

Statistische modellen en data-analyse: project 3

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

Juni 2014

Inleiding

De dataset “Regios.xlsx” bevat data over 180 steden en gemeenten. De opgave bestaat erin een regressiemodel op te stellen om te voorspellen of een observatie een stad of gemeente is. Dit model wordt dan vergeleken met een discriminant analyse voor classificatie.

Ik heb één waarde gecorrigeerd in de dataset, namelijk voor het ‘Migratie saldo’ van de gemeente Meerhout, dit was 23 wat erg opvalt in vergelijking met de andere waarden, dit moest dan ook verbeterd worden naar 0.23. (Thomas Neven heeft hierop gewezen.)

Opgave 1: Gebruik een regressiemodel om de respons ‘Stad’ te verklaren aan de hand van de andere beschikbare variabelen.

Omdat de respons een binaire variabele is heb ik logistische regressie gebruikt om een model op te stellen.

Opgave 2: Geeft dit regressiemodel een aanvaardbare fit?

De residual deviance van het model is 114.85, dit is kleiner dan $\chi^2_{167,0.05}$ dus geeft het model een goede fit.

Opgave 3: Is de variabele ‘Regio’ nuttig in het regressiemodel of kan deze variabele beter weggelaten worden?

De likelihood ratio test toont aan dat de variabele Regio niet belangrijk is in het model. $Partialdeviance = 0.0079 < \chi^2_{1,0.95} = 3.84$ dus we aanvaarden H_0 .

Opgave 4: Zijn er invloedrijke observaties of andere afwijkingen die best opgelost worden om een beter model te bekomen?

Uit de plot van de deviance residuals, figuur 1, is duidelijk dat de enige uitschieter ‘Mesen’ is.

