

Korrespondenssianalyysi - kaavaliite (v. 1.03)

Jussi Hirvonen

4.4.2018

Paperin versiot		
Versio	muutokset	päivämäärä
0.1	harjoittelua - drawmatrix	14.7.2017
1.01	harjoittelua - matriisiyhtälöt ja ca-peruskaavat	5.8.2017
1.02	pienää korjailua, turhan poistoa	28.1.2018
1.03	lisäilläään yksinkertaisen ca:n kaavoja, taulukoita R-paketilla furniture	4.4.2018

Sisältö

1	Kaavat ja matemaattisen merkinnät - työkalut	1
2	Yksinkertaisia kaavoja leipätekstiin	1
3	Matriisit ja niiden havainnollistaminen	2
4	Matriisiyhtälöt ilman kaavioita	3

1 Kaavat ja matemaattisen merkinnät - työkalut

Tähän voisi kerätä Latex'in suositukset ja tavat esittää kaavoja, tässä on pientä epäselvyyttä mulla vielä. Esimerkkejä löytyy varsinaisesta luonnosdokkarista. Myös kaavito sun muut, kuvien insertointi on varsinaisessa luonnosdokkarissa.

2 Yksinkertaisia kaavoja leipätekstiin

Yksinkertaisen korrespondenssianalyysin esittelyssä tarvitaan muutama kaava, vaikka koitan niitä liitteen ulkopuolella välttää.

Taulukko 1: Alle kouluikäinen lapsi todennäköisesti kärsii, jos hänen äitinsä käy työssä

	1	2	3	4	5	Test	P-V
	n = 295	n = 600	n = 593	n = 889	n = 732		
maa						Chi Square: 588.1	<.0
FI	47 (15.9%)	188 (31.3%)	149 (25.1%)	423 (47.6%)	303 (41.4%)		
HU	219 (74.2%)	288 (48%)	225 (37.9%)	190 (21.4%)	75 (10.2%)		
SE	29 (9.8%)	124 (20.7%)	219 (36.9%)	276 (31%)	354 (48.4%)		

Kahden luokittelumuuttujan ristiintaulukointi (kontigenssitalu)

En käytä kaavaliitteen notaatiota leipätekstissä, vaan yksinkertaisempaa tapaa. Aloietetaan taulukolla.

Koitetaan furniture-paketilla vääntää taulukko. Tämä menee oikealta yli (output = latex) .

Toinen koe (output = latex2):

Taulukko 2: Alle kouluikäinen lapsi todennäköisesti kärsii, jos hänen äitinsä käy työssä

	V6				
	1	2	3	4	5
	n = 295	n = 600	n = 593	n = 889	n = 732
maa					
FI	47 (15.9%)	188 (31.3%)	149 (25.1%)	423 (47.6%)	303 (41.4%)
HU	219 (74.2%)	288 (48%)	225 (37.9%)	190 (21.4%)	75 (10.2%)
SE	29 (9.8%)	124 (20.7%)	219 (36.9%)	276 (31%)	354 (48.4%)

Yritetään vielä yksinkertaisempaa:

Bigskip

3 Matriisit ja niiden havainnollistaminen

Drawmatrix toimii, mutta vaatii säätöä. Voisi olla leipätekstissä hyvä matriisiyhtälöiden havainnollistamiseen.

Kummallinen kohdistus.

$$\left(\begin{array}{|c|} \hline A \\ \hline \end{array} \begin{array}{|c|} \hline B \\ \hline \end{array}^{-1} \right) \begin{array}{|c|} \hline C \\ \hline \end{array}$$

Taulukko 3: Alle kouluikäinen lapsi todennäköisesti kärsii, jos hänen äitinsä käy työssä

	V6					
	1	2	3	4	5	Total
	n = 4	n = 4	n = 4	n = 4	n = 4	n = 4
maa						
FI	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)
HU	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)
SE	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)
Total	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)	1 (25%)

Ainakin SVD - osuudessa voi hyödyntää tätä:

$$\left(\begin{array}{c|c|c|c} S & s & s & D \\ \hline & j & k-1 & k \end{array} \right) U$$

4 Matriisiyhtälöt ilman kaavioita

Tässä lähteenä Greenacren kirja (ca in practice) ja sen liite Theory of CA. Muistiinpanoja löytyy, joissa viitataan myös Biplots in practice - kirjaan. Ei valmis, lähinnä kaavojen kirjoittelun harjoittelua.

Korrespondenssianalyysin perusyhtälöt:

Datamatriisin \mathbf{N} alkiot ovat ei-negatiivisia (eli nollat sallittuja) ja samassa mitta-asteikossa (jos mitta-asteikko on intervalli- tai suhdeasteikko mittayksiköiden on oltava samoja), ja n on taulukon alkoiden summa. GDA-kirjassa on tarkennettu tätä vaatimusta ei-negatiivisuudesta.

Korrespondenssimatriisi \mathbf{P} saadaan jakamalla matriisin \mathbf{N} alkiot niiden summalla n (tai ehkä parempi merkintä N). Merkitään matriisin \mathbf{P} rivisummien vektoria \mathbf{r} ja sarakesummien vektoria \mathbf{c} . Korrespondenssianalyysin termein nämä vektorit ovat rivi- ja sarakemassojen vektoreita, ja niitä vastaavat diagonaalimatriisit ovat \mathbf{D}_r ja \mathbf{D}_c .

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

$$(1) \quad \mathbf{S} = \mathbf{D}_r^{-1/2}(\mathbf{P} - \mathbf{r}\mathbf{c}^T)\mathbf{D}_c^{-1/2}$$

tai

$$\mathbf{S} = \mathbf{D}_r^{1/2}(\mathbf{D}_r^{-1}\mathbf{P}\mathbf{D}_c^{-1} - \mathbf{1}\mathbf{1}^T)\mathbf{D}_c^{-1/2} \quad .$$

Toinen esitystapa on hyödyllinen, kun tarkastellaan CA:n yhteyksiä muihin läheisiin menetelmiin (erityisesti kai log ratio analysis of compositional data?). Ehkäpä siksi, että matriisin alkiolle elementtimuodossa saadaan vastaavasti kaksi esitystapaa. Ensimmäinen on

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

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 $\alpha_1 \geq \alpha_2 \geq \dots$.

Näin saadaan standardikoordinaatit ja principal-koordinaatit riveille ja sarakkeille.

Rivien standardikoordinaatit

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

Sarakkeiden standardikoordinaatit

$$(3) \quad \Gamma = \mathbf{D}_c^{-1/2} \mathbf{V}$$

Rivien principal-koordinaatit

$$(4) \quad \mathbf{F} = \mathbf{D}_r^{-1/2} \mathbf{U} \mathbf{D}_\alpha = \Phi \mathbf{D}_\alpha$$

Sarakkeiden principal-koordinaatit

$$(5) \quad \mathbf{G} = \mathbf{D}_c^{-1/2} \mathbf{V} \mathbf{D}_\alpha = \Gamma \mathbf{D}_\alpha$$