi-dataa luetteloinnin 2 mukaisella tavalla. Data on generoitu parametrien arvoilla  $\pi = 0.3$  ja  $\lambda = 1$ , joten tiedetään funktion `cureEM` toimivan, kun sen laskemat arvot vastasivat datan generoinnissa käytettyjä arvoja.

```

1 #Luodaan testidataa
2 library(LaplacesDemon)
3 tmax <- 2
4 n <- 1000
5 cure_fraction <- 0.7
6 cured <- rbern(n, cure_fraction)
7 testidata <- data.frame(time = rexp(n)*(1-cured) + tmax*cured,
8                          cured = 1-cured)
9 testidata$censored <- as.integer(testidata$time < tmax)
10 testidata$time[which(testidata$censored == FALSE)] <- tmax
11 testidata$weight <- 1

```

### Luettelointi 2: Testidatan generoiminen.

Tämän lisäksi laskettiin R:n valmiilla optimointifunktiolla `optim` logaritmisien uskottavuusfunktion (39) maksimoivat parametrien arvot luetteloinnissa 3 esitetyllä koodilla. Tällä tavoin varmistettiin, että funktio `cureEM` tuottaa saman ratkaisun kuin jokin muukin optimointimenetelmä. Mixture cure -mallille halutaan löytää uskottavuusfunktion maksimi ja funktio `optim` ratkaisee annetun funktion minimin, joten täytyy muistaa, että

$$\arg \max f(x) = \arg \min -f(x).$$

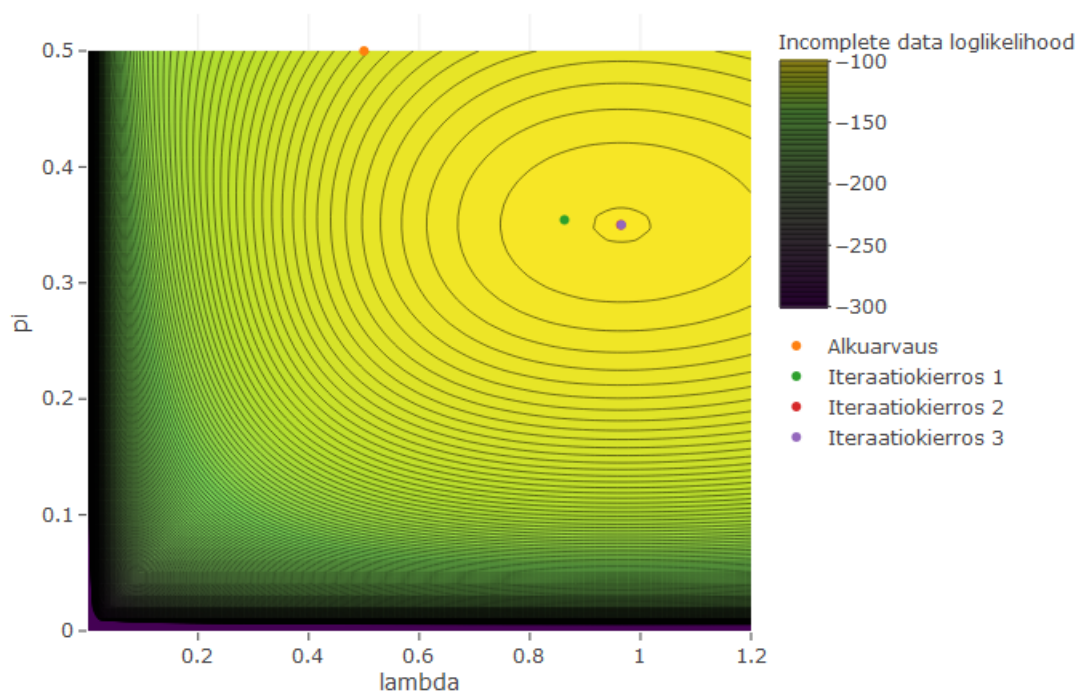
```

1 #Logaritminen uskottavuusfunktio
2 loglikelihood <- function(params){
3
4   lambda <- params[1]
5   pi <- params[2]
6   delta <- testidata$censored
7   t <- testidata$time
8   cured <- testidata$cured
9
10  -sum(delta*(log(pi) + log(lambda) - lambda*t) +
11        (1-delta)*((cured == 0)*log(1-pi) + (cured == 1)*(log(pi) -
12                    lambda*t)))
13 }
14 #Optimointi funktiolla optim()
15 #Alkuarvaus parametreille lambda ja pi
16 init <- c(0.5, 0.5)
17 optim(init, fn = loglikelihood)

```

### Luettelointi 3: Täydellisen datan logaritminen uskottavuusfunktio ja valmiilla optimointifunktiolla optimaalisten parametrien ratkaiseminen

Havainnollistetaan algoritmin konvergenssia kuvalla 11. Siitä nähdään, että epätäydellisen datan logaritmisien uskottavuusfunktion suurin arvo saavutetaan suunnilleen arvoilla  $\lambda = 1$  ja  $\pi = 0.35$ . Lisäksi nähdään, että tällä datalla luetteloinnin 1 funktio konvergoi jo muutaman iteraatiokierroksen jälkeen uskottavuusfunktion optimiin.



Kuva 11: Generoidun datan konvergenssi epätäydellisen datan logaritmisesti uskottavuusfunktion maksimiin.

## 5 Pelidataan soveltaminen

Käytetään osiossa 2 esiteltyä teoriaa peleistä kerättyyn dataan. Monissa peleissä on mahdollisuus ostotapahtumien tekemiseen peliä pelatessa ja niiden tekeminen on pelin tekijän kannalta jopa toivottavaa. Kaikki pelaajat eivät kuitenkaan ostotapahtumaa tee ja osa pelaajista tekee useita ostoja. Pelaajat voidaan siis jaotella kahteen osapopulaatioon, jotka osion 2 käsitteiden mukaisesti ovat alttiiden pelaajien joukko ja immuunien pelaajien joukko. Ensin mainittuun siis kuuluvat ne, jotka tekevät ostotapahtuman jossain vaiheessa, ja jälkimmäiseen kuuluvat ne, jotka eivät ostotapahtumaa tee, vaikka niitä seurattaisiin äärettömän pitkään. Tilanne on siis osiossa 4.1 esitellyn mallin mukainen. Ostotapahtuman tekeillä pelaajilla alttiusindeksi  $\zeta = 1$  ja ostotapahtuman tekemiselle immuuneille pelaajille  $\zeta = 0$ .

Pelistä on kerätty dataa tietyllä aikavälillä, joten dataa on sensuroitu tyyppin I sensuroinnin mukaisesti. Kaikilla pelaajilla on siis kalenterin mukaisessa ajassa sama sensurointiajankohta – havainnoinnin päättymisen ajankohta. Aiemmin esitelty sensurointi-indeksi  $\delta$  toimii tässäkin tapauksessa samalla tavalla:  $\delta = 0$ , jos pelaaja on sensuroitu, ja  $\delta = 1$ , jos pelaaja on tehnyt ostotapahtuman.

Joissain peleissä pelaajat aloittavat hyvin lyhyen aikavälin aikana, joten seuranta-aika jokaiselle pelaajalle on likimain sama. Toisissa peleissä aloitukset tapahtuvat pidemmän aikavälin aikana. Aloitusajankohdan jakautumiseen vaikuttaa se, että peliä testatessa sitä mainostetaan. Pelaajia saattaa silloin tulla hyvinkin lyhyen ajan sisällä tai sitten niitä tulee koko mainoskampanjan ajan tasaisesti. Aikaisemmin julkaistuille peleille uusia pelaajia tulee yleensä melko tasaista tahtia. Pelidatan tapauksessa havaittu aika  $t$  on aika aloituksesta joko ostotapahtuman tekemiseen tai datan keräämisen päättymiseen.

Osion 4.1 mallissa käytetyt parametrit,  $\lambda$  ja  $\pi$ , ovat tässä tapauksessa konversiotahdi ja monetisaatioprosentti, ja ne lasketaan osion 4 mukaisella tavalla. Monetisaatiolla tarkoitetaan ostotapahtuman tekemistä. Pelaaja siis monetisoituu tehdessään ensimmäisen ostotapahtuman. Konversiotahdilla tarkoitetaan monetisoitumisnopeutta. Vikaantumisaikojen oletetaan noudataavan eksponenttijakaumaa, joten konversiotahdi  $\lambda$  on vakio. Yksittäiselle pelille malli siis ennustaa, kuinka iso osa pelaajista tekee ostotapahtuman ja kuinka nopeasti ensimmäinen ostotapahtuma tehdään. Tässä tutkielmassa tavoitteena on tutkia, missä vaiheessa ja kuinka tarkasti malli pystyy lopullisen monetisaatioprosentin ennustamaan. Osan testeistä tein Hipster sheep-nimisestä pelistä [1] kerätylle datalle ja osan tein mallin mukaiselle generoidulle datalle.

## 5.1 Käytetyn datan kuvailu ja muokkaus

Käyttämäni data on kerätty Hipster sheep -nimisestä pelistä 11.12.2014 ja 4.1.2017 välisenä aikana. Data koostuu kolmesta taulukosta, joiden nimet ovat *sheepsters\_device*, *sheepsters\_end\_session* ja *sheepsters\_store*.

