

Sakke Purolainen

Korrespondenssianalyysi aineiston analysoinnin tukena

Informaatioteknologian ja viestinnän tiedekunta
Kandidaattitutkielma
Toukokuu 2022

Tiivistelmä

Sakke Purolainen: Korrespondenssianalyysi

Kandidaattitutkielma

Tampereen yliopisto

Matematiikan ja tilastollisen data-analyysin kandidaattiohjelma

Toukokuu 2022

Korrespondenssianalyysi on monimuuttujamenetelmä, jota käytetään mallintamaan monimutkaista kategorista dataa. Tutkielman tarkoituksena on perehtyä korrespondenssianalyysin ominaisuuksiin, sekä havainnollistaa menetelmää käytännön tilanteissa.

Tutkielmassa esitellään ensin korrespondenssianalyysin idea. Sen jälkeen tarkastellaan menetelmän matemaattista taustaa vektoreiden ja matriisien avulla. Tarkastelussa keskitytään etenkin rivi- ja sarakeprofileihin, riippumattomuuden testaamiseen ja kuvan muodostamiseen. Lopuksi sovelletaan sitä kahteen eri kyselytutkimusaineistoon R-ohjelmistoa apuna käyttäen sekä tehdään sen pohjalta tulkintoja. Tutkielmassa esitellään myös lyhyesti imputaatiomenetelmä, jolla pyritään paikkaamaan aineiston puuttuvia arvoja. Imputaatiota hyödynnetään ensimmäiseen kyselytutkimusaineistoon. Tutkielma osoittaa, että korrespondenssianalyysi on erittäin hyvä menetelmä mallintamaan kategorista dataa. Visuaalisuutensa ansiosta analyysin tulkitseminen on myös melko helppoa. Lisäksi tutkielmassa huomataan aineiston riittävyden merkitys: menetelmän toimivuuden kannalta on tärkeää, että aineisto sisältää runsaasti havaintoja erityisesti silloin, jos kategorioita on useita.

Tämän julkaisun alkuperäisyys on tarkastettu Turnitin OriginalityCheck -ohjelmalla.

Sisällys

1	Johdanto	4
2	Menetelmä	5
2.1	Rivi- ja sarakeprofiilit	5
2.2	Riippumattomuuden testaaminen	8
2.3	Kuvan muodostaminen	10
3	Soveltaminen aineistoihin	13
3.1	Aineistot	13
3.2	Datan muokkaus	13
3.3	Puuttuvan tiedon paikkaaminen (Imputaatio)	14
3.4	Korrespondenssianalyysi pankkidataan	15
3.5	Korrespondenssianalyysi uraseurantadataan	16
4	Johtopäätökset	21
	Lähteet	22
A	Kuvaajien R-koodi	23

1 Johdanto

Korrespondenssianalyysistä ei juuri suomenkielistä materiaalia internetistä löydy. Englannin kielellä materiaalia on jonkin verran, mutta monet englanninkielisistä artikkeleistakin ovat hyvin pintapuolisia. Vaikka korrespondenssianalyysi on hyvin lähellä pääkomponenttianalyysia, ei se jostain syystä ole läheskään yhtä yleisesti käytetty.

Le Roux ja Rouanet (2010) mainitsevat, että korrespondenssianalyysia on kehitetty usean henkilön toimesta, mukaan lukien R. A. Fisherin, joka oli yksi merkittävimmistä henkilöistä tilastotieteessä 1900-luvun alkupuolella. Ranskalainen tilastotieteilijä Jean-Paul Benzecri kuitenkin kehitti ja popularisoi menetelmän vasta 1960- ja 1970-luvuilla, ensin Ranskassa ja myöhemmin muualla Euroopassa. Hän myös antoi menetelmälle nykyään yleisesti käytetyn nimen *analyse de correspondance* (Le Roux & Rouanet, 2010.)

Korrespondenssianalyysi on exploratiivinen menetelmä monimutkaiselle kategoriselle datalle (Glynn & Ribbson, 2014). Sillä pyritään esittämään tietoa kaksisuuntaisella ristitaulukolla. Taulukko sisältää kahden tai useamman kategorisen muuttujan frekvenssit. Korrespondenssianalyysilla saadaan rakennettua malli, joka kuvaa kahden muuttujan välistä vuorovaikutusta sekä rivien ja sarakkeiden yhteyttä toisiinsa (Rencher, 2002.)

Työn tavoite on esitellä korrespondenssianalyysin yleinen periaate sekä menetelmän matemaattinen tausta. Lopuksi menetelmää käytetään kahteen dataan ja tehdään niistä tulkinnat analyysin pohjalta.

2 Menetelmä

Tavanomaisen kaksisuuntaisen korrespondenssianalyysin lisäksi on olemassa korrespondenssianalyysin yleistys moniulotteisiin tauluihin (Rencher, 2002). Käytännön esimerkit käsittelevät tässä tutkielmassa molempia tapauksia. Testatessa merkitsevyyttä voidaan käyttää joko χ^2 -testiä tai lineaarista log-lineaarista mallia (Rencher, 2002). Molempien käyttö perustuu asymptoottisiin tuloksiin (Rencher, 2002).

Mikäli ristitaulukossa jokin solun frekvensseistä on hyvin pieni tai nolla, χ^2 -approksimaatio ei ole kovinkaan tyydyttävä. Tässä tapauksessa joitakin muuttujia voi yhdistää keskenään, jotta saadaan nostettua solun frekvenssiä. Korrespondenssianalyysi voikin olla hyödyllinen työkalu tunnistamaan muuttujia, jotka ovat samankaltaisia, ja tästä syystä voitaisiin mahdollisesti yhdistää (Rencher, 2002.)

Menetelmä esittää tulokset kaksiulotteisena kuvaajana, joka visualisoi muuttujien yhteyttä intuitiivisesti (Glynn & Ribbson, 2014). Mikäli kahden rivin pisteet ovat lähellä toisiaan, niiden profiilit ovat rivien suhteen samankaltaisia. Yhtä lailla, mikäli sarakkeiden pisteet ovat lähellä toisiaan, ovat niiden profiilit samankaltaisia rivien suhteen (Rencher, 2002.) Menetelmän vahvuus onkin juuri visuaalisessa esityksessä.

