

Statistische modellen en data-analyse: project 1

Toon Nolten, r0258654
<toon.nolten@student.kuleuven.be>

April 2014

Inleiding

De toegewezen dataset die ik moest analyseren was "Turnhout.xlsx". Dit is een kleine dataset met een aantal gegevens over de gemeenten rond Turnhout, bvb.: aantal mannelijke inwoners, gemiddelde inkomens. Om de resultaten te verkrijgen waar dit verslag op gebaseerd is, moet de code uitgevoerd worden, wegens tijdsgebrek heb ik die niet in het verslag kunnen verwerken.

Opgave 1: Principaal component analyse

- 1 Voer een PCA analyse uit op jullie gegevens. Argumenteer waarom je de analyse baseert op de correlatie- of de covariantiematrix van de gegevens. Bepaal het aantal componenten dat je verkiest te behouden en verklaar je keuze.**

Ik heb gekozen voor de analyse op basis van de correlatie matrix omdat de schalen van de verschillende variabelen veel verschillen (aantallen, euros, procenten). Op basis van de summary en screeplot van de pca met correlatie matrix, behoudt ik 7 principale componenten, zo behouden we 98.6%, 99% van de informatie met ongeveer de helft van de variabelen.

- 2 Bekijk de loadings van de weerhouden componenten en tracht de PCs te interpreteren aan de hand van de loadings. Ga na welke originele variabelen het belangrijkste zijn voor elk van de weerhouden PCs.**

(afgeleid uit de loadings)

De eerste PC is ongeveer een gemiddelde waarbij de variabelen waarbij geteld wordt een groter gewicht hebben, deze variabelen komen in de andere

componenten nagenoeg niet meer voor. Belangrijkste variabele: Aantal overlijdens (meeste met zeer weinig marge)

De tweede component bevat immigratie, de verhouding jonge mensen, de kostprijs van huizen en de verhouding 80+/60+: Immigranten zijn meestal jongere mensen of ze hebben veel kinderen; jongere mensen en de prijs van huizen zijn ook sterk verbonden, zij kunnen zich immers meestal geen dure villa's veroorloven. Belangrijkste variabele: Gemiddelde verkoopprijs van woonhuizen

PC3 bevat voornamelijk gemiddeld inkomen, werkzaamheidsgraad, werkloosheid, grote belastingaangiften en immigratie (ook 80+/60+). Een verband tussen werkzaamheid, werkloosheid, aangiften en inkomens is te verwachten. Immigratie gebeurt misschien vooral daar waar de lonen gunstig zijn. Belangrijkste variabele: Gemiddeld inkomen per aangifte

PC4 bestaat uit immigratie, verhouding jonge mensen, inkomen, werkloosheid en aantal geboorten. Nogmaals het verband tussen immigranten en jongeren, nu ook aantal geboorten. En ook, werkloosheid en inkomen. Belangrijkste variabele: Immigratie saldo

PC5 immigratie, 80+/60+, werkzaamheid, verkoopprijs, inkomen. Een verband tussen werkzaamheid, inkomen en verkoopprijs van huizen is te verwachten. Deze zijn ook gerelateerd aan de verhouding bejaarden en een beetje aan immigratie. Belangrijkste variabele: Werkzaamheidsgraad

PC6 immigratie, jonge mensen, inkomen, werkzaamheid, werkloosheid en verkoopprijs. Opnieuw het verband jonge mensen, immigratie. Ook, inkomen, verkoopprijs en immigratie, werkzaamheid, werkloosheid. Belangrijkste variabele: Gemiddelde verkoopprijs van woonhuizen

PC7 werkloosheid, werkzaamheid, inkomen, jonge mensen. Opnieuw jonge mensen, werkloosheid, werkzaamheid en inkomen. Belangrijkste variabele: Werkloosheidsgraad

3 Onderzoek de impact van de centrumstad op het resultaat van de PCA analyse.

(afgeleid uit biplots tussen verschillende PCs)

De eerste principale component wordt gedicteerd door de centrumstad, Turnhout. De andere componenten worden nagenoeg niet beïnvloed door de centrumstad. Dit is te verklaren aan de hand van de samenstelling van de PCs: De eerste PC bevat aantallen en die zijn voor een grote centrumstad groter, de andere PCs bestaan vooral uit verhoudingen en gemiddelden waarop de absolute hoeveelheid inwoners weinig invloed heeft.