Taulukossa *sheepsters\_device* on seitsemän muuttujaa – *id*, *nr*, *version*, *timestamp*, *platform*, *country* ja *language* – ja 20 295 havaintoa. Taulukossa 4 näytetään viisi ensimmäistä riviä tästä datasta. Tämä data kertoo pelaajan ensimmäisen pelin alkamisajankohdan sekä mitä versiota, millä alustalla (android tai ios), missä maassa ja millä kielellä pelaaja on peliä pelannut. Taulukko koostuu siis ensimmäisellä pelikerralla kerättävistä tiedoista.

Taulukko 4: Sheepsters\_device-datan ensimmäiset viisi riviä.

id	nr	version	timestamp	platform	country	language
gp_0010e805-3a50-47ef-a932-4752945ecec	1	1.18	2015-10-23 10:19:22	android	GB	en
gp_0014cc9a-2b20-49ba-a1a1-48d0679dfec2	1	1.18	2015-10-24 09:32:23	android	GB	en
gp_001598a0-4988-4431-a465-b42af85bec93	1	1.21	2015-12-14 14:42:33	android	CN	zh
gp_0029862d-3236-4b5d-bf19-be399b98d4fb	1	1.22b	2015-12-09 10:14:41	android	US	en
gp_0031cfeb-ee08-479e-a07e-382c94c0b764	1	1.18	2015-10-23 09:15:14	android	US	en

Taulukossa *sheepsters\_end\_session* on kuusi saraketta – *id*, *nr*, *version*, *timestamp*, *number* ja *duration* – ja 299 408 riviä. Tällä kertaa muuttuja *timestamp* kuvaa pelin lopettamisajankohtaa. Muuttuja *duration* nimensä mukaisesti kuvaa pelin kestoa. Taulukossa 5 esitellään tämän datan viisi ensimmäistä riviä.

Taulukko 5: Sheepsters\_end\_session-datan ensimmäiset viisi riviä.

id	nr	version	timestamp	number	duration
gp_0029862d-3236-4b5d-bf19-be399b98d4fb	6	1.22b	2015-12-09 10:14:50	1	9
gp_00d19508-90fe-431e-a117-c69fc382748d	6	1.32	2016-07-23 06:32:43	1	0
gp_010702a0-ca52-4f0c-8443-f8ffa4fbd9bf	6	1.35	2016-09-23 12:58:15	1	2
gp_013f2ec1-33d6-4e1d-b51c-17ad208db600	6	1.31	2016-07-12 02:49:49	1	0
gp_0161e2d0-4804-4363-bf5c-8bdaf08f6dec	6	1.33	2016-10-22 13:58:00	1	2

Taulukossa *sheepsters\_store* puolestaan on ostotapahtuman tehneet pelaajat, ostotapahtuman ajankohta ja ostetun tavaran nimi. Muuttujia on siis viisi: *id*, *nr*, *version*, *timestamp* ja *purchase*. Havaintoja tässä taulukossa on yhteensä vain 1 018 ja niistä viisi ensimmäistä esitetään taulukossa 6.

Dataa pitää muokata ennen kuin se pystytään syöttämään listauksessa 1 esitetylle funktiolle ja pystytään laskemaan arvot osiossa 4.1 esitellyn mallin parametreille  $\lambda$  ja  $\pi$ . Ensin poistetaan yksilöt, joilla jokin havaituista ajoista on havainnointivälin ulkopuolella. Tämän jälkeen otetaan *sheepsters\_end\_session* taulukosta vain kunkin pelaajan viimeinen havaittu lopetusajankohta. Vastaavasti ostotapahtumien taulukosta otetaan vain

Taulukko 6: Sheepsters\_store-datan viisi ensimmäistä riviä.

id	nr	version	timestamp	purchase
gp_ef98c216a189eb83	6	1.9	2015-05-07 09:24:08	jarful_of_diamonds
macosx_eaed619e-f0e7-468b-9993-8ad944fd0ecc	12	1.22b	2015-12-07 14:15:38	jarful_of_diamonds
macosx_eaed619e-f0e7-468b-9993-8ad944fd0ecc	12	1.22b	2015-12-07 14:22:37	jarful_of_diamonds
macosx_855cb5e4-01ca-4081-91f8-30cfad0b8e9d	14	1.35b	2016-08-29 10:55:20	backbag_of_diamonds
gp_3bfb710c2bf9bd5a	15	1.13b	2015-07-13 07:51:24	jarful_of_diamonds

ensimmäisten ostotapahtumien ajankohdat. Uniikit pelaajat on tunnistettu muuttujien *id* ja *version* pariin avulla. Osa pelaajista on pelannut useaa eri versiota, joten tästä syystä pelkästään muuttujan *id* käyttäminen tuottaisi ongelmia.

Aloitus- ja lopetusaikojen taulukot yhdistetään siten, että otetaan vain ne muuttujien *id* ja *version* yhdistelmät, jotka esiintyvät kummassakin taulukossa. Saatuun taulukkoon lisätään vielä tieto ostotapahtuman ajankohdasta niille pelaajille, jotka sellaisen tekivät, ja lopuilla pelaajilla muuttujan *timestamp\_store* arvo on puuttuva. Tämän jälkeen poistetaan datasta vielä ne pelaajat, joilla aloitus- ja lopetusajankohdat eivät ole oikeassa järjestyksessä tai ostotapahtuman ajankohta on ennen aloittamista tai vasta lopettamisen jälkeen. Oikea järjestys on siis  $timestamp\_start < timestamp\_store < timestamp\_end$ , kun ajankohtien vertailuun käytetään kirjastossa `anytime` olevaa funktiota `anytime`, joka muuttaa ajankohdat sekuntimääräksi vuoden 1970 alusta laskettuna.

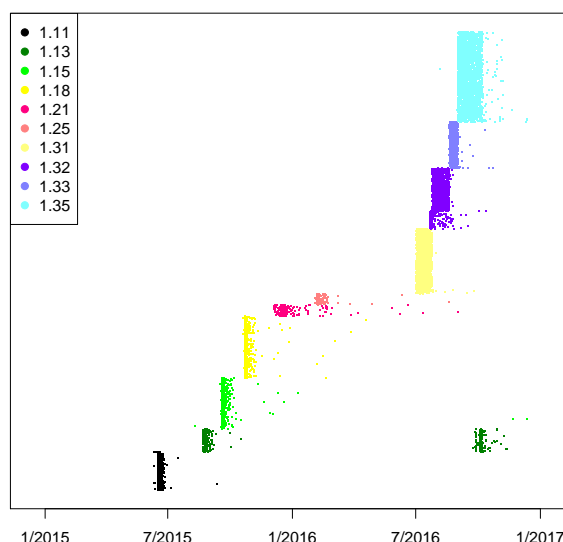
Pelaajat ovat pelanneet yhteensä 68:aa eri versiota. Data rajataan siten, että tarkastellaan vain pelaajia, jotka pelasivat jotakin kymmenestä suosituimmasta versiosta. Taulukossa 7 esitetään näiden versioiden pelaajamäärät ja ostotapahtuman tehneiden pelaajien lukumäärät. Huomataan, että kovinkaan moni pelaaja ei ole tehnyt ostotapahtumaa.

Taulukko 7: Suosituimpien versioiden pelaajien ja ostotapahtumien määrät.

Versio	1.11	1.13	1.15	1.18	1.21	1.25	1.31	1.32	1.33	1.35
Pelaajia	1006	601	1338	1604	309	287	1691	1582	1211	2366
Ostajia	1	1	2	6	6	5	24	21	18	35

Myös pelaajien aloitustahdeissa on eroja, kuten nähdään kuvasta 12. Versioissa 1.31, 1.32, 1.33 ja 1.35 aloitusajat ovat jakautuneet laajemmalle aikavälille kuin muissa versioissa. Versiossa 1.13 puolestaan on mielenkiintoista, että moni on alkanut pelaamaan sitä reilu vuotta myöhemmin kuin ensimmäiset sitä versiota pelanneet. Tämän tutkielman myöhemmissä vaiheissa keskitytään vain versioon 1.18 ja sitä myöhempiin versioihin, sillä käytettä-

vä malli olettaa, että ostotapahtumia on tehty, ja versiota 1.18 aiemmissa versioissa niitä ei juurikaan ole.



Kuva 12: Suosituimpia versioita pelaajaavien pelaajien aloitusajat.

Osiossa 4.1 kuvailussa tilanteessa havaittu data  $\mathbf{y} = (t, \delta)$ , joten selvitetään havaitut ajat  $t_i$  ja sensurointi-indeksit  $\delta_i$  Hipster sheep -pelin datalle. Pelaajan havaittu aika on ensimmäisen pelin aloitusajankohdasta joko ostotapahtuman ajankohtaan tai sensurointiajankohtaan mennessä kulunut aika. Havaitun ajan laskemiseen käytetään R:n valmista funktiota `difftime` ja yksikön määritin olevan päiviä. Funktio `difftime` laskee kahden merkkijonona annetun ajankohdan välisen ajan. Sensurointi-indeksit ovat kuten aiemminkin ja ne määritetään havaitun ajan laskemisen yhteydessä.

## 5.2 Mallin oletuksien toteutuminen

Nyt käytettävä data on saatu mallin vaatimaan muotoon, joten tarkastellaan seuraavaksi osiossa 4.1 esitetyn mallin oletuksien toteutumista. Ensimmäinen oletus on, että osa pelaajista on immuuneja ja osa alttiita tarkastelun kohteena olevalle tapahtumalle, eli tässä tapauksessa ostotapahtuman tekemiselle. Toinen oletus on, että havaittujen ostotapahtumien ajankohdat noudattavat eksponenttijakaumaa.