2.1 Rivi- ja sarakeprofiilit

Rencher (2002) määrittelee rivi- ja sarakeprofiilit kirjassaan seuraavasti: Alla on ristitaulukko 2.1, jossa on a riviä ja b saraketta. Jokainen taulukon alkio n kuvaa frekvenssejä kaikille rivien ja sarakkeiden kombinaatioille. Reunajakaumien summat on esitetty summanotaatiolla: $n_{i.} = \sum_{j=1}^b n_{ij}$ ja $n_{.j} = \sum_{i=1}^a n_{ij}$. Kaikkien summa on merkitty yksinkertaisesti n (Rencher, 2002.)

Ristitaulukon frekvenssit n_{ij} voidaan muuttaa suhteellisiksi frekvensseiksi p_{ij} yksinkertaisesti jakamalla n : $p_{ij} = n_{ij}/n$. Suhteellisten frekvenssien matriisia kutsutaan korrespondenssimatriisiksi ja sitä merkitään kirjaimella \mathbf{P} (Rencher, 2002.)

Taulukko 2.1. Suhteellisten frekvenssien ristitaulukko, jossa on a riviä ja b saraketta

	1	2	...	b	Summa
1	p_{11}	p_{12}	\cdots	p_{1b}	$p_{1.}$
2	p_{21}	p_{22}	\cdots	p_{2b}	$p_{2.}$
\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
a	p_{a1}	p_{a2}	\cdots	p_{ab}	$p_{a.}$
Summa	$p_{.1}$	$p_{.2}$	\cdots	$p_{.b}$	1

Viimeinen sarake taulussa 2.1 sisältää rivisummat $p_{i.} = \sum_{j=1}^b p_{ij}$. Merkitään sarakevektoria vektorilla \mathbf{r} . \mathbf{r} saadaan

$$(2.1) \quad \mathbf{r} = \mathbf{P}\mathbf{j} = (p_{1.}, p_{2.}, \dots, p_{a.})' = (n_{1.}/n, n_{2.}/n, \dots, n_{a.}/n)',$$

missä \mathbf{j} on $a \times 1$ vektori numero ykkösiä. Samaan tapaan viimeinen rivi taulukossa 2.1 sisältää sarakesummat $p_{.j} = \sum_{i=1}^a p_{ij}$. Merkitään sarakevektoria vektorilla \mathbf{c} ja se saadaan

$$(2.2) \quad \mathbf{c}' = \mathbf{j}'\mathbf{P} = (p_{.1}, p_{.2}, \dots, p_{.b}) = (n_{.1}/n, n_{.2}/n, \dots, n_{.b}/n),$$

missä \mathbf{j}' on $1 \times a$ vektori numero ykkösiä. Korrespondenssimatriisi ja reunajakaumien summat taulukossa 2.1 voidaan esittää

$$\begin{pmatrix} \mathbf{P} & \mathbf{r} \\ \mathbf{c}' & 1 \end{pmatrix} = \left(\begin{array}{cccc|c} p_{11} & p_{12} & \cdots & p_{1b} & p_{1.} \\ p_{21} & p_{22} & \cdots & p_{2b} & p_{2.} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ p_{a1} & p_{a2} & \cdots & p_{ab} & p_{a.} \\ \hline p_{.1} & p_{.2} & \cdots & p_{.b} & 1 \end{array} \right)$$

Seuraavaksi jokainen \mathbf{P} :n rivi ja sarake muutetaan profiiliksi. Tällöin i :nnen rivin profiili \mathbf{r}'_i , $i = 1, 2, \dots, a$, esitetään jakamalla i :nnes rivi taulusta 2.1 sen reunajakauman summalla:

$$(2.3) \quad \mathbf{r}'_i = \left(\frac{p_{i1}}{p_{i.}}, \frac{p_{i2}}{p_{i.}}, \dots, \frac{p_{ib}}{p_{i.}} \right)$$

Elementit jokaisessa \mathbf{r}'_i ovat suhteellisia frekvenssejä ja siten niiden summa on 1:

$$(2.4) \quad \mathbf{r}'_i \mathbf{j} = \sum_{j=1}^b \frac{n_{ij}}{n_{i.}} = \frac{n_{i.}}{n_{i.}} = 1$$

Määritellään:

$$(2.5) \quad \mathbf{D}_r = \text{diag}(\mathbf{r}) = \begin{pmatrix} p_{1.} & 0 & \cdots & 0 \\ 0 & p_{2.} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & p_{a.} \end{pmatrix}$$

Nyt riviprofilien matriisi \mathbf{R} voidaan esittää myös:

$$(2.6) \quad \mathbf{R} = \mathbf{D}_r^{-1} \mathbf{P} = \begin{pmatrix} \mathbf{r}'_1 \\ \mathbf{r}'_2 \\ \vdots \\ \mathbf{r}'_a \end{pmatrix} = \begin{pmatrix} \frac{p_{11}}{p_{1.}} & \frac{p_{12}}{p_{1.}} & \cdots & \frac{p_{1b}}{p_{1.}} \\ \frac{p_{21}}{p_{2.}} & \frac{p_{22}}{p_{2.}} & \cdots & \frac{p_{2b}}{p_{2.}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{p_{a1}}{p_{a.}} & \frac{p_{a2}}{p_{a.}} & \cdots & \frac{p_{ab}}{p_{a.}} \end{pmatrix}$$

Samaan tapaan, j :nnes sarakeprofiili \mathbf{c}_j , $j = 1, 2, \dots, b$, määritellään jakamalla j :nnes sarake sen reunajakauman summalla:

$$(2.7) \quad \mathbf{c}'_j = \left(\frac{p_{j1}}{p_{j.}}, \frac{p_{j2}}{p_{j.}}, \dots, \frac{p_{jb}}{p_{j.}} \right)$$

Elementit jokaisessa \mathbf{c}'_i ovat suhteellisia frekvenssejä ja siten niiden summa on 1:

$$(2.8) \quad \mathbf{j}' \mathbf{c}_j = \sum_{i=1}^a \frac{n_{ij}}{n_{j.}} = \frac{n_{.j}}{n_{j.}} = 1.$$

Määritellään

$$(2.9) \quad \mathbf{D}_c = \text{diag}(\mathbf{c}) = \begin{pmatrix} p_{.1} & 0 & \cdots & 0 \\ 0 & p_{.2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & p_{.b} \end{pmatrix}$$

Nyt sarakeprofilien matriisi \mathbf{R} voidaan esittää myös:

$$(2.10) \quad \mathbf{R} = \mathbf{P} \mathbf{D}_c^{-1} = (\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_b) = \begin{pmatrix} \frac{p_{11}}{p_{.1}} & \frac{p_{12}}{p_{.2}} & \cdots & \frac{p_{1a}}{p_{.a}} \\ \frac{p_{21}}{p_{.1}} & \frac{p_{22}}{p_{.2}} & \cdots & \frac{p_{2a}}{p_{.a}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{p_{a1}}{p_{.1}} & \frac{p_{a2}}{p_{.2}} & \cdots & \frac{p_{ab}}{p_{.a}} \end{pmatrix}$$