Opgave 2: Testen van de multivariate normaliteit

- 1 **Onderzoek of de veronderstelling van multivariate normaliteit aannemelijk zou zijn voor jullie gegevens. Gebruik hiervoor ook de scores van de PCA analyse.**

Univariate marginalen: het aantal mannen, vrouwen, geboorten, overlijdens en kleine belastingaangiften is licht rechtsscheef. De andere marginalen komen redelijk goed overeen met de normale verdeling.

Bivariate marginalen: een aantal variabelen vertoont een sterke correlatie en bij een aantal variabelen zijn er mogelijk uitschieters maar over het algemeen zien de bivariate verdelingen er redelijk normaal verdeeld uit.

Mahalanobis afstanden lijken redelijk chi-kwadraat verdeeld. Wijken niet genoeg af van de rechte om te veronderstellen dat ze niet chi-kwadraat verdeeld zijn. De horizontale bij $\chi^2_{0.025}$ ligt onder alle punten, dit is niet wat ik verwachtte te zien en ik ben niet zeker waarom dat zo is.

- 2 **Als de multivariate normaliteit niet aannemelijk zou blijken te zijn, onderzoek dan of je de data kan aanpassen om beter aan deze veronderstelling te voldoen.**

De multivariate normaliteit lijkt wel aannemelijk.

Appendix: Code

```

1 # Voor deze opgave analyseren jullie gegevens afkomstig van de Vereniging
2 # van Vlaamse Steden en Gemeenten (http://www.vvsg.be). Deze vereniging
3 # brengt gegevens uit verschillende databanken samen en levert uitgebreide
4 # informatie over demografische en economische indicatoren per gemeente,
5 # regio, provincie, etc. Binnen het Vlaamse gewest worden 13 centrumsteden
6 # herkend, nl. Aalst, Antwerpen, Brugge, Genk, Gent, Hasselt,
7 # Kortrijk, Leuven, Mechelen, Oostende, Roeselare, Sint-Niklaas en Turnhout.
8 # Elke groep analyseert de gegevens van een regio bestaande uit een
9 # centrumstad en zijn omliggende gemeenten.
10
11 # De datasets bevatten de meest recente gegevens (afkomstig van 2011 of 2012)
12 # voor de volgende variabelen:
13
14 #+-----+
15 #|          | Totaal aantal mannen          |
16 #|          | Totaal aantal vrouwen        |
17 #|          | Aantal geboorten             |
18 #| Demografisch | Aantal overlijdens           |
19 #|          | Immigratie saldo (Inwijkelingen-uitwijkelingen) |
20 #|          | 80+/60+ verhouding (in%)     |
21 #|          | (0-19)/totaal verhouding (in %) |
22 #+-----+
23 #|          | Aantal belastingaangiften < 10.000 euro |
24 #|          | Aantal belastingaangiften > 50.000 euro |
25 #| Economisch  | Gemiddeld inkomen per aangifte |
26 #|          | werkzaamheidsgraad (% werkenden van (18-64+) bevolking) |
27 #|          | werkloosheidsgraad (% werkzoekenden van (18-64+) bevolking) |
28 #|          | Gemiddelde verkoopprijs van woonhuizen |
29 #+-----+
30
31 # De groepen worden samengesteld tijdens het practicum.
32 # De datasets worden als volgt verdeeld:
33
34 #      +-----+-----+
35 #      | Groep | Dataset |
36 #      +-----+-----+
37 #      | 1     | Aalst   |
38 #      | 2     | Gent    |
39 #      | 3     | Hasselt |
40 #      | 4     | Kortrijk|
41 #      | 5     | Leuven  |
42 #      | 6     | Mechelen|
43 #      | 7     | Sint-   |
44 #      |      | Niklaas|
45 # ==>| 8     | Turnhout|
46 #      +-----+-----+