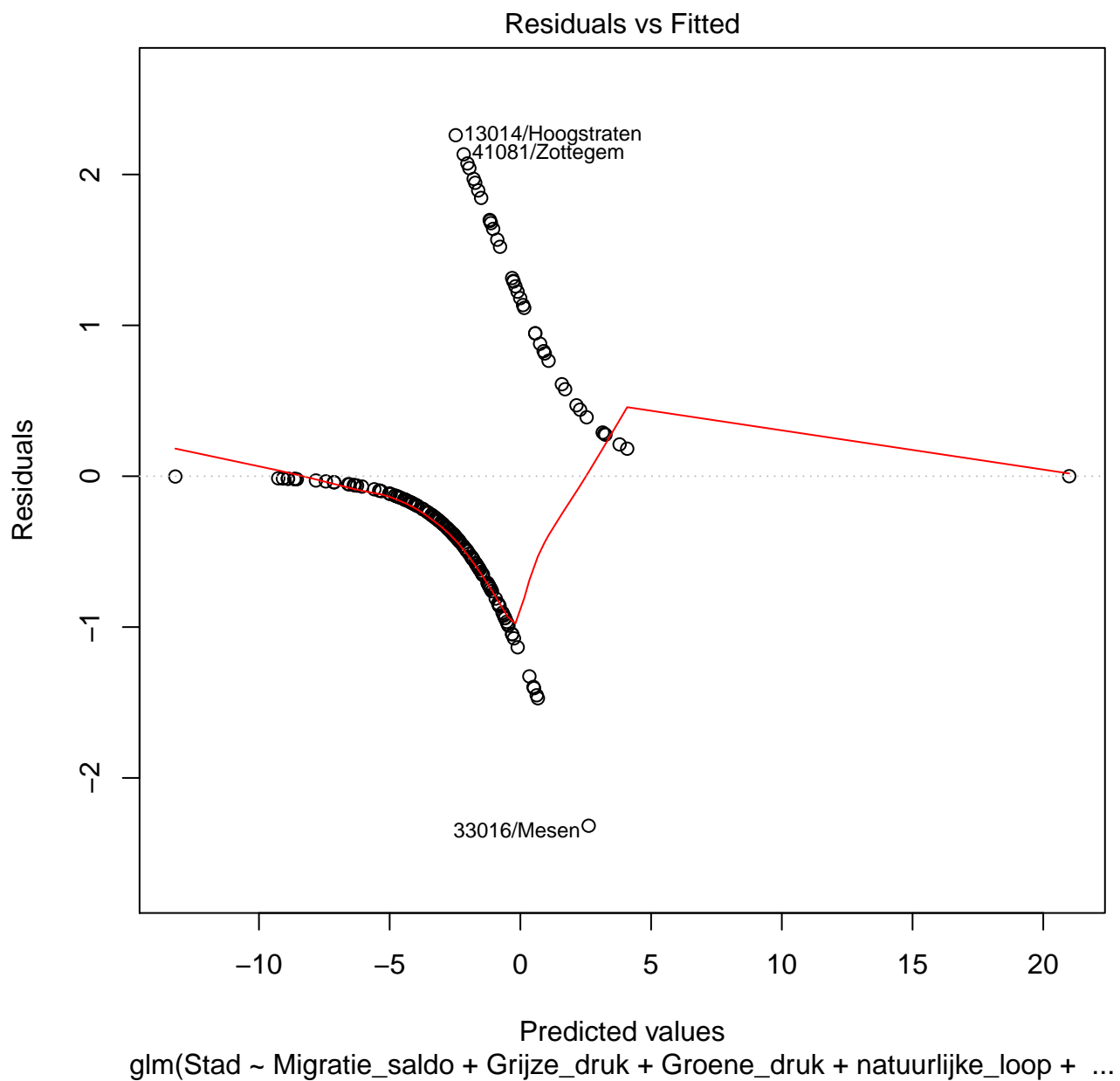


Figure 1: Plot of Deviance Residuals

Opgave 5: Hoe goed kan je met dit model steden van gemeenten onderscheiden?

Het model heeft een lage (14.4%) prediction error, tabel 1. Van de 180 observaties worden 4 gemeenten geclassificeerd als steden en 22 steden geclassificeerd als gemeentes. Dit in de aanname dat de misclassificatie voor beiden dezelfde kost heeft (threshold 0.5).

Voorspelling	Gemeente	Stad
Gemeente	137	4
Stad	22	17

Table 1: Classification table for Logistic Regression

Opgave 6: Via discriminantanalyse kan je ook steden van gemeenten proberen te onderscheiden. Levert dit een beter resultaat op dan met voorgaand model?

Het levert een vergelijkbaar resultaat op, tabel 2, prediction error 15.6%. Maar de threshold voor classificatie bij de logistische regressie is momenteel niet afgesteld op de data en kan mogelijk nog voor verbetering zorgen.