Vektori \mathbf{r} on määritelty (2.1) \mathbf{P} :n rivisummien sarakevektoriksi. Se voidaan esittää sarakeprofiilien painotettuna keskiarvona:

$$(2.11) \quad \mathbf{r} = \sum_{j=1}^b p_{.j} \mathbf{c}_j$$

Samaan tapaan vektori \mathbf{r} on määritelty (2.2) \mathbf{P} :n sarakesummien rivivektoriksi ja se voidaan esittää riviprofiilien painotettuna keskiarvona:

$$(2.12) \quad \mathbf{c}' = \sum_{i=1}^a p_i \mathbf{r}'_i$$

Huomaa, että: $\sum_{j=1}^b p_{.j} = \sum_{i=1}^a p_i = 1$, tai

$$(2.13) \quad \mathbf{j}' \mathbf{r} = \mathbf{c}' \mathbf{j} = 1,$$

missä ensimmäinen \mathbf{j} on $a \times 1$ ja toinen on $b \times 1$. Tästä syystä $p_{.j}$:t ja p_i :t ovat sopivia painoja painotetussa keskiarvossa (2.12) ja (2.13) (Rencher, 2002.)

2.2 Riippumattomuuden testaaminen

Kuten esittelyssä todettiin, ristitaulukon dataa voidaan käyttää selvittämään kahden kategorisen muuttujan välistä yhteyttä. Jos kahta muuttujaa merkitään x :llä ja y :llä, oletus riippumattomuudesta voidaan esittää todennäköisyyksinä samalla tavalla kuin Rencher (2002) kirjassaan seuraavasti:

$$(2.14) \quad P(x_i y_j) = P(x_i) P(y_j), \quad i = 1, 2, \dots, a; \quad j = 1, 2, \dots, b,$$

missä x_i and y_j vastaa i :ttä riviä ja j :ttä saraketta ristitaulukossa. Käyttämällä taulukon (2.1) notaatiota voidaan estimoida (2.14)

$$(2.15) \quad p_{ij} = p_i \cdot p_{.j}, \quad i = 1, 2, \dots, a; \quad j = 1, 2, \dots, b.$$

χ^2 -testisuure x :n ja y :n riippumattomuuden testaamiseksi on

$$(2.16) \quad \chi^2 = n \sum_{i=1}^a \sum_{j=1}^b \frac{(p_{ij} - p_i \cdot p_{.j})^2}{p_i \cdot p_{.j}},$$

mikä on suurin piirtein (asymptoottisesti) χ^2 -jakautunut vapausastein $(a-1)(b-1)$.
Kaksi vaihtoehtoista muotoa (2.16) ovat

$$(2.17) \quad \chi^2 = n \sum_{i=1}^a np_{i.} \sum_{j=1}^b \left[\left(\frac{p_{ij}}{p_{i.}} - p_{.j} \right)^2 / p_{.j} \right],$$

$$(2.18) \quad \chi^2 = n \sum_{j=1}^b np_{.j} \sum_{i=1}^a \left[\left(\frac{p_{ij}}{p_{.j}} - p_{i.} \right)^2 / p_{i.} \right],$$

vektori ja matriisi muotona (2.17) ja (2.18) voidaan kirjoittaa

$$(2.19) \quad \chi^2 = \sum_{i=1}^a np_{i.} (\mathbf{r}_i - \mathbf{c})' \mathbf{D}_c^{-1} (\mathbf{r}_i - \mathbf{c}),$$

$$(2.20) \quad \chi^2 = \sum_{j=1}^b np_{.j} (\mathbf{c}_j - \mathbf{r})' \mathbf{D}_r^{-1} (\mathbf{c}_j - \mathbf{r}),$$

missä \mathbf{r} , \mathbf{c} , \mathbf{r}_i , \mathbf{c}_j , \mathbf{D}_r ja \mathbf{D}_c on määritelty kohdissa (2.1), (2.2), (2.3), (2.7), (2.5) ja (2.9). Kohdassa (2.19) verrataan \mathbf{r}_i :tä \mathbf{c} :hen jokaisella i :llä ja \mathbf{c}_j :tä \mathbf{r} :ään jokaisella j :llä. Täten kumpikin näistä on yhtä pätevä testaamaan riippumattomuutta, kun verrataan p_{ij} ja $p_i p_j$ kaikilla i, j , koska kaikki χ^2 määritelmät (2.16)-(2.20) ovat yhtäpitäviä.

Täten, jos muuttujat x ja y ovat riippumattomia, voidaan olettaa että, ristitaulukon riveillä on samankaltaiset profiilit tai yhtä lailla, sarakkeilla on samankaltaiset profiilit, jos muuttujat ovat riippumattomia. Riviprofiileja voidaan verrata keskenään vertaamalla jokaista riviprofilia \mathbf{r}'_i riviprofilien painotettuun keskiarvoon \mathbf{c}' (2.2). Tämä vertailu on tehty kohdassa (2.19). Samalla tavalla voidaan verrata sarakeprofiileja (2.20).

χ^2 -testisuure kohdassa (2.16) voidaan esittää vektori- ja matriisimuodossa

$$(2.21) \quad \chi^2 = n * tr[\mathbf{D}_r^{-1}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1}(\mathbf{P} - \mathbf{rc}')']$$

$$(2.22) \quad = n \sum_{i=1}^k \lambda_i^2$$

missä $\lambda_1^2, \lambda_2^2, \dots, \lambda_k^2$ ovat nolosta poikkeavia $\mathbf{D}_r^{-1}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1}(\mathbf{P} - \mathbf{rc}')'$:n ominaisarvoja ja

$$(2.23) \quad k = \text{rank}[\mathbf{D}_r^{-1}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1}(\mathbf{P} - \mathbf{rc}')'] = \text{rank}(\mathbf{P} - \mathbf{rc}').$$

$\text{rank}(\mathbf{P} - \mathbf{rc}')$ on tavallisesti $k = \min[(a - 1), (b - 1)]$. On selvää, että aste on vähemmän kun $\min(a, b)$, sillä

$$(2.24) \quad (\mathbf{P} - \mathbf{rc}')\mathbf{j} = \mathbf{P}\mathbf{j} - \mathbf{rc}'\mathbf{j} = \mathbf{r} - \mathbf{r} = \mathbf{0}$$

(Rencher, 2002.)