Käytetään artikkelissa [14] esiteltyä Akaike-informaatiokriteeriä immuunien yksilöiden olemassaolon selvittämiseen. Sen avulla voidaan selvittää dataa parhaiten kuvaava malli useiden vaihtoehtojen joukosta. Lasketaan vaihtoehtoina oleville malleille AIC-arvot kaavalla

$$\text{AIC} = 2k - 2 \log(\hat{L}),$$

jossa vakio  $k$  on mallin parametrien lukumäärä ja muuttuja  $\hat{L}$  on suurimman uskottavuuden estimaatti, ja valitaan malleista se, jolla saadaan pienin AIC-arvo.

Tässä tapauksessa vertailtavia malleja on kaksi. Ensimmäisessä näistä oletetaan, että  $\pi = 1$ , mikä tarkoittaa, että ostotapahtumalle immuuneja pelaajia ei ole eli kaikki pelaajat tekevät ostotapahtuman jossain vaiheessa. Tässä tapauksessa mixture cure -mallin käyttäminen ei olisi järkevää. Toinen vertailtava malli on mixture cure -malli, jolloin datan tulee olla sen oletuksien mukaista eli osan pelaajista tulee olla immuuneja monetisoitumiselle. Tällöin tapahtuman tekevien pelaajien osuus  $\pi < 1$ .

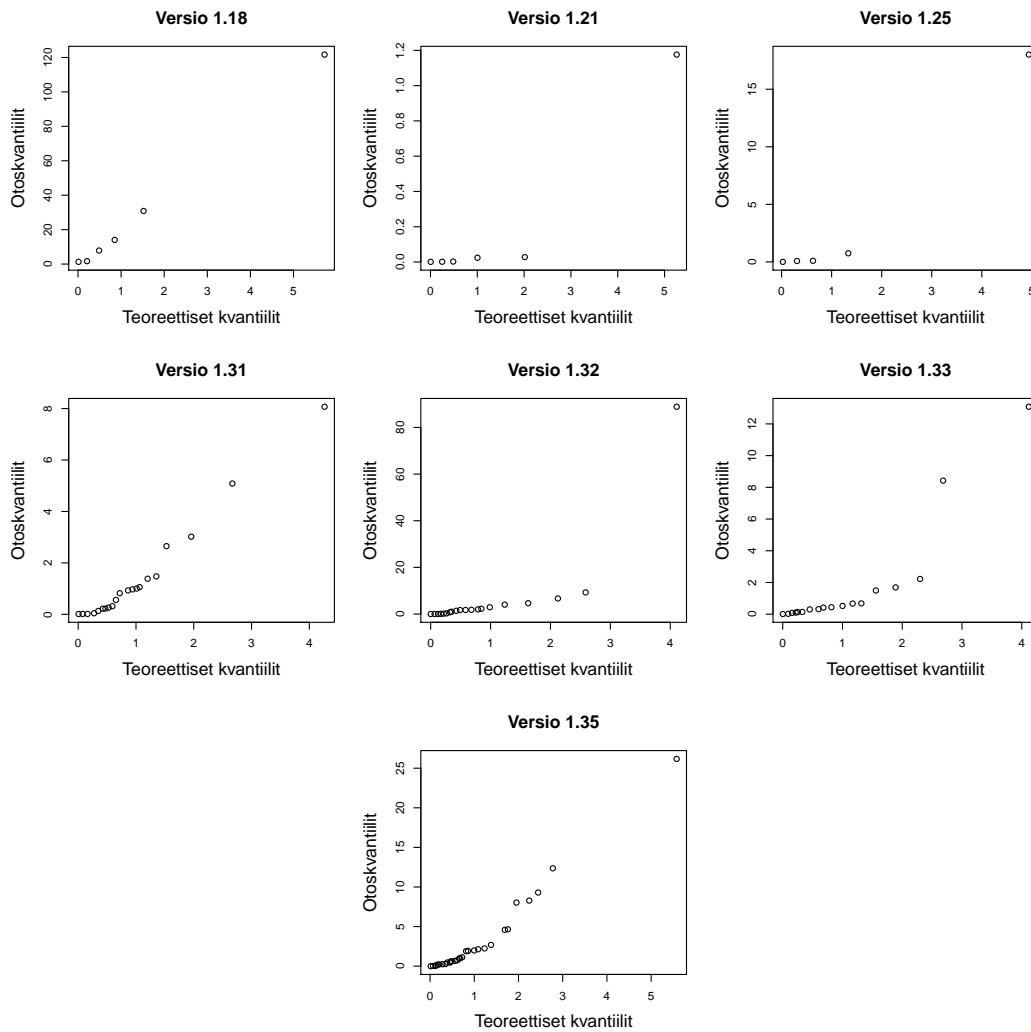
Dataan sovitettava funktio on nyt epätäydellisen datan uskottavuusfunktio (37), sillä täydellisen datan uskottavuusfunktion käyttämiseen tarvittaisiin tieto pelaajien alttiusindeksin  $\zeta$  arvoista. Lasketaan AIC-arvot jokaiselle tarkasteltelun kohteena olevalle versiolle erikseen, jotta tiedetään jokaisen osadatan noudattavan haluttua mallia. Taulukossa 8 esitetään lasketut arvot ja niistä nähdään, että jokaiselle versiolle pienemmän AIC-arvon tuottaa immuuneja pelaajia sisältävä malli; tämän kriteerin pohjalta mixture cure -mallin käyttäminen on siis järkevää.

Taulukko 8: Tarkasteltaville versioille lasketut AIC-arvot.

Versio	$\pi < 1$	$\pi = 1$
1.18	135.6777	154.0082
1.21	114.0396	132.1508
1.25	97.14472	110.3026
1.31	395.2878	444.7831
1.32	383.3518	417.9001
1.33	337.9899	364.7231
1.35	565.7867	604.2286

Tarkastellaan vielä mallin toista oletusta: havaittujen vikaantumisaikojen pitäisi noudattaa eksponenttijakaumaa. Verrataan kunkin version vikaantu-

misaikojen jakaumaa eksponenttijakauman mukaiseen dataan kuvan 13 kvantiilikuvaajien avulla. Kuvaajissa pitäisi näkyä pisteiden muodostama suora, jotta otos noudattaisi eksponenttijakaumaa. Kullakin versiolla on suoralta poikkeavia arvoja enintään kaksi, joten voidaan olettaa, että mallin toinenkin oletus toteutuu jokaiselle tarkasteltavalle versiolle. Eksponenttijakauma vaikuttaakin olevan vähintäänkin hyvä arvio vikaantumisaikojen jakaumasta, vaikka joillekin pelaajille konversio kestää kauan.



Kuva 13: Vikaantumisaikojen kvantiilikuvaajat.

### 5.3 Mallin käyttäminen

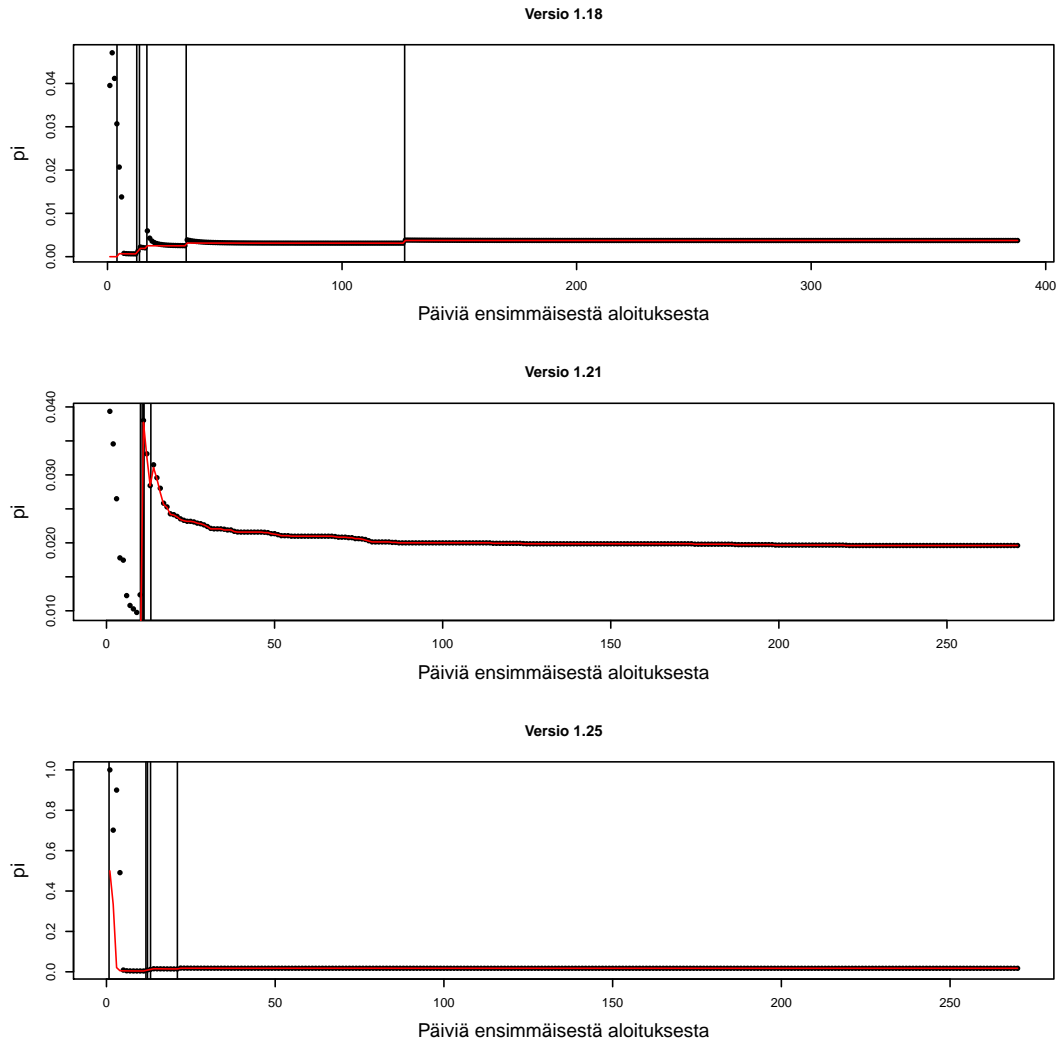
Mixture cure -mallin oletuksien toteutuminen on varmistettu, joten seuraavaksi voidaan varsinaisesti käyttää tätä mallia Hipster sheep -pelin datan aiemmin määriteltyihin osadatoihin. Selvitetään siis kyseisen pelin suosituimpien versioiden datoilta ostotapahtuman tekevien pelaajien osuutta kuvaavan parametrin  $\pi$  ja ostotahtia kuvaavan parametrin  $\lambda$  arvot. Sensurointiajankohdtaa muuttamalla voidaan selvittää, missä vaiheessa EM-algoritmilla lasketut parametrien arvot ovat samat kuin mitkä ne olisivat myöhemmillä sensurointiajankohdilla.

