

Statistische modellen en data-analyse: project 3

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

September 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.)

1 Opgaven

1.1 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. Dit heb ik gedaan door eerst een model op te stellen met alle variabelen als predictors:

```
glm(Stad ~., family=binomial, data=turnhout)
```

Op dat model heb ik dan een stapsgewijze verbetering uitgevoerd gebaseerd op de *AIC*:

```
step(turnhout.glm)
```

Dit levert het volgende model op:

```
glm(Stad ~ Grijsze_druk + Groene_druk + Gemiddeld_inkomen_per_aangifte  
+ `%-leefloners`, family=binomial, data=turnhout)
```

Als we ook interactiecomponenten in rekening brengen kunnen we het model nog verbeteren maar dit gaat ten koste aan de eenvoud en de interpreteerbaarheid van het model en levert in dit geval onvoldoende verbetering ($\pm 5\%$ minder classificatiefouten).

(Er bestaan betere modelselectiemethoden maar die laten we hier buiten beschouwing, cf. het R pakket ‘glmnet’)

1.2 Geeft dit regressiemodel een aanvaardbare fit?

Dit gaan we na met een ‘Goodness of fit’ test waarbij de nul-hypothese betekent dat we het model aanvaarden. De residual deviance van het model is 118.57, dit is kleiner dan $\chi^2_{175,0.05} = 145.4058$ dus op een betrouwbaarheidsniveau van 5% besluiten we dat het model een goede fit geeft.

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

De variabele ‘Regio’ zit niet in het model, ik heb dus het model dat gebruik maakt van alle predictor variabelen opgesteld om deze vergelijking op te doen.

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 , de variabele is niet nuttig.

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

Op de plot van de deviance residuals, figuur 1, is de enige duidelijke uitschieter ‘Mesen’. In dit geval betekent dit dat ‘Mesen’ een twijfelgeval is. Hiervoor is er niet echt een oplossing als we een model willen opstellen voor alle steden en gemeenten in België, omdat we niet kunnen bepalen of dit een echte uitschieter is of dat de steekproef een bias heeft.

1.5 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 3 gemeenten geclassificeerd als steden en 23 steden geclassificeerd als gemeentes. Dit in de aanname dat de misclassificatie voor beiden dezelfde kost heeft (threshold 0.5). Wat wel opvalt is dat de meeste steden verkeerd ingedeeld zijn, dit is omdat er veel meer gemeentes dan steden in de steekproef aanwezig zijn, daardoor worden nieuwe observaties sneller ingedeeld als gemeentes.

Voorspelling	Gemeente	Stad
Gemeente	138	3
Stad	23	16

Table 1: Classification table for Logistic Regression

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

Het gebruikte model is gebaseerd op dezelfde variabelen als in het verbeterde logistische model:

```
lda(Stad ~ Grijs_druk + Groene_druk + Gemiddeld_inkomen_per_aangifte  
+ `_%_leefloners`, data = turnhout)
```

Het levert een vergelijkbaar resultaat op, tabel 2, prediction error 15%. 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	140	1
Stad	26	13