2.3 Kuvan muodostaminen

Rencher (2002) esittää kuvan muodostamisen kirjassaan seuraavasti: lasketaan rivi- ja sarakepisteet kaksiulotteiselle mallille datasta ristitaulukkaan. Aluksi skaalataan $\mathbf{P} - \mathbf{rc}'$, jotta saadaan

$$(2.25) \quad \mathbf{Z} = \mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{rc}')\mathbf{D}_c^{-1/2},$$

jonka elementit ovat

$$(2.26) \quad z_{ij} = \frac{p_{ij} - p_{i.}p_{.j}}{\sqrt{p_{i.}p_{.j}}},$$

ositetaan \mathbf{Z} käyttämällä singulaariarvohajotelmaa

$$(2.27) \quad \mathbf{Z} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}'.$$

\mathbf{U} matriisin sarakkeet ovat $\mathbf{Z}\mathbf{Z}'$ matriisin skaalatut ominaisvektorit. Taas \mathbf{V} matriisin sarakkeet ovat normalisoidut $\mathbf{Z}'\mathbf{Z}$ matriisin ominaisvektorit ja $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_k)$, missä $\lambda_1^2, \lambda_2^2, \dots, \lambda_k^2$ ovat $\mathbf{Z}'\mathbf{Z}$:n ja $\mathbf{Z}\mathbf{Z}'$:n ominaisarvot. Ominaisvektorit \mathbf{U} ja \mathbf{V} vastaavat ominaisarvoja $\lambda_1^2, \lambda_2^2, \dots, \lambda_k^2$. $\lambda_1, \lambda_2, \dots, \lambda_k$ kutsutaan \mathbf{U} :n yksittäis arvoiksi.

Nyt voidaan laskea hajotelma $\mathbf{P} - \mathbf{rc}'$:lle:

$$\begin{aligned}
D_r^{-1/2}(P - rc')D_c^{-1/2} &= U\Lambda V' \\
(2.28) \quad P - rc' &= D_r^{-1/2}U\Lambda V'D_c^{-1/2} \\
&= A\Lambda B' = \sum_{i=1}^k \lambda_i a_i b_i'
\end{aligned}$$

Missä $A = D_r^{-1/2}U$, $B = D_c^{-1/2}V$, a_i ja b_i ovat A :n ja B :n indeksiä i vastaavia sarakkeita.

Kohdassa (2.28) $P - rc'$:n rivit on esitetty B' :n rivien lineaarikombinaationa. Koordinaatit $P - rc'$:n i :nnelle riville löytyy $A\Lambda$:n i :nneltä riviltä. Vastaavasti $P - rc'$:n i :nnen sarakkeen koordinaatit löytyy $\Lambda B'$:n i :nneltä sarakeelta.

Jotta saadaan koordinaatit rivienpoikkeamille $r'_i - c'$, $R - jc'$:ssä ja sarakepoikkeamat $c_j - r$, $C - rj'$:ssä, esitetään molemmat matriisit $P - rc'$:n funktiona.

$$(2.29) \quad R - jc' = D_r^{-1}(P - rc'),$$

$$(2.30) \quad C - rj' = D_c^{-1}(P - rc').$$

Näin rivipoikkeamien koordinaatit $R - jc'$:ssä saadaan matriisioperaationa

$$(2.31) \quad X = D_r^{-1}AA$$

missä koordinaatit löytyvät X :n sarakkeilta.

Samalla tavalla sarakepoikkeamien koordinaatit $C - rj'$:ssä saadaan matriisioperaationa

$$(2.32) \quad Y = D_c^{-1}BA$$

missä koordinaatit löytyvät Y :n sarakkeilta.

Nyt, jotta voidaan piirtää koordinaatit riviprofilien poikkeamalle $r'_i - c'$, $i = 1, 2, \dots, a$ kahdessa ulottuvuudessa, tällöin

$$(2.33) \quad X_1 = \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{a1} & x_{a2} \end{pmatrix}$$

Vastaavasti, jotta voidaan piirtää koordinaatit sarakeprofiilien poikkeamalle $\mathbf{c}'_j - \mathbf{r}'_j$, $j = 1, 2, \dots, b$ kahdessa ulottuvuudessa, tällöin

$$(2.34) \quad \mathbf{Y}_1 = \begin{pmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ \vdots & \vdots \\ y_{b1} & y_{b2} \end{pmatrix}$$

Koska \mathbf{A} ja \mathbf{B} omaavat samat yksittäisarvot $\lambda_1, \lambda_2, \dots, \lambda_k$, voidaan ne esittää samassa kuvassa. Eetäisyys kahden riviprofiilin välillä on

$$(2.35) \quad d_{ij}^2 = (\mathbf{r}_i - \mathbf{r}_j)' \mathbf{D}_c^{-1} (\mathbf{r}_i - \mathbf{r}_j)$$

Jos kaksi rivipistettä ovat lähellä toisiaan, nämä pisteet voidaan yhdistää yhdeksi kategoriaksi, mikäli on tarvetta kasvattaa χ^2 aproksimaatiota. Rivi- ja sarakepisteen etäisyys ei ole niinkään merkitsevää, mutta niiden läheisyys on. Tällöin nämä kaksi kategoriata esiintyvät yleisemmin yhdessä, kun voisi ajatella niiden esiintyvän sattumalta, mikäli ne olisivat riippumattomia (Rencher, 2002.)

On näytetty, että rivi- ja sarakepisteiden metriikka on sama. Täten molemmat pisteparvet voidaan esittää samassa kuvassa. Glynn ja Robinson (2014) esittävät kirjassaan, että kuvan jokainen ulottuvuus kuvaa prosenttiosuutta varianssin rakenteesta. Jos kuvataan yhtä ulottuvuutta, saadaan yksinkertaisesti jana, joka kuvaa datapisteiden eroa etäisyytenä toisistaan. Yleensä tämä kuitenkin kuvaa melko huonosti datapisteiden yhteyttä toisiinsa. Jos lisätään toinen dimensio eli y-akseli, saadaan kaksiulotteinen pisteparvi, joka on erittäin yleinen esitystapa korrespondenssianalyysissä sekä monissa muissa spesifikaatioiden vähentämistekniikoissa. Yleensä kaksi ensimmäistä dimensiota kattavat suurimman osan datan varianssista, ja ovat täten riittävä. Mitä lähempänä pisteet ovat toisiaan, sitä enemmän voidaan ajatella niillä olevan riippuvuutta (Glynn & Robinson, 2014.)