Määritetään kullekin tarkasteltavalle versiolle sensurointiajat päivän välein ensimmäisen pelaajan aloitusajasta viimeisen pelaajan lopettamisaikaan. Kullakin sensurointiajalla datasta piirrettävät kuvien 9 ja 10 mukaiset kuvat näyttäisivät hieman erilaiselta, mutta perusmuodoltaan ne kuitenkin olisivat kuvien 9 ja 10 kaltaiset. Sensurointiaikana voidaan käyttää todellista sensurointihetkeä eli havainnoinnin päättymishetkeä, jos halutaan verrata lyhemmillä tarkasteluväleillä laskettuja parametrien arvoja mahdollisimman pitkällä tarkasteluvälillä laskettuihin arvoihin. Tällä datalla se kuitenkin on turhaa, sillä kuvista 14 ja 15 nähdään, että jokaisella versiolla on jo käytettyjen sensurointiaikojen joukossa sellaisia ajankohtia, että monetisaatioprosentti  $\pi$  ei enää muutu.

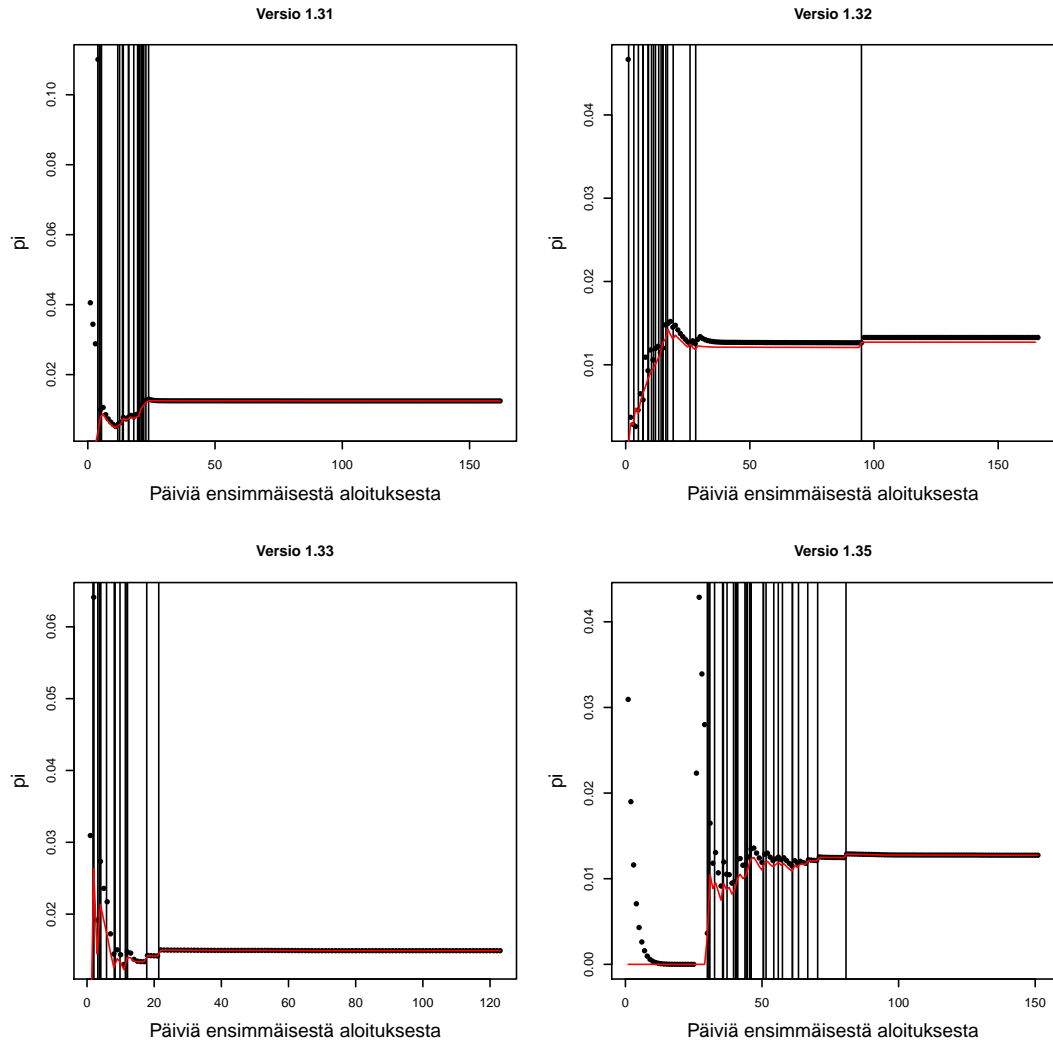
Kuvista 14 ja 15 nähdään, että suurin osa ostotapahtuman tekevästä pelaajista tekee sen ensimmäisen kuukauden kuluessa siitä, kun ensimmäinen pelaaja on kyseistä versiota alkanut pelata. Versio 1.35 poikkeaa muista versioista tässä asiassa, sillä kukaan sen pelaajista ei tee ostotapahtumaa ensimmäinen 30 päivän aikana, mutta sen jälkeen ostotapahtumia tehdään enemmän kuin missään muussa versiossa. Kuvasta 12 nähdään, että yksi aloitusaika on selvästi muita ennen, joten tämä poikkeus johtunee siitä. Kyseinen aloitusaika saattaa olla virheellinen, mutta asiaa ei saada varmistettua, joten annetaan kyseisen pelaajan olla datassa mukana.

Kuvista 14 ja 15 nähdään myös, että `cureEM`-funktiolla laskettu ostotapahtuman tekevien pelaajien osuus on melko sama kuin todellinen alttiiden yksilöiden osuus kullakin sensurointiajalla. Todellinen alttiiden yksilöiden osuus on kuitenkin jokaisella ajanhetkellä arvioitua osuutta pienempi. Mixture cure -malli huomioi myös mahdollisesti myöhemmin ostotapahtuman tekevät pelaajat, joten monetisaatioprosentin kuuluukin olla isompi kuin senhetkisestä datasta laskettu ostotapahtuman tehneiden pelaajien osuus.

Kuvasta 16 nähdään, että absoluuttisella sensurointiajalla laskettujen parametrien arvoilla piirretty välttöfunktion (36) kuvaaja vastaa erittäin hyvin osiossa 2.2 esitellyn epäparametrisen Kaplan–Meier-mallin kuvaajaa. Riittävän suurelle otokselle parametrisen mallin kuuluukin tuottaa samanlainen

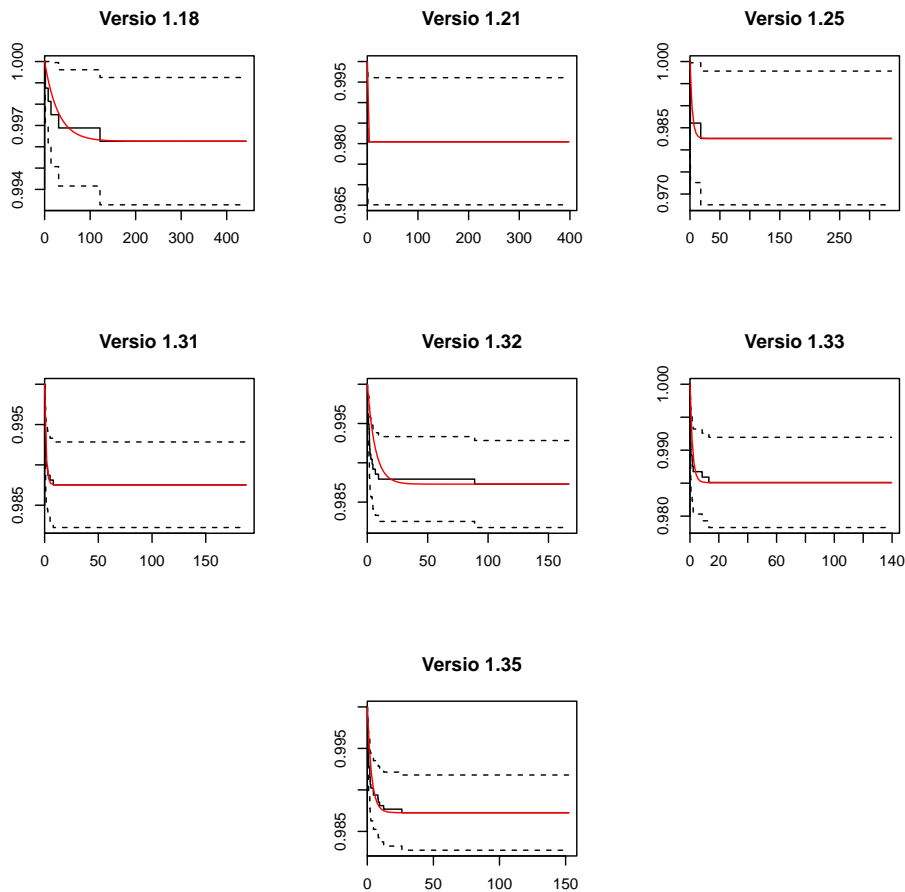


Kuva 14: Versioille 1.18, 1.21 ja 1.25 arvioitu monetisaatioprosentti päivän välein olevilla sensurointiajoilla. Pystyviivat kuvaavat ajanhetkiä, jolloin ostotapahtuman tehneiden pelaajien määrä on muuttunut. Punainen viiva kuvaa ostotapahtuman tehneiden pelaajien osuutta sen hetken kaikista pelaajista.



