

# Kaavat bookdown-paketilla

*Jussi Hirvonen*

*3.8.2018*

## Kaavat ja matemattiset merkinnät

Kaavat on esitettävä bookdown-paketin määrittäyksillä. Viittausnimien on oltava yksikäsitteisiä koko dokumentissa, jos käytetään “merge and knit” menetelmää. Jos taas jokainen lapsidokumentti on “itsenäinen” (“knit and merge”), tämä koskee vain kyseistä dokumenttia (kts. Bookdown - webkirja).

## Kahden luokittelumuuttujan taulukko

Kahden luokittelumuuttujan riippuvuutta voidaan testata  $\chi^2$  - testillä. Testisuure saadaan laskemalla yhteen jokaisen solun havaittujen ja odotettujen (riippumattomuushypoteesi) frekvenssien erotukset muodossa

$$\chi^2 = \frac{(\text{havaittu} - \text{odotettu})^2}{\text{odotettu}} (\#eq : khii21) \quad (1)$$

Tämä voidaan esittää ca:han sopivammalla tavalla parilla muunnoksella, jolloin saamme riveittäin vastaavat termit rivisummalla painotettuna:

$$\text{rivisumma} \times \frac{(\text{havaittu riviprofiili} - \text{odotettu riviprofiili})^2}{\text{odotettu riviprofiili}} (\#eq : khii22) \quad (2)$$

Kun jaamme nämä tekijät havaintojen kokonaismäärällä  $n$ , rivisumma muuntuu rivin massaksi, ja niiden summa muotoon  $\frac{\chi^2}{n}$ .

$$\frac{\chi^2}{n} = \phi^2 (\#eq : inert1) \quad (3)$$

Tunnusluku  $\phi^2$  on korrespondenssianalyysissä kokonaisinertia (total inertia). Se kuvaa, kuinka paljon varianssia taulukossa on ja on riippumaton havaintojen lukumäärästä. Tilastotieteessä tunnusluvulla on useita vaihtoehtoisia nimiä (esim. mean square contingency coefficient), ja sen neliöjuurta kutsutaan  $\phi$  - kertoimeksi.

Tässä siirrytään kahden luokittelumuuttujan taulukosta suhteellisten frekvenssien taulukkoon, ja pieni pohdinta taulukoista yleensä olisi paikallaan. Kaavojen

@ref(eq:khii21) ja @ref(eq:khii22) yhteyden pitäisi olla selkeä. Frekvenssitaulukossa (jossa kaikki taulukon luvut on jaettu havaintojen lukumäärällä n) riviprofilien 1 ja 3 (euklidinen) etäisyys on

$$\sqrt{(p_{11} - p_{31})^2 + (p_{12} - p_{32})^2 + (p_{13} - p_{33})^2 + (p_{14} - p_{34})^2 + (p_{15} - p_{35})^2} (\#eq : khii23) \quad (4)$$

Rivien  $\chi^2$  - etäisyys on painotettu euklidinen etäisyys, jossa painoina ovat riviprofilin odotetut arvot. Ne ovat riippumattomuushypoteesin mukaisesti riviprofilien keskiarvoprofilin vastaavat alkioit  $r_i$ .

$$\sqrt{\frac{(p_{11} - p_{31})^2}{r_1} + \dots + \frac{(p_{15} - p_{35})^2}{r_5}}$$

Inertia voidaan esittää rivien ja **keskiarvorivin** (sentroidin)

$$\chi^2$$

-etäisyyksien neliöiden painotettuna summana, jossa painoina ovat rivien massat  $m_i$  ja summa lasketaan yli rivien  $i$ .

$$\phi^2 = \sum_i (massa\ m_i) \times (profiilin\ i\ \chi^2 - etaisyys\ sentroidista)^2$$

Kaavat.tex - dokumentissa on tässä kohdassa testailtu R:n furniture - paketin taulukoita latex- ja latex2 - output-formaateilla. Ne voi liittää LateX-dokumenttiin, jossa on käytössä paketti booktabs. Bookdownissa luultavasti tämä on tarpeeton, kable riittänee.