3 Soveltaminen aineistoihin

3.1 Aineistot

Tutkielmassa sovelletaan korrespondenssianalyysia pankkidataan. Pankkidatassa on suoritettu kysely 36:lle pankissa työskentelevälle henkilölle. Kysely on eräänlainen työpaikkatyytyväisyyskysely, jossa on yli 60 erilaista kysymystä. Osa kysymyksistä on henkilötietoja ja suurin osa työhyvinvointiin liittyviä kysymyksiä, joihin on vastattu antamalla numero 1-5.

Toinen aineisto on Aarresaari-verkoston laatima ja toteuttama uraseurantakysely vuonna 2015 valmistuneille. Siinä on tutkittu 5 vuotta aiemmin valmistuneiden ylemmän korkeakoulututkinnon ja ns. päättävän alemman korkeakoulututkinnon suorittaneiden sijoittumista työelämään. Kyselyssä on kartoitettu valmistuneiden tyytyväisyyttä tutkintoonsa, työuran kokonaisuutta sekä tämänhetkistä työtilannetta. Data koostuu 88 erilaisesta kysymyksestä. Vastaajia kyselyyn on ollut 6830.

Tutkielmassa käytetään R-ohjelmistoa kuvien laatimiseen. R-koodin lähteenä on käytetty pääasiassa kirjan *Modern Applied Statistics with S. Fourth Edition* (Venables & Ripley, 2002) materiaalia.

3.2 Datan muokkaus

Pankkidatassa on 31 havaintoa ja 63 muuttujaa. Ensin data siistitään sellaiseen muotoon, että jäljellä ovat ainoastaan solut, joihin on kerätty dataa. Data sisältää useita puuttuvia arvoja. Koska havaintoja on niin vähän, ei ole järkevää käyttää ainoastaan niitä havaintoja, joissa ei esiinny puuttuvaa dataa. Niinpä onkin mielekästä generoida jotkin arvot puuttuvan datan tilalle. Tässä tapauksessa käytetään imputaatiota ja vielä täsmällisemmin, usean arvon imputaatiota. Pankkidatassa muuttujia on runsaasti, vaikkakin kategorisia muuttujia on hieman vähänlaisesti. Esimerkkiin on valittu kaksi muuttujaa: koulutusaste sekä ”Tärkein palkkauksesta, urakehityksestä”. Eri koulutusasteet ovat: peruskoulu, ammattikoulu, ylioppilas, keskiaste sekä korkeakoulu. Mahdolliset vastausvaihtoehdot kysymykseen ovat: Urakehitys on mahdollista työyhteisössäni, Palkkaukseni on oikeudenmukainen, Minulle tarjotaan riittävästi kouluttautumismahdollisuuksia, Ammattitaitoani hyödynnetään parhaalla mahdollisella tavalla, Uskon työskenteleväni pankissa vuonna 2003, Henkilöstö valitaan

objektiiivisin perustein.

Uraseurantadata sisältää myös puuttuvia arvoja, mutta koska havaintoja on niin runsaasti, voidaan analyysissä käyttää niitä arvoja, joista ei puutu dataa. Datasta on valittu muuttujiksi koulutusala, työmarkkinatilanne ja työnantaja. Koulutusalamuuttuja koostuu luonnontieteistä, teknillisestä alasta, teologiasta, oikeustieteellisestä sekä taidealasta. Työmarkkinatilanne kuvaa sitä, millainen työsuhde valmistunella on viiden vuoden jälkeen. Se koostuu vaihtoehdoista: vakituinen kokopäivätyö, määräaikainen kokopäivätyö, osa-aikatyö, itsenäinen yrittäjä, useita rinnakaisia työsuhteita, työllistetty/harjoittelija, työn työnhakija, työvoimakoulutus, päätoiminen opiskelu, perhevapaa (työsuhteessa), perhevapaa (ei työsuhteessa), työskentely apurahalla, työvoiman ulkopuolella, muu tilanne. Työnantajamuuttuja koostuu vaihtoehdoista: kunta tai kuntaryhmä, valtio, suuri yritys, pieni tai keskisuuri yritys, oma yritys, järjestö, säätiö tai seurakunta, yliopisto, ammattikorkeakoulu, muu työnantaja.

3.3 Puuttuvan tiedon paikkaaminen (Imputaatio)

Usean arvon imputaatioissa puuttuvat arvot paikataan joukolla M arvoja. Tyypillisesti ($5 \leq M \leq 10$) on sopiva määrä estimoimaan otosvaihtelua. Nummen (2019) mukaan imputaation vaiheet voidaan jakaa seuraavaan kolmeen vaiheeseen.

1. Luodaan M kappaletta imputaatioita.
2. Lasketaan M kappaletta parametriestimaatteja ja niille keskihajonnat.
3. Saadut estimaatit yhdistetään, jotta saadaan yksittäinen estimaatti parametrille. (Tässä tapauksessa valitaan jokin sopiva imputaatio ja käytetään sitä. Jätetään yhdistämisvaihe tekemättä, sillä se ei ole tässä tapauksessa niin oleellista (Nummi, 2019.))