Kuva 15: Versioille 1.31, 1.32, 1.33 ja 1.35 arvioitu monetisaatioprosentti päivän välein olevilla sensurointiajoilla. Pystyviivat kuvaavat ajanhetkiä, jolloin ostotapahtuman tehneiden pelaajien määrä on muuttunut. Punainen viiva kuvaa ostotapahtuman tehneiden pelaajien osuutta sen hetken kaikista pelaajista.

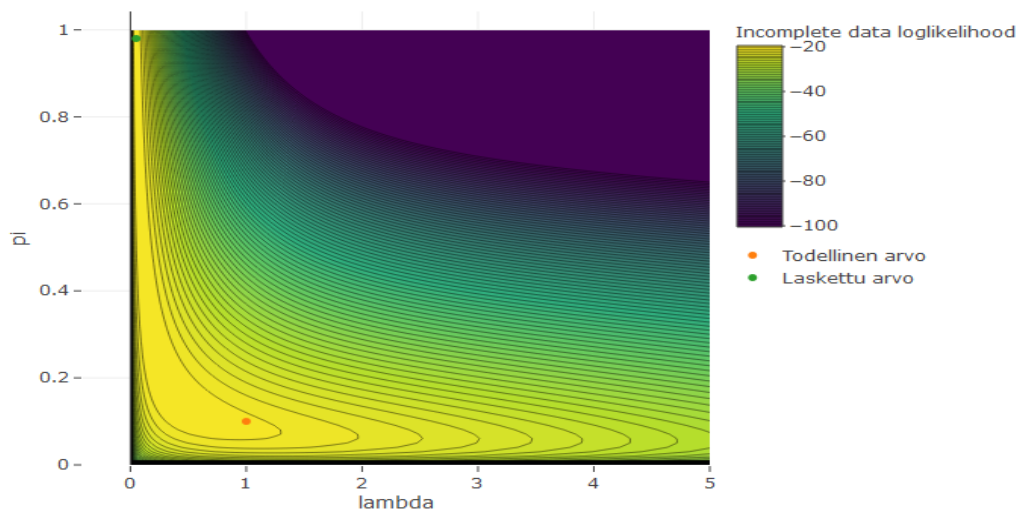
vastaus kuin epäparametrisen mallin, kunhan parametrisen mallin oletukset toteutuvat. Huomataan myös, että välttöfunktio laskee osalla versioista jyrkemmin kuin toisilla. Mitä nopeammin aloituksen jälkeen pelaajat tekevät ensimmäisen ostotapahtuman, sitä jyrkemmin välttöfunktion arvo laskee. Version 1.32 kuvaajasta nähdään, että punainen käyrä ei täysin mukaile epäparametrisen mallin käyrää. Tämä johtuu siitä, että version 1.32 datassa on yksi pelaaja, joka monetisoituu vasta pidemmän ajan kuluttua. EM-algoritmillä laskettu parametrin  $\lambda$  arvo lasketaan eräänlaisena keskiarvona, joten tämä myöhemmin monetisoituva pelaaja kasvattaa sen arvoa selvästi eikä saatu käyrä tästä syystä täysin vastaa epäparametrisen mallin käyrää.



Kuva 16: Kullekin versiolle piirretty välttöfunktion (36) kuvaaja absoluuttisella sensurointiajalla lasketuilla parametrien  $\pi$  ja  $\lambda$  arvoilla sekä Kaplan–Meier-mallin tuottamat välttöfunktion estimaatit 95 % luottamusväleineen.

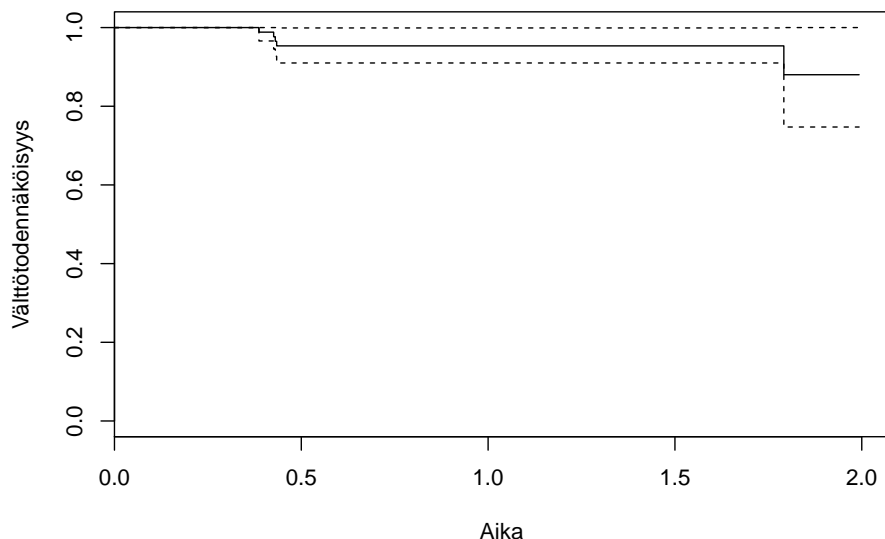
Joillain datajoukoilla tulee ongelmia funktion **cureEM** konvergenssin kanssa, jos käyttää hyvin pientä lopetuskriteeriä (esim.  $10^{-9}$ ) estimaattien laskemisessa. Malli ennustaa ensimmäisen ostotapahtuman jälkeen, että lähes kaikki pelaajat tekevät ostotapahtuman hitaalla konversiotahtilla. Tällaisessa tapauksessa ainoastaan muuttuja  $\lambda$  on mallia selittävä muuttuja eikä ole mahdollista erottaa elinaika-analyysin eksponenttijakauman mukaista mallia mixture cure -mallista.

Tämä ongelma ilmeni esimerkiksi version 1.15 datalla sensurointiajan ollessa 53 päivää ensimmäisen pelaajan aloituksesta. Ilmiö on kuitenkin selvemmin nähtävissä generoidulla datalla piirretyssä kuvaajassa 17. Kuvaaja on piirretty pienelle datalle ja ostotapahtumia on ollut vain muutamia. Kuvan 17 piirtämiseen käytetyn datan Kaplan–Meier-estimaatin kuvaaja esitetään kuvassa 18.



Kuva 17: Generoidulle testidatalle piirretty epätäydellisen datan logaritminen uskottavuusfunktion arvot muuttujien  $\lambda$  ja  $\pi$  funktiona. Uskottavuusfunktion arvoa vastaavat värit esitetään kuvan vieressä. Yhden alueen sisällä olevat uskottavuusfunktion arvot eroavat enintään yhdellä yksiköllä alueen suurimmasta arvosta.

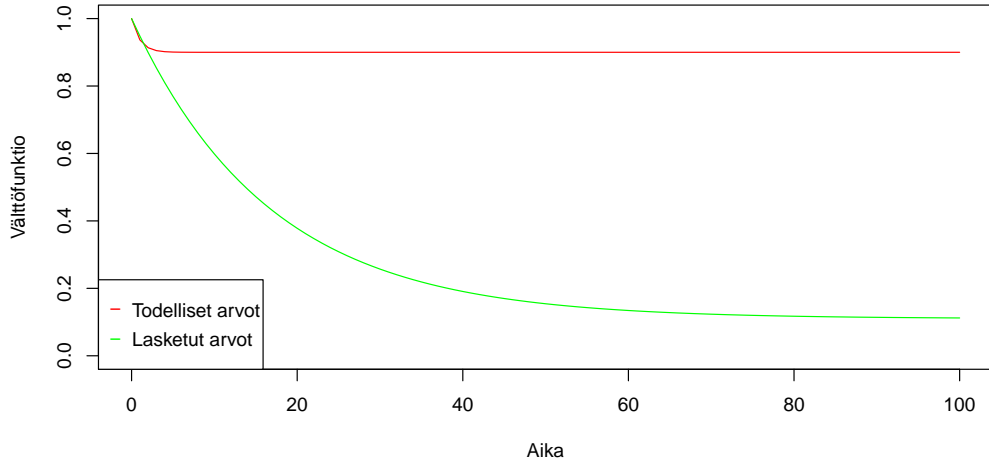
Huomataan, että kuvassa 17 esitetyt todellisten parametrien arvojen ja laskettujen parametrien arvojen pisteet ovat epätäydellisen datan logaritminen uskottavuusfunktion isoimpien arvojen kaistaleella. Kuvan 18 perusteella vaikuttaa siltä, että vain pieni osa yksilöistä monetisoi tuu kohtuullisella tahdilla  $\lambda$  tai sitten kaikki yksilöt monetisoi tuvat erittäin hitaalla tahdilla.