## Matriisit ja niiden havainnollistaminen

Näissä ei ole vielä numeroita ja viitetietoa. Ei kutoudu rmd-tiedostosta, mutta tekee tex-tiedoston ja siitä saa luotua pdf - tiedoston. drawmatrix - paketti ei ihan tunnu toimivan, toistaiseksi ei mukana.

### drawmatrix - kaavoja

Yksinkertainen korrespondenssianalyysi on kahden luokittelumuuttujan määrittelyn frekvenssitaulukon analyysiä. Taulukon rivit ovat havaintoyksiköiden (individuals, havaintoyksikkö) aggregoituja summia, sarakkeet muuttujia.

Analyyssissä osa riveistä tai sarakkeista voidaan jättää pois ratkaisun laskennasta ns. passiiviksi, ja esittää kartalla täydentävinä pisteinä (supplementary points). Ne eivät vaikuta ratkaisuun, eli teknisesti niiden massa on nolla, mutta pisteiden esityksen (projektion) tarkkuus voidaan arvioida. Täydentävien profilien on kuitenkin oltava yhteismitallisia taulukon datan kanssa. Mikä tahansa ei käy

(kts. CAinP, vast.luku). Pinotut tai yhdistetyt matriisit (“stacked matrices”). Yksinkertainen korrespondenssianalyysi on kahden luokittelumuuttujan määrittämän taulukon (kontingenssitaulukko) analyysiä, mutta tutkimusasetelmaa voi melko helposti muuttaa useamman muuttujan analyysiksi. Menetelmän materiaattinen perusta ja ratkaisualgoritmi (SVD) toimivat, tulkinta vain muuttuu. Itse asiassa menetelmän yleisyys tekee sen vääränkin käytön mahdolliseksi.

Yksinkertaisin laajennus on lisätä alkuperäisen taulukon alle toinen taulukko. Rivit ovat esimerkissä maittan summattuja vastauksia, ja niiden alle voidaan lisätä joku toinen luokittelumuuttuja. Havaintojen määrä yhdistetyssä (“pinotussa”) taulussa kaksinkertaistuu. Miksi tämä ei vaikuta tuloksiin vääristävästi??

Merkitään edellisten analyysien kuuden maan ja viiden vastausvaihtoehdon taulukkoa matriisilla  $\mathbf{A}_{IJ}$ , missä  $I$  on rivien ja  $J$  sarakkeiden lukumäärä. Taulukoidaan ikäluokan (1 - 6) ja sukupuolen ( $f$  = nainen,  $m$  = mies) vuorovaikutusmuuttuja ( $f1, \dots, f6$  ja  $m1, \dots, m6$ ) samojen vastausvaihtoehtojen kanssa. Jos tätä taulukkoa merkitään matriisilla  $\mathbf{B}_{I'J}$ , voimme muodostaa yhdistetyn matriisin

#### **drawmatrix - kaavoja**

Miten päällekkäisten matriisien ympärille saisi sulut?

Rivien lukumäärä on molemmissa matriiseissa sama, koska luokkia sattuu olemaan kuusi sekä maa- että ikä- ja sukupuoli - luokittelumuuttujissa. Kun matriisit ovat dimensioiltaan ja myös muuttujien sisällön kannalta samankaltaiset, niitä kutsutaan yhteensopiviksi (“matched matrix”). Tällöin yksinkertaista korrespondenssianalyysissä voi soveltaa tutkimusongelmaan, jossa halutaan erotella jonkun ryhmän sisäinen vaihtelu ryhmien välisestä vaihtelusta. (Greenacren ehdottama ABBA - analyysi).

#### **drawmatrix - kaavoja**

ABBA on erityistapaus yleisemmästä moniulotteisen taulukon (multiway table) analyysistä, jossa useita kahden muuttujan taulukoita “pinotaan” päällekkäin ja rinnakkain. Voimme ottaa yhden kysymyksen vastausten lisäksi analyysiin mukaan useamman kysymyksen vastaukset laajentamalla kahden päällekkäisen matriisin taulukkoa oikealle.

Teknisesti analyysi on yksinkertainen korrespondenssianalyysi, miten tämä tulkitaan?