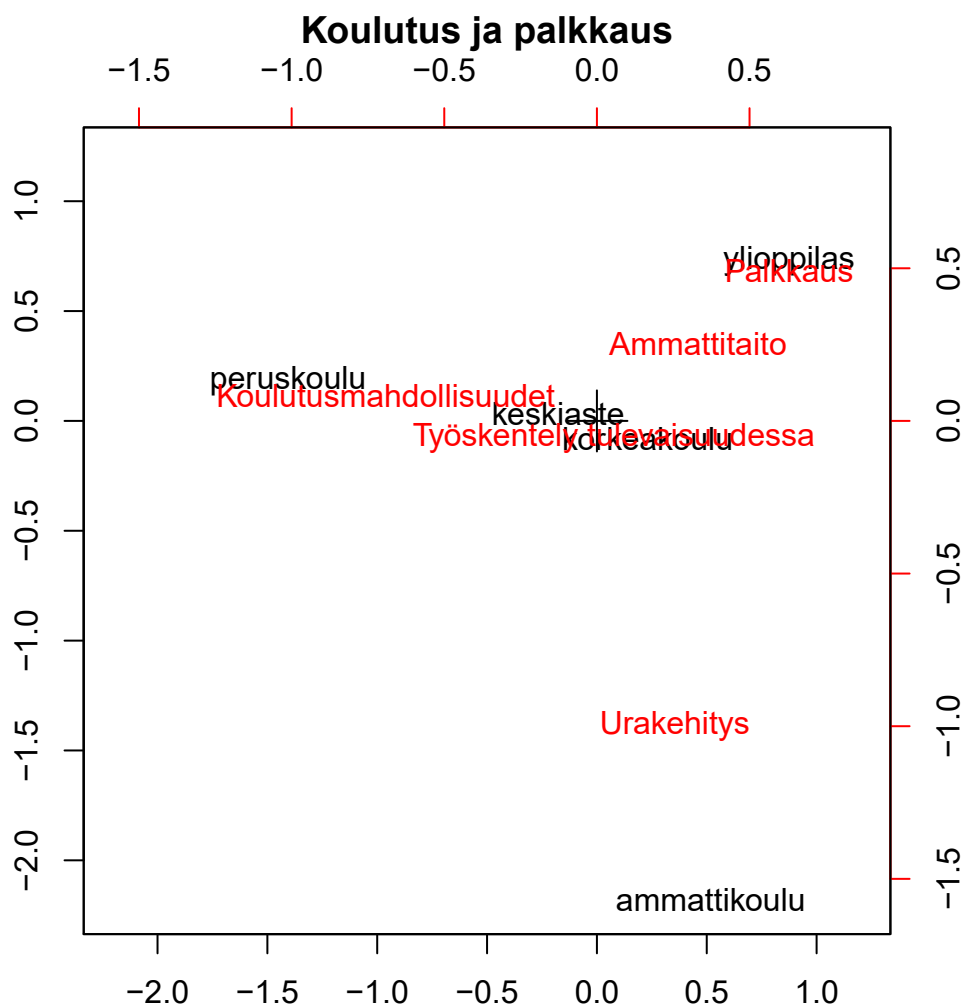
Työssä käytetty imputaation R-koodi on seuraavanlainen:

```
library ( mice )  
imp<-mice( pankkidata , m=5 , maxit = 10 , method="pmm"  
 , seed = 500 , print=FALSE )  
completePankkidata<-complete ( imp , 2 )  
complete . cases ( completePankkidata )
```

Näistä imputoiduista arvoista on valittu yksi paikkaamaan puuttuvaa tietoa.

3.4 Korrespondenssianalyysi pankkidataan

Sovelletaan korrespondenssianalyysia pankkidataan. Selvitetään pitävätkö eri koulutustaustoista tulevat työntekijät eri asioita tärkeimpänä heidän palkkauksessaan.



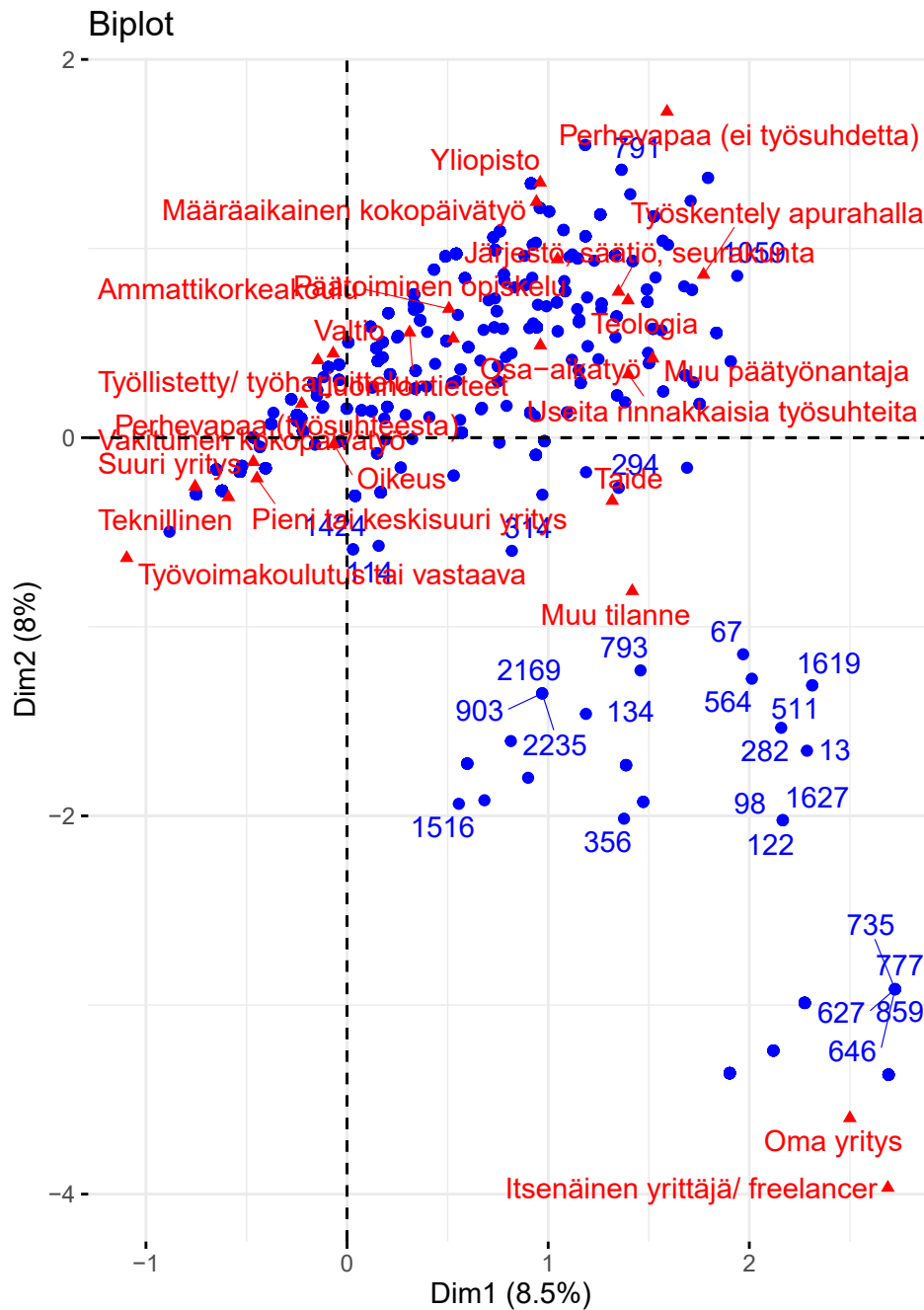
Kuva 3.1. Kuvassa on sovitettu molemmat muuttujat samaan kuvaan. Punaisella tekstillä näkyvät vastaukset ja mustalla tekstillä koulutustausta.