Table 2: Classification table for Linear Discriminant Analysis

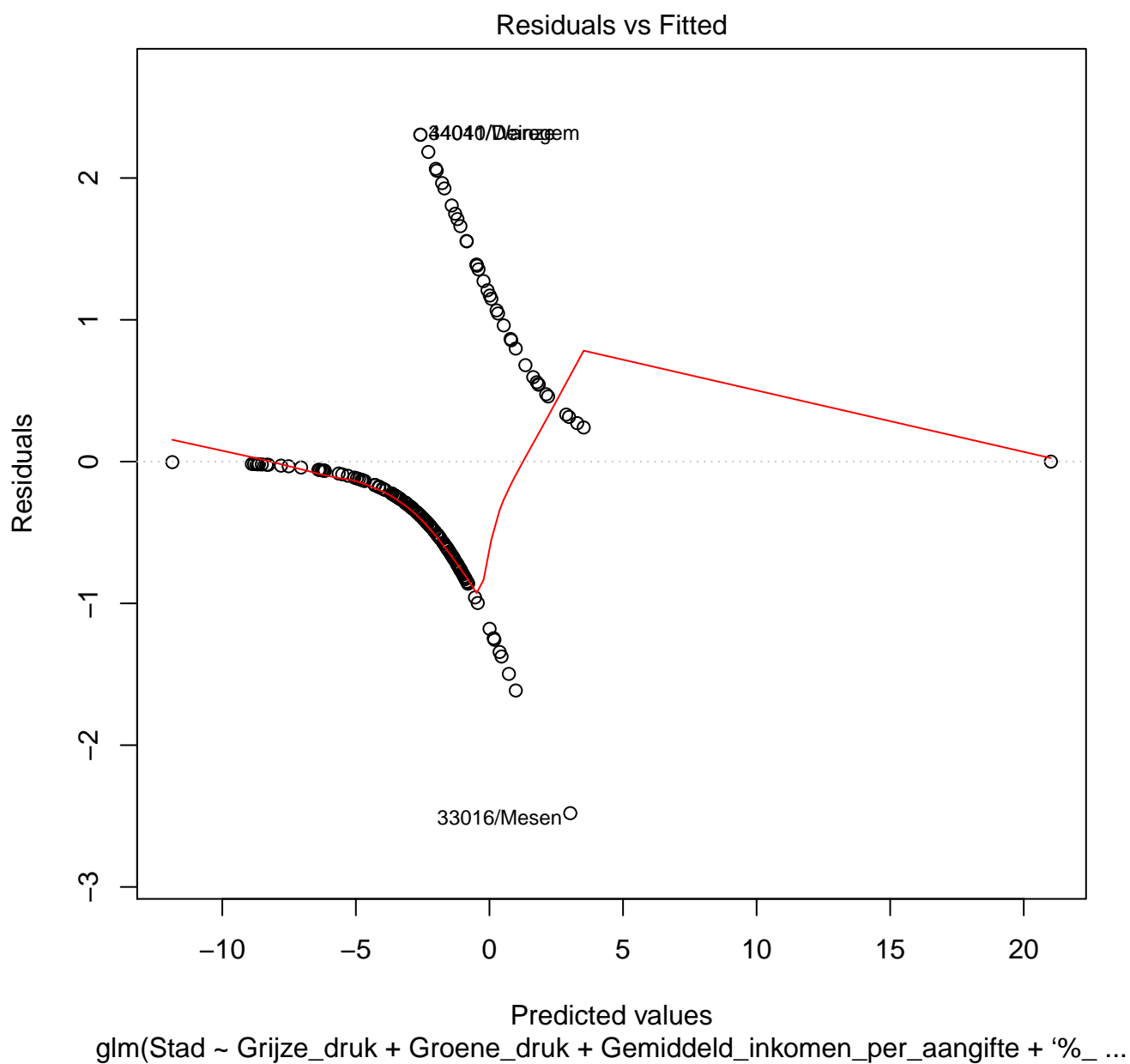


Figure 1: Plot of Deviance Residuals

2 Appendix: Code

```
1 # De dataset regio's.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 # /          / % belastingsaangiften > 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 ~., family=binomial, data=turnhout)
74 turnhout.glms = step(turnhout.glm, ~.)
75
76 cat('\n', 'Summary_turnhout.glms:', '\n',
77     '-----', '\n')
78 print(summary(turnhout.glms))
79
80 cat('\n', 'Anova_turnhout.glms:', '\n',
81     '-----', '\n')
82 print(anova(turnhout.glms, test="Chisq"))
83
84
85 # 2. Geeft dit regressiemodel een aanvaardbare fit?
86
87 cat('\n', 'Acceptable_fit', '\n',
88     '-----', '\n')
89 cat('_', 'Chi^2_175,0.05:', qchisq(0.05, 175), '\n') # = 145.4058
90
91 # Ja, chi^2_175,0.05 = 145.4058, de model deviance is 118.57 dus
92 # we aanvaarden het model.
93
94 # 3. Is de variabele `Regio' nuttig in het regressiemodel of kan
95 #   deze variabele beter weggelaten worden?
96
97 turnhout.noregio = update(turnhout.glm, .~-Regio)
98
99 cat('\n', 'Likelihood_ratio_test_for_Regio', '\n',
100    '-----', '\n')
101 cat('_', 'Partial_Deviance:',
102     turnhout.noregio$deviance - turnhout.glm$deviance, '\n')
103 cat('_', 'Chi^2_1,0.95:', qchisq(0.95, 1), '\n')
104
105 # De likelihood ratio test toont aan dat de variabele Regio niet
106 # belangrijk is in het model. 0.0079 < 3.84 dus we aanvaarden H_0.
107
108
109 # 4. Zijn er invloedrijke observaties of andere afwijkingen die best
110 #   opgelost worden om een beter model te bekomen?
111
112 ##pdf('deviance_residuals.pdf')
113 #plot(turnhout.glms, which=1) # residuals plot
114 ##dev.off()
115
116 # Uit de plot van de deviance residuals is duidelijk dat de enige
117 # uitschieter `Mesen' is.

```

```

119 # 5. Hoe goed kan je met dit model steden van gemeenten
120 #     onderscheiden?
121
122 logitpred = predict(turnhout.glms) > 0.5
123 cat('\n', 'Classification table:', '\n',
124     '-----', '\n')
125 logt = table(turnhout$Stad, logitpred)
126 print(logt)
127 cat('\n', 'Prediction error:',
128     round((logt[1,2]+logt[2,1])/(sum(logt))*100, 2), '%\n')
129
130 # Het model heeft een lage (14.4%) prediction error. Van de 180
131 # observaties worden 3 gemeenten geclassificeerd als steden en 23
132 # steden geclassificeerd als gemeentes. Dit in de aanname dat
133 # de misclassificatie voor beiden dezelfde kost heeft
134 # (threshold 0.5) .
135
136 # 6. Via discriminantanalyse kan je ook steden van gemeenten
137 #     proberen te onderscheiden. Levert dit een beter resultaat op dan
138 #     met voorgaand model?
139
140 library(MASS)
141 turnhout.lda = lda(Stad
142                    ~Grijze_druk
143                    +Groene_druk
144                    +Gemiddeld_inkomen_per_aangifte
145                    +`%-leefloners`, data=turnhout)
146
147 cat('\n', 'Summary turnhout.lda:', '\n',
148     '-----', '\n')
149 print(turnhout.lda)
150
151 ldapred = predict(turnhout.lda, turnhout)$class
152 cat('\n', 'Classification table:', '\n',
153     '-----', '\n')
154 ldat = table(turnhout$Stad, ldapred)
155 print(ldat)
156 cat('\n', 'Prediction error:',
157     round((ldat[1,2]+ldat[2,1])/(sum(ldat))*100, 2), '%\n')
158
159 # Het levert een vergelijkbaar resultaat op, prediction error 15%.
160 # Maar de threshold voor classificatie bij de logistische regressie
161 # is momenteel niet afgesteld op de data en kan mogelijk nog
162 # voor verbetering zorgen.

```