#### **drawmatrix - kaavoja**

### **Korrespondenssianalyysin perusyhtälöt ja kaavat**

#### **viitetiedot puuttuvat kaavoista**

Tässä lähteenä Greenacren kirja (ca in practice) ja sen liite Theory of CA. Muistiinpanoja löytyy, joissa viitataan myös Biplots in practice - kirjaan. Kevään

2017 kurssin luentokalvoja on myös käytetty. Lisäillään vielä käsitteitä LeRouxin ja Rouanetin kirjasta.

Datamatriisilla  $\mathbf{N}$  on  $I$  riviä ja  $J$  saraketta ( $I \times J$ ). Alkiot ovat ei-negatiivisia (eli nollat sallittuja) ja samassa mitta-asteikossa. Jos mitta-asteikko on intervallitai suhdeasteikko, mittayksiköiden on oltava samoja (esim. euroja, metrejä). Taulukon alkioden summa on  $\sum_i \sum_j n_{ij} = n$ , missä  $i = 1, \dots, I$  ja  $j = 1, \dots, J$ . GDA-kirjassa on tarkennettu tätä vaatimusta ei-negatiivisuudesta.

Korrespondenssimatriisi  $\mathbf{P}$  saadaan jakamalla matriisin  $\mathbf{N}$  alkiot niiden summalla  $n$ . Merkitään matriisin  $\mathbf{P}$  rivisummien vektoria  $\mathbf{r}$  ( $= (r_1, \dots, r_I)$ ) ja sarakesummien vektoria  $\mathbf{c}$  ( $= (c_1, \dots, c_J)$ ). Niitä vastaavat diagonaalimatriisit ovat  $\mathbf{D_r}$  ja  $\mathbf{D_c}$ .

Korrespondenssianalyysin perusrakenne (algoritmi?) on tämä. Singulaariarvohajoitelma (singular value decomposition) tuottaa ratkaisun kun sitä sovelletaan standardoituun residuaalimatriisiin  $\mathbf{S}$ .

$$\mathbf{S} = \mathbf{D_r}^{-1/2}(\mathbf{P} - \mathbf{r}\mathbf{c}^T)\mathbf{D_c}^{-1/2} \quad (5)$$

Residuaalimatriisi voidaan esittää myös ns. kontingenssi-suhdelukujen (contingency ratio) avulla.

$$\mathbf{D_r}^{-1}\mathbf{P}\mathbf{D_c}^{-1} = \left( \frac{p_{ij}}{r_i c_j} \right)$$

$$\mathbf{S} = \mathbf{D_r}^{1/2}(\mathbf{D_r}^{-1}\mathbf{P}\mathbf{D_c}^{-1} - \mathbf{1}\mathbf{1}^T)\mathbf{D_c}^{-1/2} \quad .$$

Toinen esitystapa on hyödyllinen, kun tarkastellaan CA:n yhteyksiä muihin läheisiin menetelmiin (log ratio analysis of compositional data, moniulotteinen skaalaus (?), lineaarinen diskriminanttianalyysi, kanoninen korrelaatioanalyysi, pääkomponenttianalyysi, kaksoiskuvat, yleensä SVD-perusteiset dimensioiden vähentämisen menetelmät).

$$s_{ij} = \frac{p_{ij} - r_i c_j}{\sqrt{r_i c_j}}$$

ja toinen

$$s_{ij} = \sqrt{r_i} \left( \frac{p_{ij}}{r_i c_j} \right) \sqrt{c_j} \quad .$$

Mitäköhän tuosta pitäisi nähdä? Selitykset löytyvät em. teorialiitteestä.