Kuva 18: Kuvaajan 17 piirtämiseen käytetyn generoidun datan Kaplan-Meier-estimaatin kuvaaja.

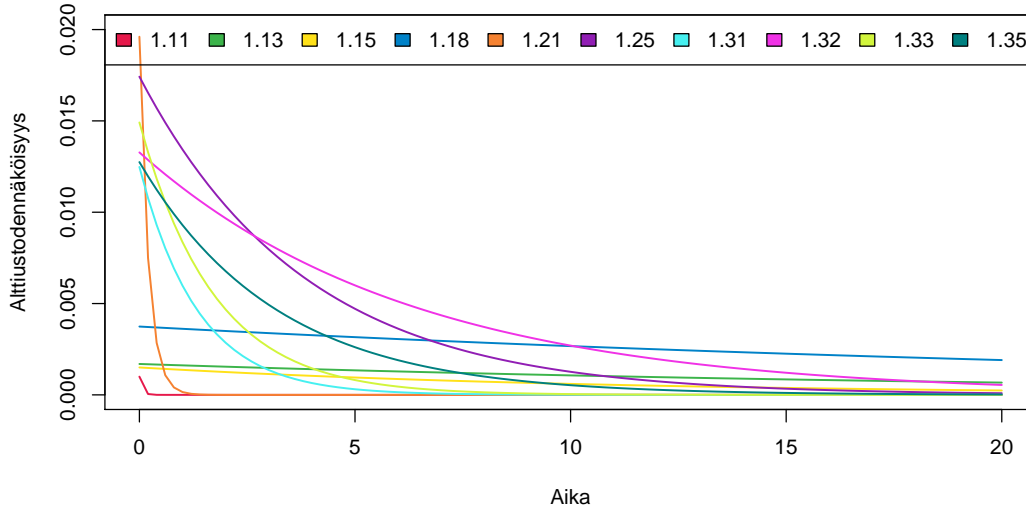
EM-algoritmin konvergenssin hitaus tälle datalle johtuu juuri uskottavien parametrien arvojen alueen laajuudesta. Datasta ei näin lyhyellä seuranta-ajalla pysty erottamaan mixture cure -mallia tavallisesta eksponenttijakoumasta. Kuvasta 19 nähdään, että lyhyellä seuranta-ajalla välttöfunktioit lasketuilla ja todellisilla parametrien arvoilla saavat samanlaisia arvoja, mutta mitä pidempi seuranta-aika on, sitä isompi ero välttöfunktioiden arvoissa on.

Myös ensimmäisen ostotapahtuman tekevän pelaajan peliaika ennen ostotapahtuman tekemistä ilmeisesti vaikuttaa siten, että malli ennustaa lähes kaikkien pelaajien tekevän ostotapahtuman, jos ensimmäinen ostotapahtuma tehdään vain hieman ennen sensurointiajankohtaa. Mitä enemmän aikaa kuluu ostotapahtuman tekemisestä sensurointiajankohtaan, sitä paremmin malli ennustaa mallissa olevan immuunien pelaajien osuuden. Kaikilla versioittaisilla datoilla ei ole ongelmaa konvergenssin hitauden kanssa. Riippuu myös versiosta, että paljonko sensurointiajankohtaa pitää siirtää myöhemmäksi, jotta malli tunnistaa immuunien pelaajien osuuden olemassaolon. Suurimman uskottavuuden estimaattien konvergenssin perusteella kuitenkin on olemassa sellainen ajanhetki, että malli tunnistaa immuunien yksilöiden olemassaolon, vaikka se tapahtuisikin hitaasti.



Kuva 19: Välttöfunktion (36) arvot todellisilla parametrien arvoilla  $\pi_g = 0.1$  ja  $\lambda_g = 1$  ja funktiolla `cureEM` lasketuilla parametrien arvoilla  $\lambda \approx 0.06$  ja  $\pi \approx 0.89$ .

Osiassa 3 mainittiin, että EM-algoritmilla saadaan laskettua jonkinlaiset uskottavat arvot puuttuvalle datalle. Tässä mallissa puuttuvaa dataa on alttiusindeksi  $\zeta$ . Osiossa 4.1 esiteltiin kaava (40), jolla lasketaan todennäköisyys, että yksilö on altis tarkastelun kohteena olevalle tapahtumalle, vaikka se ei vielä ole sitä tehnyt. Nämä todennäköisyydet ovat luetteloinnissa 1 riveillä 18 laskettavat kertoimet. Sensuroituille pelaajille ne muuttuvat jokaisella laskukierroksella ja sensuroimattomille yksilöille se on aina yksi, sillä sensuroimaton pelaaja on jo tehnyt ostotapahtuman. Kuvassa 20 on kullekin versiolle piirretty alttiustodennäköisyys ajanfunktiona. Huomataan, että aluksi todennäköisyydet vaihtelevat melko paljon versioittain. Ajan edetessä se kuitenkin lähetsyy nollaa jokaisella versioilla. Tästäkin kuvasta nähdään, että versioiden 1.11, 1.13 ja 1.15 tarkastelu ei edes olisi ollut järkevää, sillä jo heti aluksi alttiustodennäköisyys on hyvin pieni.



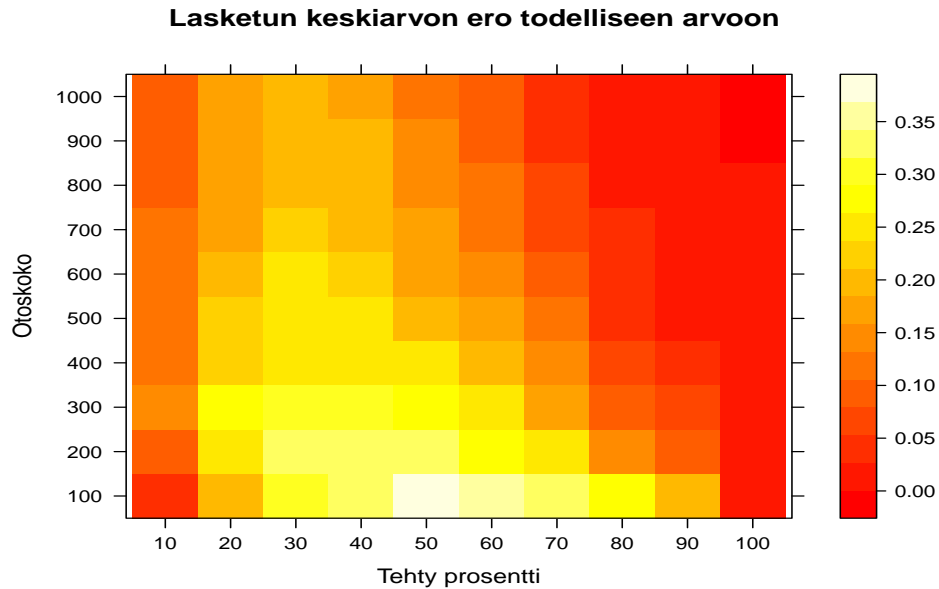
Kuva 20: Funktio (40) piirrettynä kullekin versiolle ajan suhteen arvolla  $j = 1$  ja absoluuttisella sensurointiajalla lasketuilla parametrien  $\lambda$  ja  $\pi$  optimaalisilla arvoilla.

## 5.4 Mallin muuttujien vaikutus laskettuihin estimaatteihin

Keskitytään lopuksi vielä mallin toimintaan vaikuttavien muuttujien vaikutuksen voimakkuuden tarkasteluun simuloimalla mallinmukaisia tilanteita. Tutkitaan otoskoon  $n$ , sensurointiajan määrittävän ostotapahtumien prosentin  $p$  ja datan generoinnissa käytettävän alttiustodennäköisyyden  $\pi_g$  vaikutusta EM-algoritmilla laskettavaan arvoon  $\pi$ .