Voorspelling	Gemeente	Stad
Gemeente	136	5
Stad	23	16

Table 2: Classification table for Linear Discriminant Analysis

Appendix: Code

```

1 # De dataset regios.xlsx bevat voor de steden en gemeenten uit 8
2 # Vlaamse regio 's gegevens voor de volgende variabelen in het jaar
3 # 2010:
4
5 # +-----+-----+
6 # |          | Relatief migratie saldo (% van totale bevolking) |
7 # | Demografisch | Grijze druk (% 65+ t.o.v. 20-64) |
8 # |          | Groene druk (% -20 t.o.v. 20-64) |
9 # |          | natuurlijke loop (+1 bij aangroei, -1 bij daling) |
10 # +-----+-----+
11 # |          | % belastingaangiften > 50.000 euro |
12 # |          | Gemiddeld inkomen per aangifte |
13 # |          | werkzaamheidsgraad (% werkenden van (18-64) |
14 # | Economisch | bevolking) |
15 # |          | werkloosheidsgraad (% werkzoekenden van (18-64) |
16 # |          | bevolking) |
17 # |          | Gemiddelde verkoopprijs van woonhuizen |
18 # +-----+-----+
19 # |          | % leefloners t.o.v. totale bevolking |
20 # | Armoede | % geboorten in kansarme gezinnen |
21 # |          | % bejaarden met inkomensgarantie (i.e. uitkering) |
22 # +-----+-----+
23 # | Geografisch | Regio |
24 # |          | Stad |
25 # +-----+-----+
26
27 # De variabele 'Regio' is een indicator die aangeeft welke steden en
28 # gemeenten tot eenzelfde regio behoren. De variabele 'Stad' geeft
29 # weer of de observatie een stad (1) of gemeente (0) is.
30 # Je beantwoordt de onderstaande vragen door individueel gepaste
31 # analyses uit te voeren met R. De bespreking van de resultaten en
32 # de nodige figuren verwerk je in een schriftelijk rapport dat
33 # maximaal uit 6 bladzijden mag bestaan (12pt lettergrootte). In een
34 # appendix voeg je de gebruikte R-code toe. Rapporteer enkel
35 # resultaten en interpretaties, herhaal geen theorie uit de cursus!
36 # Het rapport dien je in pdf vorm in ten laatste op *12 juni 2014*
37 # via Toledo. Dit project telt mee voor 5 punten van het
38 # eindresultaat.

```

```

40 # Opgaven:
41 # Het doel is om te onderzoeken of steden van gemeenten kunnen
42 # onderscheiden worden op basis van de beschikbare variabelen.
43
44 library(gdata)
45 turnhout = read.xls(xls='Regios.xlsx', sheet=1, pattern='druk')
46 rownames(turnhout) = turnhout[,1]
47 turnhout = turnhout[,-1]
48 colnames(turnhout) = c('Migratie_saldo',
49                        'Grijze_druk',
50                        'Groene_druk',
51                        'natuurlijke_loop',
52                        '%_aangiften_>_50.000_euro',
53                        'Gemiddeld_inkomen_per_aangifte',
54                        'werkzaamheidsgraad',
55                        'werkloosheidsgraad',
56                        '%_leefloners',
57                        '%_geboorten_in_kansarme_gezinnen',
58                        '%_bejaarden_met_inkomensgarantie',
59                        'Gemiddelde_verkoopprijs_van_woonhuizen',
60                        'Regio',
61                        'Stad')
62 for(column in 1:ncol(turnhout)) {
63     turnhout[,column] = as.numeric(as.character(turnhout[,column]))
64 }
65
66 cat('\n', 'Summary_Regios.xlsx:', '\n',
67     '_____', '\n')
68 print(summary(turnhout))

```

```

70 # 1. Gebruik een regressiemodel om de respons 'Stad' te verklaren
71 #   aan de hand van de andere beschikbare variabelen.
72
73 turnhout.glm = glm(Stad
74                    ~Migratie_saldo
75                    +Grijze_druk
76                    +Groene_druk
77                    +natuurlijke_loop
78                    +'%_aangiften_>_50.000_euro'
79                    +Gemiddeld_inkomen_per_aangifte
80                    +werkzaamheidsgraad
81                    +werkloosheidsgraad
82                    +Gemiddelde_verkoop prijs_van_woonhuizen
83                    +'%_leefloners'
84                    +'%_geboorten_in_kansarme_gezinnen'
85                    +'%_bejaarden_met_inkomensgarantie'
86                    +Regio, family=binomial, data=turnhout)
87 cat('\n', 'Summary_turnhout.glm:', '\n',
88     '_____', '\n')
89 print(summary(turnhout.glm))
90
91 cat('\n', 'Anova_turnhout.glm:', '\n',
92     '_____', '\n')
93 print(anova(turnhout.glm, test="Chisq"))
94
95
96 # 2. Geeft dit regressiemodel een aanvaardbare fit?
97
98 cat('\n', 'Acceptable_fit', '\n',
99     '_____', '\n')
100 cat('\n', 'Chi^2_167,0.05:', 'qchisq(0.05, 167), '\n') # = 138.1184
101
102 # Ja,  $\chi^2_{167,0.05} = 138.1184$ , de model deviance is 114.85 dus
103 # we aanvaarden het model.
104
105 # 3. Is de variabele 'Regio' nuttig in het regressiemodel of kan
106 #   deze variabele beter weggelaten worden?
107
108 turnhout.noregio = update(turnhout.glm, ~.-Regio)
109
110 cat('\n', 'Likelihood_ratio_test_for_Regio', '\n',
111     '_____', '\n')
112 cat('\n', 'Partial_Deviance:',
113     turnhout.noregio$deviance - turnhout.glm$deviance, '\n')
114 cat('\n', 'Chi^2_1,0.95:', 'qchisq(0.95, 1), '\n')
115
116 # De likelihood ratio test toont aan dat de variabele Regio niet
117 # belangrijk is in het model.  $0.0079 < 3.84$  dus we aanvaarden  $H_0$ .