```

```

48 library(gdata)
49 turnhout = read.xls(xls='Turnhout.xlsx', sheet=1, pattern='Mannen')
50 rownames(turnhout) = turnhout[,1]
51 turnhout = turnhout[,-1]
52 colnames(turnhout) = c('Mannen',
53                        'Vrouwen',
54                        'n_geboorten',
55                        'n_overlijdens',
56                        'Immigratie_saldo',
57                        '80+/60+_(in%)',
58                        '(0-19)/totaal_(in%)',
59                        'Aantal_aangiften_<_10.000_euro',
60                        'Aantal_aangiften_>_50.000_euro',
61                        'Gemiddeld_inkomen_per_aangifte',
62                        'Werkzaamheidsgraad',
63                        'Werkloosheidsgraad',
64                        'Gemiddelde_verkoop prijs_van_woonhuizen')
65 for(column in 1:ncol(turnhout)) {
66     turnhout[,column] = as.numeric(as.character(turnhout[,column]))
67 }
68
69 cat('\n', 'Summary_turnhout.xlsx:', '\n',
70     '_____', '\n')
71 print(summary(turnhout))
72
73 # Opgave 1: PCA
74 # 1. Voer een PCA analyse uit op jullie gegevens. Argumenteer waarom je
75 # de analyse baseert op de correlatie- of de covariantiematrix van de
76 # gegevens. Bepaal het aantal componenten dat je verkiest te behouden
77 # en verklaar je keuze.
78
79 pca_cov_turnhout = precomp(turnhout, scale=FALSE)
80 pca_cor_turnhout = precomp(turnhout, scale=TRUE)
81
82 cat('\n', 'Pca_on_covariance_matrix:', '\n',
83     '_____', '\n')
84 print(summary(pca_cov_turnhout))
85 #screeplot(pca_cov_turnhout, type='lines')
86
87 cat('\n', 'Pca_on_correlation_matrix:', '\n',
88     '_____', '\n')
89 print(summary(pca_cor_turnhout))
90 #screeplot(pca_cor_turnhout, type='lines')
91
92 ### Kies de correlatie matrix omdat de schalen van de verschillende
93 ### variabelen veel verschillen (aantallen, euros, procenten)
94 ### Behoud 7 principale componenten, zo behouden we 98.6, ~99% van de
95 ### informatie met ~1/2 van de variabelen.

```

```

97 # 2. Bekijk de loadings van de weerhouden componenten en tracht de PCs
98 # te interpreteren aan de hand van de loadings. Ga na welke originele
99 # variabelen het belangrijkste zijn voor elk van de weerhouden PCs.
100
101 cat('\n_', 'Loadings_pca_on_correlation_matrix:', '\n_',
102      '\n_', '\n')
103 print(pca_cor_turnhout$rotation[,1:7])
104
105 #### De eerste PC is ongeveer een gemiddelde waarbij de variabelen waarbij
106 #### getelt wordt een groter gewicht hebben.
107 #### Aantal overlijdens (meeste met zeer weinig marge)
108 #### De tweede component bevat immigratie, de verhouding jonge mensen,
109 #### de kostprijs van huizen en de verhouding 80+/60+:
110 #### Immigranten zijn meestal jongere mensen of ze hebben veel kinderen;
111 #### jongere mensen en de prijs van huizen zijn ook sterk verbonden, zij
112 #### kunnen zich immers meestal geen dure villa's veroorloven.
113 #### Gemiddelde verkoopprijs van woonhuizen
114 #### PC3 bevat voornamelijk gemiddeld inkomen, werkzaamheidsgraad,
115 #### werkloosheid, grote belastingaangiften en immigratie (ook 80+/60+).
116 #### Een verband tussen werkzaamheid, werkloosheid, aangiften en inkomens is
117 #### te verwachten. Immigratie gebeurt misschien vooral daar waar de lonen
118 #### gunstig zijn.
119 #### Gemiddeld inkomen per aangifte
120 #### PC4 bestaat uit immigratie, verhouding jonge mensen, inkomen,
121 #### werkloosheid en aantal geboorten.
122 #### Nogmaals het verband tussen immigranten en jongeren, nu ook aantal
123 #### geboorten. En ook, werkloosheid en inkomen.
124 #### Immigratie saldo
125 #### PC5 immigratie, 80+/60+, werkzaamheid, verkoopprijs, inkomen.
126 #### Een verband tussen werkzaamheid, inkomen en verkoopprijs van huizen is te
127 #### verwachten. Deze zijn ook gerelateerd aan de verhouding bejaarden en
128 #### een beetje aan immigratie.
129 #### Werkzaamheidsgraad
130 #### PC6 immigratie, jonge mensen, inkomen, werkzaamheid, werkloosheid en
131 #### verkoopprijs. Opnieuw het verband jonge mensen, immigratie. Ook, inkomen,
132 #### verkoopprijs en immigratie, werkzaamheid, werkloosheid.
133 #### Gemiddelde verkoopprijs van woonhuizen
134 #### PC7 werkloosheid, werkzaamheid, inkomen, jonge mensen.
135 #### Opnieuw jonge mensen, werkloosheid, werkzaamheid en inkomen.
136 #### Werkloosheidsgraad