Tarkastellaan ensin otoskoon  $n$  ja seuranta-ajan vaikutusta monetisaatio-prosentin  $\pi$  estimaattiin. Seuranta-aika määritetään siten, että sensurointi tapahtuu siinä vaiheessa, kun  $p$  prosenttia alttiista yksilöistä on monetisoitunut. Generoidaan 1000 mallinmukaista otosta arvolla  $\pi_g = 0.05$  kullekin  $(n, p)$ -parille. Lasketaan kullekin otokselle parametrien  $\pi$  ja  $\lambda$  arvot funktiolla `cureEM`. Otetaan keskiarvo saaduista parametrin  $\pi$  arvoista ja lasketaan sille 95 %:n luottamusväli. Datan generoinnin yhteydessä määritetään, että  $\pi = 0$  ja  $\lambda = 0$  niille otoksille, joissa sensuroimattomia yksilöitä ei ole ollenkaan. Lopulliset parametrien arvot riippuvat alkuarvauksesta, jos kukaan ei ole tapahtumaa tehnyt, joten määrittämällä niille arvo 0, pyritään välttämään parametrien jakauman vinoutuminen. Saatuja tuloksia havain-

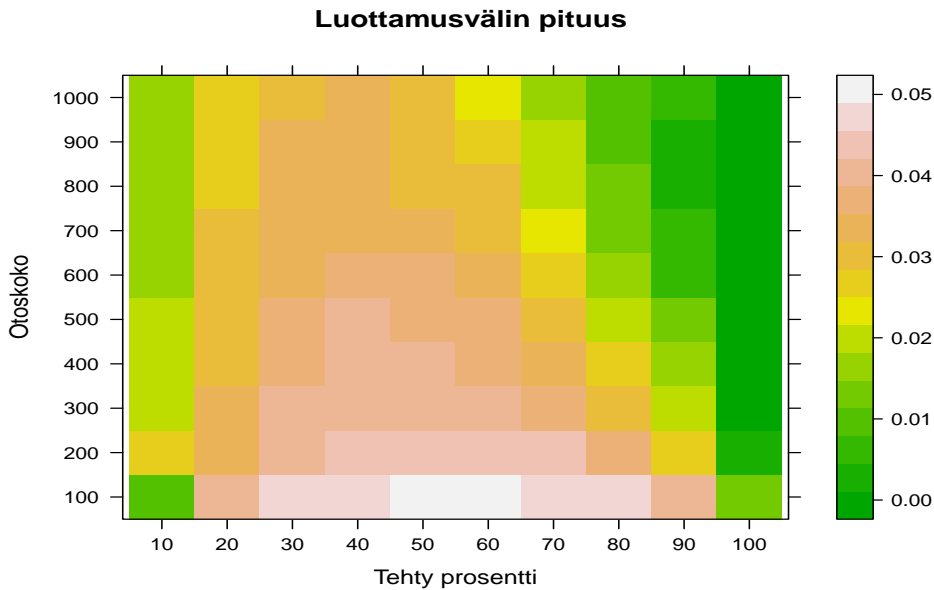
nollistetaan kuvilla 21 ja 22.



Kuva 21: Funktiolla `cureEM` laskettujen parametrin  $\pi$  estimaattien keskiarvojen ero todelliseen arvoon  $\pi_g = 0.05$ .

Kuvassa 21 havainnollistetaan laskettujen estimaattien  $\pi$  keskiarvon eroa datan generoinnissa käytettyyn alttiiden yksilöisen osuuteen  $\pi_g = 0.05$ . Huomataan, että isoin ero on arvoilla  $n = 100$  ja  $p = 50\%$ . Pienemmillä prosentteilla  $p$  ero on pienempi todennäköisesti siitä johtuen, että niissä tapauksissa on useita otoksia, joissa ei ollenkaan ole sensuroituja yksilöitä. Näin pienellä otoskokoalla todennäköisesti suurimmalle osalle otoksia tehtyjen tapahtumien määrä näillä sensurointiajoilla on nolla, joten se tietysti pienentää keskiarvoa huomattavasti. Kuvasta 21 nähdään, että mitä isompi osa alttiista yksilöistä on jo monetisoitunut, sitä lähempänä lasketut arvot ovat todellista arvoa.

Näin pienellä alttiiden yksilöiden määrällä  $\pi_g$  ei kovinkaan monella 1000 otoksen simulaatiolla saatu sellaista arviota parametrille  $\pi$ , että ero todelliseen arvoon  $\pi_g$  olisi pienempi kuin todellisen arvon suuruus 0.05. Otoskoko  $n$  kasvattamalla tällaisia arvoja saadaan useammin. Tällaisia arvoja ei saada kuin arvolla  $p = 100\%$ , jos otoskoko  $n \leq 400$ . Kun  $n \geq 600$ , tällaisia arvoja saadaan jo siinä vaiheessa, kun  $p = 80\%$ . Otoskoko täytyy vielä kasvattaa, jotta saadaan lähes oikeansuuruinen estimaatti, pienemmällä tehtyjen tapahtumien prosentilla  $p$ . Kuitenkaan seuranta-aikaa ei voida määrittää arvoilla  $p \leq 70\%$ , vaikka otoskoko olisi 1000 yksilöä, jos saatu keskiarvo halutaan melko lähelle todellista arvoa. On kuitenkin mahdollista, että vielä isommal-

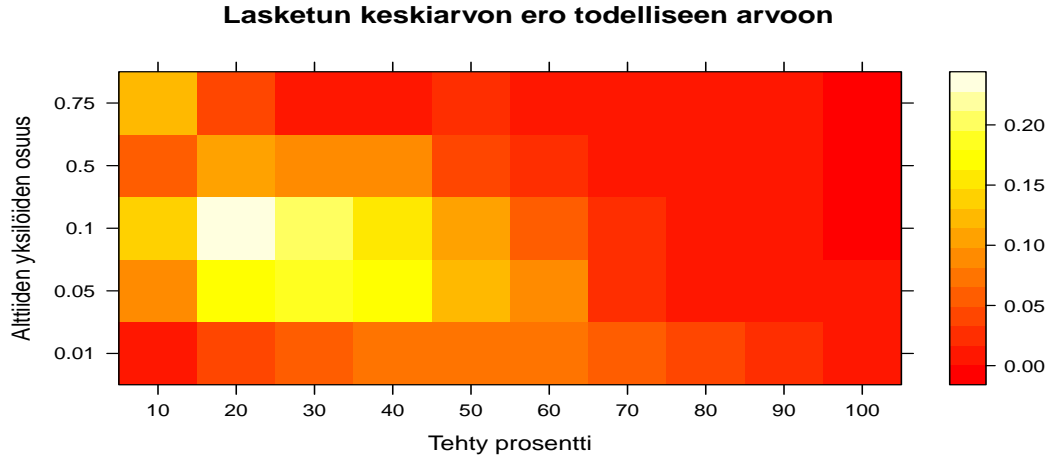


Kuva 22: Funktiolla `cureEM` laskettujen parametrin  $\pi$  estimaattien 95 %:n luottamusvälin pituus.

la otoskoolla ennustamiseen tarvittavaa tehtyjen tapahtumien prosenttia  $p$  saadaan vielä pienemmäksi.

Kuvasta 22 nähdään, että simuloinnilla laskettujen parametrin  $\pi$  estimaattien 95 %:n luottamusvälin pituus vaihtelee jonkin verran. Pienellä otoskoolla  $n$  se on isompi kuin isolla otoskoolla. Jälleen huomataan, että myös sensurointiajan määrittämiseen vaikuttava monetisoituneiden alttiiden yksilöiden osuus  $p$  vaikuttaa luottamusvälin pituuteen. Nyt kuitenkin sen vaikutus nähdään selkeämmin vasta isommilla muuttujien  $n$  ja  $p$  arvoilla.

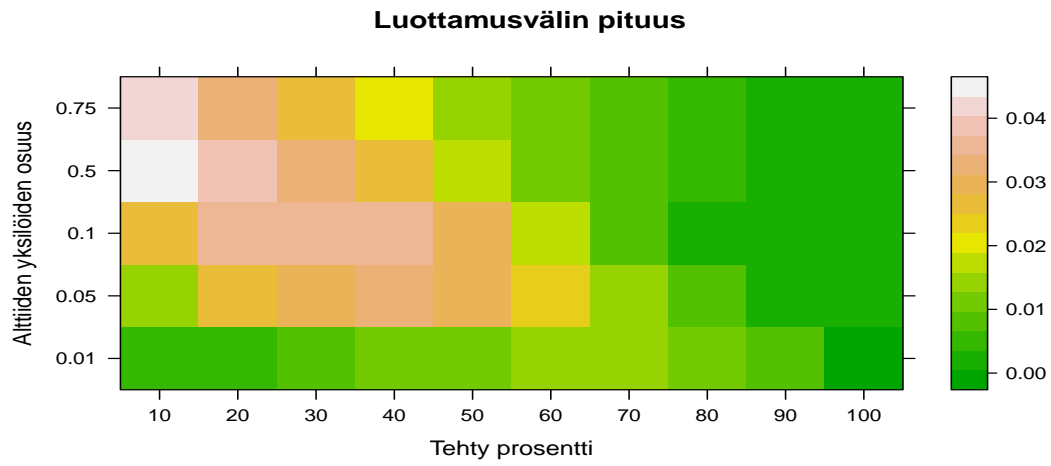
Tarkastellaan seuraavaksi parametrin  $\pi_g$  arvon vaikutusta lasketun monetisaatioprosentin  $\pi$  arvoon. Käytetään datan generoinnissa otoskokoa  $n = 1000$ . Generoidaan tuhat otosta kullekin  $(\pi_g, p)$ -parille ja lasketaan `cureEM`-funktiolla laskettujen parametrin  $\pi$  estimaattien keskiarvo ja 95 %:n luottamusväli. Tällä kertaa simuloinnissa muuttujina ovat datan generoinnissa käytettävä alttiusprosentti  $\pi_g$  ja sensurointiajan määrittävä tapahtumien osuus  $p$  kaikkien alttiiden yksilöiden määrästä. Generoidaan dataa arvoilla  $\pi_g = 0.01, 0.05, 0.1, 0.5$  ja  $0.75$ . Sensurointiajan määrittävä tehtyjen tapahtumien prosentti  $p$  saa samat arvot kuin otoskoon vaikutusta selvitetessä. Kuvissa 23 ja 24 havainnollistetaan simuloinnin tuloksia esittämällä jälleen laskettujen parametrin  $\pi$  estimaattien keskiarvon ero todelliseen arvoon  $\pi_g$  ja 95 %:n luottamusvälit.



Kuva 23: Parametrin  $\pi$  estimaattien keskiarvon ero todellisesta alttiiden yksilöiden osuudesta  $\pi_g$ .