Kuvasta on havaittavissa, että keskiaste ja korkeakoulu ovat vastanneet melko tasaisesti kaikkia vaihtoehtoja. Pelkän peruskoulun suorittaneet pitävät koulutusmahdollisuuksia tärkeimpänä. Ylioppilaat puolestaan pitävät ammattitaidon hyödyntämistä ja palkkauksen oikeudenmukaisuutta tärkeimpänä. Ammattikoulun käyneet taas näkevät urakehitys mahdollisuuksia tärkeimpänä. Kuva näyttää erittäin selkeitä tuloksia, mutta tulkintoja tehdessä täytyy ottaa erityisesti huomioon vastanneiden

henkilöiden vähyys. Esimerkiksi vain ammattikoulun käyneitä vastaajia on ainoastaan yksi. Tästä syystä tuloksille ei voida antaa kovinkaan suurta painoarvoa.

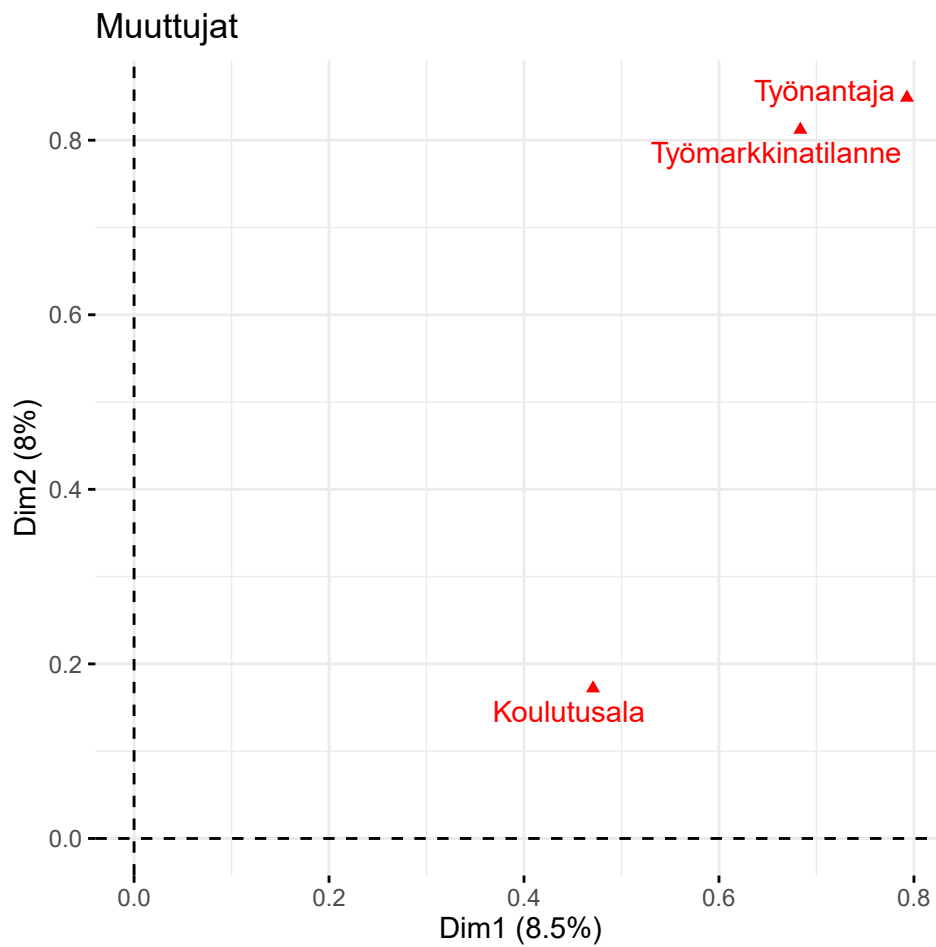
3.5 Korrespondenssianalyysi uraseurantadataan

Toisessa esimerkissä sovelletaan menetelmää uraseurantadataan. Pyritään selvittämään, millaisiin työsuhteisiin eri koulutusaloilta on päädytty.



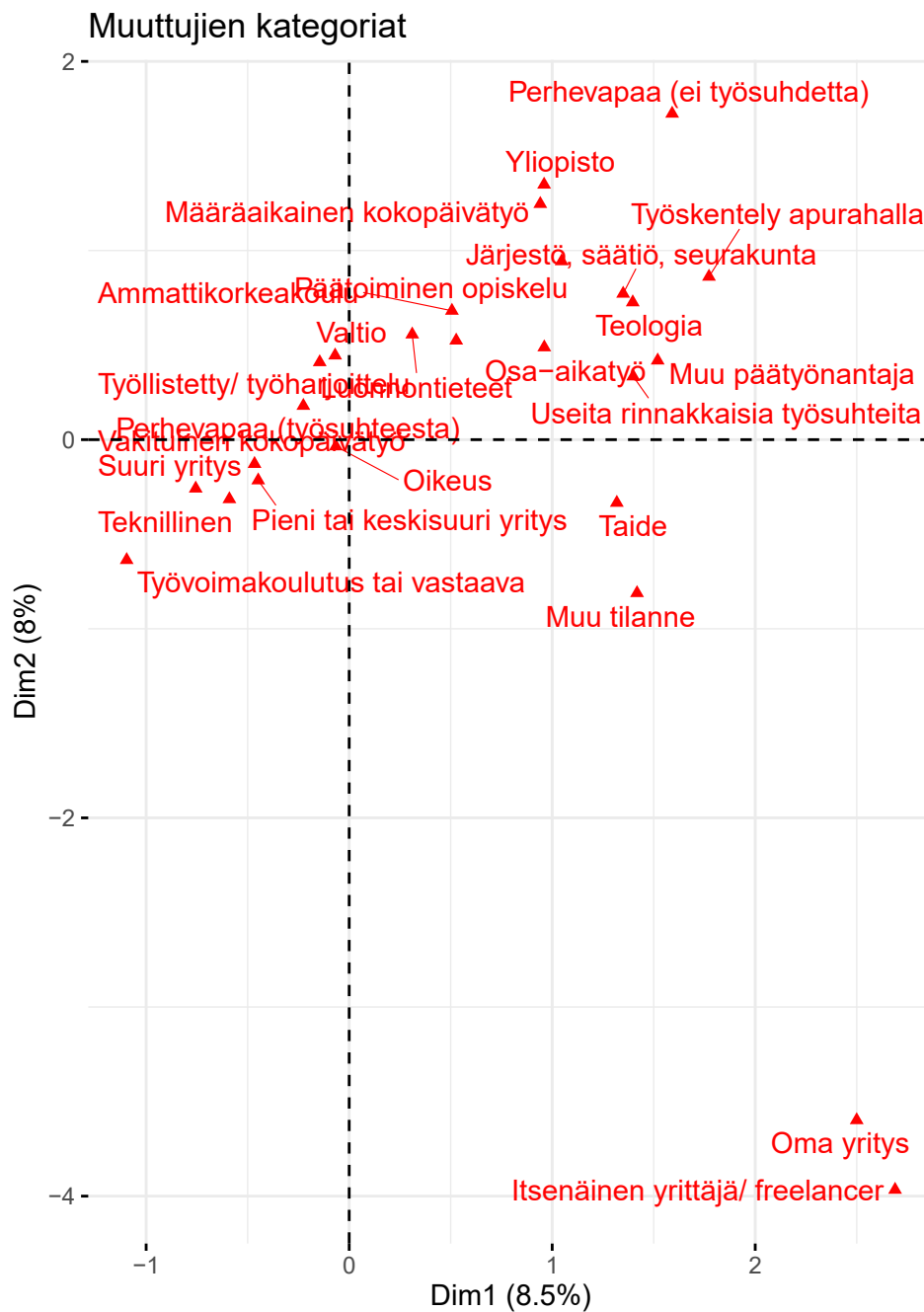
Kuva 3.2. Kuvassa on sovitettu kategoriat ja havainnot samaan kuvaan. Sinisellä värillä on merkitty havaintojen pisteet sekä monesko havainto on kyseessä. Punaisella taas on merkitty kategorioiden pisteet ja niiden nimet.

