

Kaavat bookdown-paketilla

Jussi Hirvonen

13.6.2018

Sisältö

1	Kaavat ja matemattiset merkinnät	1
1.1	Kahden luokittelumuuttuja taulukko	1
1.2	Matriisit ja niiden havainnollistaminen	2
1.3	Korrespondenssianalyysin perusyhtälöt ja kaavat	3

1 Kaavat ja matemattiset merkinnät

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

1.1 Kahden luokittelumuuttuja taulukko

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

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

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

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

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

$$\frac{\chi^2}{n} = \phi^2 \quad (3)$$

Tunnusluku ϕ^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 ϕ - kertoimeksi.

Tässä siirrytään kahden luokittelumuuttujan taulukosta suhteellisten frekvenssien taulukkoon, ja pieni pohdinta taulukoista yleensä olisi paikallaan. Kaavojen (1) ja (2) yhteyden pitäisi olla selkeä. Frekvenssitaulukossa (jossa kaikki taulukon luvut on jaettu havaintojen lukumäärällä n) riviprofilien 1 ja 3 (euklidinen) etäisyys on

$$\sqrt{(p_{11} - p_{31})^2 + (p_{12} - p_{32})^2 + (p_{13} - p_{33})^2 + (p_{14} - p_{34})^2 + (p_{15} - p_{35})^2}$$

Rivien χ^2 - etäisyys on painotettu euklidinen etäisyys, jossa painoina ovat riviprofiilin odotetut arvot. Ne ovat riippumattomuushypoteesin mukaisesti riviprofiilien keskiarvoprofiilin 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 .

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

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.

1.2 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 määrittelmän frekvenssitaulukon analyysiä. Taulukon rivit ovat havaintoyksiköiden (individuals, havaintoyksikkö) aggregoituja summia, sarakkeet muuttujia.

Analyyssissä osa riveistä tai sarakkeista voidaan jättää pois ratkaisun laskennasta ns. passiivisiksi, 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 profiilien 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 määrittämisen taulukon (kontingenssitaulukko) analyysiä, mutta tutkimusasetelmaa voi melko helposti muuttaa useamman muuttujan analyysiksi. Menetelmän matemaattinen 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 esimerkiksi maittaan summattuja vastauksia, ja niiden alle voidaan lisätä joku toinen luokittelumuuttuja. Havaintojen määrä yhdistetyssä ("pinotussa") taulussa kaksinkertaistuu. Miksi tämä ei ei vaikuta tuloksiin väärin??

Merkitään edellisten analyysien kuuden maan ja viiden vastausvaihtoehdon taulukkoa matriisilla 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 B_{IJ} , 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

1.3 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. Lisäilläään 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 intervalli- tai 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} .

$$\mathbf{S} = \mathbf{D_r}^{-1/2}(\mathbf{P} - \mathbf{r}\mathbf{c}^T)\mathbf{D_c}^{-1/2} \quad (4)$$

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 \mathbf{D}_α on diagonaalimatriisi, jonka alkiot ovat singulaariarvot suuruusjärjestyksessä $\alpha_1 \geq \alpha_2 \geq \dots$

Matriisit \mathbf{U} ja \mathbf{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 \mathbf{S} (CAinP, ss. 244) matriisien \mathbf{U} ja \mathbf{V} ensimmäisten sarakkeiden ja ensimmäisten singulaariarvojen avulla.

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

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

$$\mathbf{\Phi} = \mathbf{D}_r^{-\frac{1}{2}} \mathbf{U} \quad (5)$$

Sarakkeiden standardikoordinaatit

$$\mathbf{\Gamma} = \mathbf{D}_c^{-\frac{1}{2}} \mathbf{V} \quad (6)$$

Rivien principal-koordinaatit

$$\mathbf{F} = \mathbf{D}_r^{-\frac{1}{2}} \mathbf{U} \mathbf{D}_\alpha = \mathbf{\Phi} \mathbf{D}_\alpha \quad (7)$$

Sarakkeiden principal-koordinaatit

$$\mathbf{G} = \mathbf{D}_c^{-\frac{1}{2}} \mathbf{V} \mathbf{D}_\alpha = \mathbf{\Gamma} \mathbf{D}_\alpha \quad (8)$$

Pääakselien inertiat (principal inertias) λ_k

$$\lambda_k = \alpha_k^2, k = 1, \dots, K, K = \min\{I - 1, J - 1\} \quad (9)$$

Bilineaarinen korrespondenssimalli

Korrespondenssimatriisi \mathbf{P} voidaan esittää matriisi- ja alkiomuodossa ns. palautuskaavana (reconstitution formula).

$$\mathbf{P} = \mathbf{D}_r \left(\mathbf{1}\mathbf{1}^T + \mathbf{\Phi} \mathbf{D}_\alpha^{\frac{1}{2}} \mathbf{\Gamma}^T \right) \mathbf{D}_c \quad (10)$$

$$p_{ij} = r_i c_j \left(1 + \sum_{k=1}^K \sqrt{\lambda_k} \phi_{ik} \gamma_{jk} \right) \quad (11)$$

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

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

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

$$\mathbf{F} = \mathbf{D}_r^{-1} \mathbf{P} \mathbf{\Gamma} \quad (12)$$

$$\mathbf{G} = \mathbf{D}_c^{-1} \mathbf{P}^T \mathbf{\Phi} \quad (13)$$

Pääkoordinaatit pääkoordinaattien funktioina:

$$\mathbf{F} = \mathbf{D}_r^{-1} \mathbf{P} \mathbf{G} \mathbf{D}_\lambda^{-\frac{1}{2}} \quad (14)$$

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

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 barisentrisen painotetun keskiarvon riippuvuus, mutta mukana ovat skaalaustekijät $\frac{1}{\sqrt{\lambda_i}}$. Ne ovat jokaisessa dimensiossa eri suuruisia.

Kokeillaan vielä kaavaviitteitä: kaavojen (1) ja (2) yhteyden pitäisi olla selkeä.