Kuvasta 23 nähdään, että mitä suurempi osa yksilöistä on alttiita, sitä pienempi osa tapahtumista täytyy tapahtua, jotta pystytään ennustamaan lopullinen muuttujan  $\pi$  arvo lähes oikein. Tästä kuvasta nähdään sama kuin kuvasta 21; alttiusprosentin arvolla  $\pi_g = 0.05$  sensurointiaika täytyy määrittää siten, että 70 % alttiista yksilöistä on monetisoitunut. Arvolla  $\pi_g = 0.01$  saadut tulokset ovat hyviä, mutta se todennäköisesti on seurausta sensuroimattomien yksilöiden puutteesta. Arvolla  $\pi_g = 0.75$  ei tarvitse seurata yksilöitä kuin siihen asti, että 30 % alttiista yksilöistä tekee tapahtuman. Jo siinä vaiheessa otoksilla laskettujen arvojen  $\pi$  keskiarvo on lähellä datan generoinnissa käytettyä arvoa  $\pi_g$ . Huomataan, että yli puolet otoksesta täytyy olla alttiita yksilöitä, jotta sensurointiajankohta voidaan määrittää osuutta 50 % pienemmällä arvolla  $p$  ja tulos saadaan melko tarkasti ennustettua. Kuva 24 on melko samanlainen kuvioltaan kuin kuva 23, mikä tarkoittaa, että isommilla muuttujien  $\pi_g$  ja  $p$  arvoilla 1 000 otoksen estimaatti monetisatioprosentista  $\pi$  ei vaihtele niin paljoa kuin pienemmillä muuttujien  $\pi_g$  ja  $p$  arvoilla.

Huomataan siis, että alle tuhannen yksilön otoksilla parametrin  $\pi$  suurimman uskottavuuden estimaatti  $\hat{\pi}$  ei konvergoi todelliseen arvoon  $\pi_g$  ennen kuin lähes kaikki alttiit yksilöt ovat tehneet tarkastelun kohteena olevan tapahtuman. Alttiiden yksilöiden osuutta kasvattamalla, suurimman uskottavuuden estimaatti  $\hat{\pi}$  konvergoi paremmin todelliseen arvoon  $\pi_g$  myös ly-



Kuva 24: Parametrin  $\pi$  estimaattien 95 % luottamusväli.

hyemmällä seuranta-ajalla.

## 6 Yhteenveto

Tässä tutkielmassa johdettiin EM-algoritmin vaatimat funktiot mixture cure -mallille. Mallissa osa yksilöistä on immuuneja tarkastelun kohteena olevalle tapahtumalle ja siinä käytetään oikealta sensurointia. Pääasiassa tässä tutkielmassa käytettiin tyyppin I sensurointia, mutta vastaavalla tavalla mallin pitäisi toimia tyyppin II tai satunnaiselle sensuroinnille. Mallin konvergenssi muodostuu kahdesta osasta: EM-algoritmi konvergoi suurimman uskottavuuden estimaattiin ja suurimman uskottavuuden estimaatti konvergoi parametrien todellisiin arvoihin. Mallin johtamisen jälkeen tarkasteltiin sen soveltuvuutta pelistä kerättyyn dataan. Soveltamisosion lopuksi tutkittiin datan tiettyjen piirteiden vaikutusta siihen, että suurimman uskottavuuden estimaatti  $\hat{\pi}$  konvergoi todelliseen parametrin  $\pi$  arvoon.

Johdettua mallia sovellettiin vain yhdestä pelistä kerättyyn dataan. Ainakin kyseinen peli toteutti mallin oletukset, joten mahdollisesti tätä mallia voidaan soveltaa muistakin peleistä kerättyyn dataan. Tässä tutkielmassa keskityttiin enimmäkseen vain monetisaatioprosentin  $\pi$  ennustamiseen, mutta vastaavalla tavalla voitaisiin tutkia konversiotahdia  $\lambda$ .

Huomattiin, että malli ennustaa Hipster sheep -pelin datalle melko oikean arvon muuttujalle  $\pi$  jo siinä vaiheessa, kun vasta osa ostotapahtumista on tehty. Vaikuttaa siis siltä, että malli toimii tässä tutkielmassa käytetyn datan mukaiselle datalle. Ennustettu parametrin  $\pi$  arvo on vain hieman jo havaittua ostotapahtuman tekevien pelaajien osuutta suurempi, mikä tarkoittaa, että sensuroitujen pelaajien todennäköisyys ostotapahtuman tekemiseen on melko pieni. Tämä todennäköisyys pienenee hyvin nopeasti seuranta-ajan pidentyessä.

Mallin konvergenssiin vaikuttavia tekijöitä ovat otoskoko, sensurointiaikaan mennessä tehtyjen tapahtumien määrä sekä alttiiden yksilöiden osuus koko populaatiosta. EM-algoritmi laskee suuren määrän iteraatiokierroksia ennen konvergenssia, jos ostotapahtumia on vähän. Soveltamisosion lopussa tehdyistä testeistä huomattiin, että mitä isompi otos, sitä paremmin EM-algoritmillä laskettu suurimman uskottavuuden estimaatti vastaa todellisia arvoja. Huomattiin myös, että sitä paremmin suurimman uskottavuuden estimaatti vastaa todellisia arvoja, mitä isompi osa populaation alttiista yksilöistä on tehnyt tapahtuman ennen sensurointia. Myös alttiiden yksilöiden osuus kaikista yksilöistä vaikuttaa siten, että mitä suurempi todellinen osuus on, sitä lyhyemmällä seuranta-ajalla tulos pystytään ennustamaan. Malli siis ennustaa sitä paremmin, mitä enemmän monetisoituneita pelaajia on.

Mallin huonona puolena on, että liian pienellä datalla ja liian vähäisellä tapahtumien määrällä menetelmä ei tunnista eroa tämän mallin ja tavallisen elinaika-analyysin mallin välillä. Tällaisessa tilanteessa parametrien esti-

maatteihin liittyy suurta tilastollista vaihtelua ja malli saattaa ennustaa jopa kaikkien pelaajien tekevän ostotapahtuman. Siispä tässä tutkielmassa käytetyn mallin käyttäminen ei ole järkevää, jos otoskoko on liian pieni, vain harva alttiista pelaajista on monetisoitunut sensurointiajankohtaan mennessä tai data on sellaista, että alttiita yksilöitä ei ylipäätään ole kovin monta.

Mallia pystytään kehittämään pienemmälle datalle sopivammaksi lisäämällä siihen arvio parametrien arvoista ennen datan näkemistä. Ennen datan näkemistä muun tiedon perusteella arvioidaan parametrien arvojen prioritiheysfunktioit, jotka myöhemmin päivitetään posterioritiheysfunktioiksi datan perusteella.

## Kiitokset

Haluan kiittää ohjaajiani Markus Viljasta ja Marko Mäkelää. Erityisesti kiitos Markus Viljaselle mielenkiintoisesta aiheesta tälle tutkielmalle sekä ajatuksista mallin kehittämiseksi. Kiitos myös ohjeista todellisen datan saamiseen sekä R-koodista, jolla sain havainnollistettua dataa kuvilla 9 ja 10.

## Kirjallisuutta

- [1] Tribeflame hipster sheep. <http://www.tribeflame.com/games.html>. Viitattu: 29.4.2018.
- [2] D. R. Cox and D. Oakes. *Analysis of survival data*. Chapman and Hall Ltd, 1984.
- [3] R. G. Miller Jr. *Survival analysis*, volume 66. John Wiley & Sons, 1981.
- [4] D. F. Moore. *Applied Survival Analysis Using R*. Springer International Publisher Switzerland, 2016.
- [5] G. McLachlan and T. Krishnan. *The EM algorithm and extensions*. John Wiley & Sons, New York, 1997.
- [6] X. Liu. *Survival Analysis: Models and Applications*. Higher Education Press, 2012.
- [7] E. L. Kaplan and P. Meier. Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, 53(282): 457–481, 1958. URL <http://www.jstor.org/stable/2281868>.
- [8] J.P. Klein, H.C. van Houwelingen, J.G. Ibrahim, and T.H. Scheike. *Handbook of Survival Analysis*. Chapman & Hall/CRC Handbooks of Modern Statistical Methods. CRC Press, 2016. URL <https://books.google.fi/books?id=t1v0BQAAQBAJ>.
- [9] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977. URL <http://www.jstor.org/stable/2984875>.
- [10] C. F. J. Wu. On the convergence properties of the EM algorithm. *The Annals of statistics*, pages 95–103, 1983.
- [11] M. Mäkelä. *Optimointialgoritmit*. Turun yliopisto, luentomoniste, Turku, 2016, . URL [https://www.utu.fi/fi/yksikot/sci/yksikot/mattil/opiskelu/kurssit/Documents/moniste\\_2016.pdf](https://www.utu.fi/fi/yksikot/sci/yksikot/mattil/opiskelu/kurssit/Documents/moniste_2016.pdf). Viitattu 27.4.2018.
- [12] M. Mäkelä. *Matemaattinen optimointi I*. Turun yliopisto, luentomoniste, Turku, 2015, . URL <https://www.utu.fi/fi/yksikot/sci/yksikot/mattil/opiskelu/2013-14-kurssit/Documents/MonisteMatOpt1.pdf>. Viitattu 27.4.2018.

- [13] A. Wang, Y. Zhang, and Y. Shao. On the likelihood of mixture cure models. *Statistics & Probability Letters*, 131:51–55, 2017.
- [14] H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974. doi: 10.1109/TAC.1974.1100705.