Biplot-kuva antaa osviittaa siitä, miten eri havainnot ovat sijoittuneet eri kategorioiden välillä. Voidaan huomata, että ”oma yritys” kategoriaan päin on joitakin arvoja kallellaan, mutta suurimmaksi osaksi arvot ovat keskittyneet lähelle origoa.



Kuva 3.3. Kuvassa on esitetty muuttujien yhteys toisiinsa.

Kuvasta huomataan, että työnantajan ja työmarkkinatilanteen välillä on korrelaatiota. Kuitenkaan näitä ei ole tarpeellista yhdistää.



Kuva 3.4. Kuvassa kategorioiden yhteys toisiinsa.

Kuva 3.4 on muuten sama kuin kuva 3.2, mutta kuvassa 3.4 on eritelty ainoastaan kategoriat. Kuvasta huomataan, että teknilliseltä alalta työllistytään yleisimmin erikokoisiin yrityksiin. Myös vakituinen kokopäivätyö on erittäin yleinen tällä alalla. Oikeustieteellinen on melko keskellä, eli se on liitettävissä jossakin määrin kaikkiin kategorioihin, joskin siihen liittyy valtiollinen sektori vahvasti, mikä käy hyvin järkeen. Teologian alalla työllistytään luonnollisesti seurakuntaan, mikä selittää pisteiden lähes päällekkäisyyden. Taideala on melko irrallaan muista kategorioista. Ai-

noastaan ”muu tilanne” on sen läheisyydessä. Taideala on myös hieman kallellaan ”itsenäinen yrittäjä” suuntaan verrattuna muihin aloihin.

4 Johtopäätökset

Kuten tutkielmassa on havaittu, korrespondenssianalyysi on erittäin hyvä menetelmä visualisoimaan kategorista dataa. Vaikka menetelmän matemaattinen esitys on pitkäkö, on sen avulla kuitenkin verrattain helppo saada yksinkertainen ja havainnollistava kuvaus datan käyttäytymisestä. Yksinkertaisen korrespondenssianalyysin tekeminen R-ohjelmalla on suhteellisen helppoa, eikä kuvankaan tulkitseminen ole kovinkaan hankalaa. Kuvasta saattaa tosin muodostua melko sekava, jos muuttujia on useampia. Toisaalta korrespondenssianalyysille löytyy vaihtoehtoisiaakin menetelmiä, eikä se tunnu olevan yleisesti erityisen käytetty menetelmä.

Menetelmän ongelma kuitenkin on siinä, että menetelmää voidaan hyödyntää ainoastaan kategorisiin muuttujiin. Työssä huomattiin myös, että havaintoja pitää olla melko paljon, etenkin multiple korrespondenssianalyysia tehtäessä, jotta analyysillä saadaan aikaiseksi hyviä tuloksia. Tilanteiden joukko, mihin menetelmä soveltuu hyvin, tuntuu olevan verrattain kapea. Se saattaa osaltaan selittää sen, miksi menetelmä ei ole kovin yleinen.

Tutkimusta voisi jatkaa esimerkiksi vertailemalla korrespondenssianalyysia toisiin monimuuttujamenetelmiin, esimerkiksi pääkomponenttianalyysiin.

Lähteet

- [1] Glynn, D. & Robinson, J.(2014) *Corpus Methods for Semantics Quantative studies in polysemy and synonymy*. University of Paris VIII, University of Sussex.
- [2] Nummi, T.(2019) *Statistical analysis with missing data*. Tampere: University of Tampere Faculty of Natural Sciences.
- [3] Rencher, A. *Methods of Multivariate Analysis. Second Edition*. Brigham Young University, 2002.
- [4] Venables, V.N. & Ripley, B.D.(2002) *Modern Applied Statistics with S. Fourth Edition*.
- [5] Le Roux, B. & Rouanet, H.(2010) *Multiple Correspondence Analysis*. Los Angeles: SAGE.

A Kuvaajien R-koodi

```
con4<-table(completePankkidata\C60, completePankkidata\C51)
con4 <- as.data.frame.matrix(con4)
corresp(con4)
par("mar")
par(mar=c(2,2,2,2))
biplot(corresp(con4, nf = 2), xlim=c(-2.2, 1.2), ylim=c(-2.2, 1.2));
title("Koulutus_ ja_ palkkaus", line = 2.2)

res.mca <- MCA(dfura5, graph = TRUE)
print(res.mca)

fviz_mca_var(res.mca, choice = "mca.cor",
             title = "Muuttujat",
             repel = TRUE,
             ggtheme = theme_minimal())

fviz_mca_biplot(res.mca,
               title = ("Biplot"),
               repel = TRUE,
               ggtheme = theme_minimal())

fviz_mca_var(res.mca,
             title = ("Muuttujien_ kategoriat"),
             repel = TRUE,
             ggtheme = theme_minimal())
```