```

```

138 # 3. Onderzoek de impact van de centrumstad op het resultaat van de PCA
139 # analyse.
140
141 biplotmatrix = function (x, y, ...){
142   par(new=TRUE)
143   biplot(pca_cor_turnhout, c(x,y), pc.biplot=TRUE)
144 }
145
146 #pairs(t(seq(7)), panel=biplotmatrix)
147
148 #### De eerste principaal component wordt gedicteerd door de centrumstad,
149 #### Turnhout. De andere componenten worden nagenoeg niet beïnvloed door de
150 #### centrumstad. Dit is te verklaren aan de samenstelling van de PCs:
151 #### De eerste PC bevat aantallen en die zijn voor een grote centrumstad
152 #### groter, de andere PCs bestaan vooral uit verhoudingen en gemiddelden
153 #### waarop de absolute hoeveelheid inwoners weinig invloed heeft.
154
155 # Opgave 2: Testen multivariate normaliteit
156 # 1. Onderzoek of de veronderstelling van multivariate normaliteit
157 # aannemelijk zou zijn voor jullie gegevens. Gebruik hiervoor ook de
158 # scores van de PCA analyse.
159
160 #par(mfrow=c(3,5))
161 #qqnorm(turnhout[,1])
162 #qqline(turnhout[,1])
163 #qqnorm(turnhout[,2])
164 #qqline(turnhout[,2])
165 #qqnorm(turnhout[,3])
166 #qqline(turnhout[,3])
167 #qqnorm(turnhout[,4])
168 #qqline(turnhout[,4])
169 #qqnorm(turnhout[,5])
170 #qqline(turnhout[,5])
171 #qqnorm(turnhout[,6])
172 #qqline(turnhout[,6])
173 #qqnorm(turnhout[,7])
174 #qqline(turnhout[,7])
175 #qqnorm(turnhout[,8])
176 #qqline(turnhout[,8])
177 #qqnorm(turnhout[,9])
178 #qqline(turnhout[,9])
179 #qqnorm(turnhout[,10])
180 #qqline(turnhout[,10])
181 #qqnorm(turnhout[,11])
182 #qqline(turnhout[,11])
183 #qqnorm(turnhout[,12])
184 #qqline(turnhout[,12])
185 #qqnorm(turnhout[,13])
186 #qqline(turnhout[,13])

```

```

188 ### Univariate marginalen: het aantal mannen, vrouwen, geboorten, overlijdens
189 ### en kleine belastingaangiften is licht rechtsscheef. De andere marginalen
190 ### komen redelijk goed overeen met de normale verdeling.
191
192 #pairs(turnhout)
193
194 ### Bivariate marginalen: een aantal variabelen vertoont een sterke correlatie
195 ### en bij een aantal variabelen zijn er mogelijk uitschieters maar over het
196 ### algemeen zien de bivariate verdelingen er redelijk normaal verdeeld uit.
197
198 #y = mahalanobis(turnhout, center=sapply(turnhout, mean), cov=cov(turnhout))^2
199 #dof=13
200 #qqplot(qchisq(ppoints(500), df= dof), y,
201 #           main = expression("Q-Q plot for" ~{chi^2}[nu == dof]))
202 #qqline(y, distribution = function(p) qchisq(p, df = dof),
203 #           prob = c(0.25, 0.75), col = 2)
204 #abline(h=qchisq(p=0.975, df=dof))
205 #mtext("qqline(*, dist = qchisq(., df=13), prob = c(0.1, 0.6))")
206
207 ### Mahalanobis afstanden lijken redelijk normaal verdeeld. Wijken niet
208 ### genoeg af om anormaliteit te veronderstellen.
209 ### De horizontale by  $\chi^2_{0.025}$  ligt onder alle punten, niet zekere wat
210 ### dat betekent.
211
212 # 2. Als de multivariate normaliteit niet aannemelijk zou blijken te zijn,
213 #      onderzoek dan of je de data kan aanpassen om beter aan deze
214 #      veronderstelling te voldoen.
215
216
217
218
219
220 # Opgave 3: Clustering
221 # 1. Voer een cluster analyse uit op de gegevens. Kies op basis van de
222 #      resultaten volgens verschillende clustermethodes een ‘‘optimaal’’
223 #      aantal clusters.
224
225
226
227 # 2. Stel de clusterresultaten grafisch voor en bespreek de kwaliteit van
228 #      de gevonden clustering. Bespreek gelijkenissen en/of verschillen tussen
229 #      de methoden.
230
231
232
233 # 3. Je kan ook een clustering uitvoeren vertrekkende van de scores van de
234 #      PCA analyse. Zou dit zinvol kunnen zijn of heeft een voorafgaande PCA
235 #      analyse weining zin? Levert dit een ander resultaat voor jullie
236 #      gegevens?

```