```

```

120 # 4. Zijn er invloedrijke observaties of andere afwijkingen die best
121 #     opgelost worden om een beter model te bekomen?
122
123 ##pdf('deviance_residuals.pdf')
124 #plot(turnhout.noregio, which=1) # residuals plot
125 ##dev.off()
126
127 # Uit de plot van de deviance residuals is duidelijk dat de enige
128 # uitschieter 'Mesen' is.
129
130 # 5. Hoe goed kan je met dit model steden van gemeenten
131 #     onderscheiden?
132
133 logitpred = predict(turnhout.noregio) > 0.5
134 cat('\n', 'Classification_table:', '\n',
135     '_____', '\n')
136 print(table(turnhout$Stad, logitpred))
137 cat('\n', 'Prediction_error:', '\n',
138     round((22+4)/(137+4+22+17)*100, 2), '%\n')
139
140 # Het model heeft een lage (14.4%) prediction error. Van de 180
141 # observaties worden 4 gemeenten geclassificeerd als steden en 22
142 # steden geclassificeerd als gemeentes. Dit in de aanname dat
143 # de misclassificatie voor beiden dezelfde kost heeft
144 # (threshold 0.5) .

```

```

146 # 6. Via discriminantanalyse kan je ook steden van gemeenten
147 #     proberen te onderscheiden. Levert dit een beter resultaat op dan
148 #     met voorgaand model?
149
150 library(MASS)
151 turnhout.lda = lda(Stad
152                    ~Migratie_saldo
153                    +Grijze_druk
154                    +Groene_druk
155                    +natuurlijke_loop
156                    +'%_aangiften_>_50.000_euro'
157                    +Gemiddeld_inkomen_per_aangifte
158                    +werkzaamheidsgraad
159                    +werkloosheidsgraad
160                    +Gemiddelde_verkoop prijs_van_woonhuizen
161                    +'%_leefloners'
162                    +'%_geboorten_in_kansarme_gezinnen'
163                    +'%_bejaarden_met_inkomensgarantie'
164                    +Regio, data=turnhout)
165
166 cat('\n', 'Summary_turnhout.lda:', '\n',
167     '_____', '\n')
168 print(turnhout.lda)
169
170 ldapred = predict(turnhout.lda, turnhout)$class
171 cat('\n', 'Classification_table:', '\n',
172     '_____', '\n')
173 print(table(turnhout$Stad, ldapred))
174 cat('\n', 'Prediction_error:',
175     round((23+5)/(137+4+22+17)*100, 2), '%\n')
176
177 # Het levert een vergelijkbaar resultaat op, prediction error 15.6%.
178 # Maar de threshold voor classificatie bij de logistische regressie
179 # is momenteel niet afgesteld op de data en kan mogelijk nog
180 # voor verbetering zorgen.

```