Singulaariarvohajoitelma (singular value decomposition, SVD) matriisille  $\mathbf{S}$  on

$$\mathbf{S} = \mathbf{U}\mathbf{D}_\alpha\mathbf{V}^T$$

missä  $D_\alpha$  on diagonaalimatriisi, jonka alkiot ovat singulaariarvot suuruusjärjestyksessä  $\alpha_1 \geq \alpha_2 \geq \dots$

Matriisit  $U$  ja  $V$  ovat ortogonaalisia singulaarivektoreiden matriiseja. Singulaariarvohajoituksen merkitys dimensioiden vähentämiselle perustuu Eckart - Young - teoreemaan. Teoreema (30-luvulta?) kertoo, että saamme pienimmän neliösumman  $m$  - ulotteisen approksimaation matriisille  $S$  (CAinP, ss. 244) matriisien  $U$  ja  $V$  ensimmäisten sarakkeiden ja ensimmäisten singulaariarvojen avulla.

$$S_{(m)} = U_{(m)} D_{\alpha(m)} V_{(m)}^T$$

Korrespondenssianalyysin ratkaisualgoritmissa tätä tulosta on muokattava niin, että rivien ja sarakkeiden massat huomioidaan pienimmän neliösumman approksimaatiossa painoina.

Näin saadaan standardikoordinaatit ja principal-koordinaatit riveille ja sarakkeille.

Rivien standardikoordinaatit

$$\Phi = D_r^{-\frac{1}{2}} U \quad (6)$$

Sarakkeiden standardikoordinaatit

$$\Gamma = D_c^{-\frac{1}{2}} V \quad (7)$$

Rivien principal-koordinaatit

$$F = D_r^{-\frac{1}{2}} U D_\alpha = \Phi D_\alpha \quad (8)$$

Sarakkeiden principal-koordinaatit

$$G = D_c^{-\frac{1}{2}} V D_\alpha = \Gamma D_\alpha \quad (9)$$

Pääakselien inertiat (principal inertias)  $\lambda_k$

$$\lambda_k = \alpha_k^2, k = 1, \dots, K, K = \min\{I - 1, J - 1\} \quad (10)$$

Bilineaarinen korresepondenssimalli

Korrespondenssimatriisi  $P$  voidaan esittää matriisi- ja alkiomuodossa ns. palautuskaavana (reconstitution formula).

$$P = D_r \left( \mathbf{1}\mathbf{1}^T + \Phi D_\alpha^{-\frac{1}{2}} \Gamma^T \right) D_c \quad (11)$$

$$p_{ij} = r_i c_j \left( 1 + \sum_{k=1}^K \sqrt{\lambda_k} \phi_{ik} \gamma_{jk} \right) \quad (12)$$

Tässä viitataan s. 101 (13.4), 109 (14.9), ja 109-110 (14.10 ja 14.11). Palautuskaavoilla on monta esitystapaa bilineaarisessa mallissa.

Rivien ja sarakkeiden riippuvuus ja transitioyhtälöt. ss. 244, 108-109 skalaarivermiot.

Pääkoordinaatit standardikoordinaattien funktiona (ns. barysentrinen ominaisuus - barycentric relationships)

$$\mathbf{F} = \mathbf{D}_r^{-1} \mathbf{P} \mathbf{\Gamma} \quad (13)$$

$$\mathbf{G} = \mathbf{D}_c^{-1} \mathbf{P}^T \mathbf{\Phi} \quad (14)$$

Pääkoordinaatit pääkoordinaattien funktioina:

$$\mathbf{F} = \mathbf{D}_r^{-1} \mathbf{P} \mathbf{G} \mathbf{D}_\lambda^{-\frac{1}{2}} \quad (15)$$

$$\mathbf{G} = \mathbf{D}_c^{-1} \mathbf{P}^T \mathbf{F} \mathbf{D}_\lambda^{-\frac{1}{2}} \quad (16)$$

Yhtälöt (9) ja (10) esittävät profilipisteet ideaalipisteiden (vertex points) painotettuina keskiarvoina, painoina profiilin elementit. Asymmetriset kartat (rivien tai sarakkeiden suhteen) perustuvat näihin yhtälöihin. Yhtälöiden (11) ja (12) kahdet pääkoordinaatit ovat perusta symmetrisille kartoille. Myös niitä yhdistää barisentrinen painotetun keskiarvon riippuvuus, mutta mukana ovat skaalaustekijät  $\frac{1}{\sqrt{\lambda_i}}$ . Ne ovat jokaisessa dimensiossa eri suuruisia.

Kokeillaan vielä kaavaviitteitä: kaavojen @ref(eq:khii21) ja @ref(eq:khii22) yhteyden pitäisi olla selkeä.