

# Statistische gegevensanalyse Project

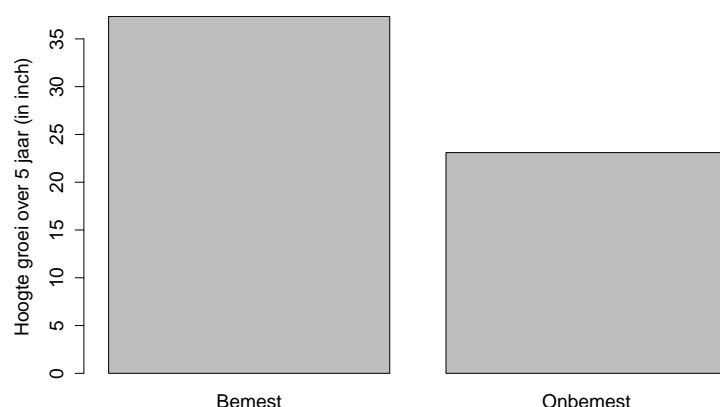
Titouan Vervack  
Bachelor of Science in Informatics

31 mei 2014

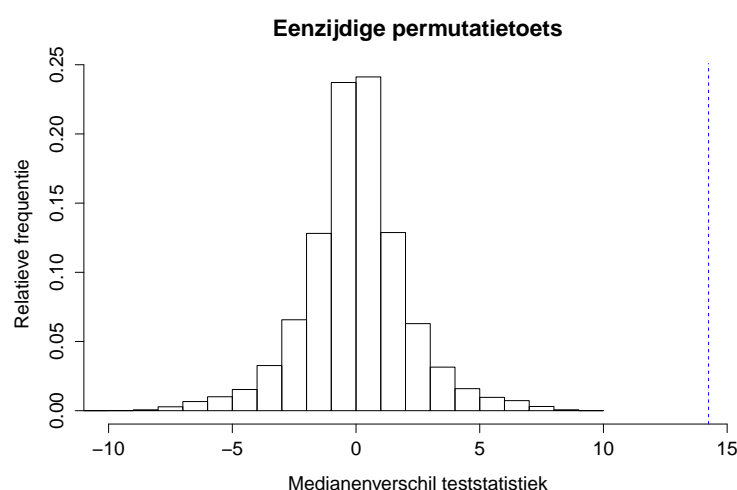
## Opgave 1

a

Om de groei in hoogte van de bomen te bekijken trekken we **Height0** af van **Height5**. Aangezien gevraagd wordt wat de invloed is van bemesting op deze hoogte, gebruiken we dit in combinatie met **Fertilizer**. Om een beeld te krijgen van de data maken we er histogrammen van, we kiezen histogrammen omdat de hoogtes continue variabelen zijn. We kunnen zien dat de data niet symmetrisch is en pieken bevat. We kiezen daarom om de mediaan te gebruiken in plaats van het gemiddelde. In figuur 1 zien we de staafdiagram die deze medianen voorstelt, we zien dat bemesting een positief effect heeft op de hoogte groei van bomen. Om zeker te zijn dat dit effect wordt veroorzaakt door de bemesting voeren we een permutatietoets uit, dit is een eenzijdige toets aangezien gevraagd wordt of de bomen sneller groeien bij bemesting.  $H_0$  stelt dat de groei bij geen of wel bemesting dezelfde is.  $H_A$  stelt dat de groei bij bemesting positief is. De p-waarde, verkregen uit de permutatietoets, is  $10^{-5}$ . Dit is een overtuigend bewijs tegen de nulhypothese waardoor we deze mogen verwerpen ten voordele van de alternatieve hypothese. Bemesting heeft dus een positief effect op de hoogte groei van bomen. De permutatieverdeling op basis van  $M = 99999$  is te zien in figuur 2.



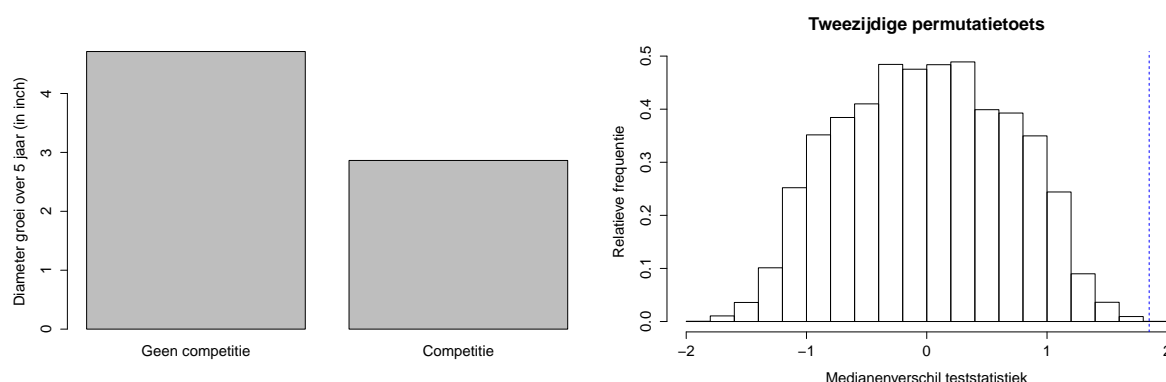
Figuur 1: Hoogte groei met bemesting of zonder bemesting



Figuur 2: Permutatieverdeling voor de hoogte groei van (on)bemeste bomen

b

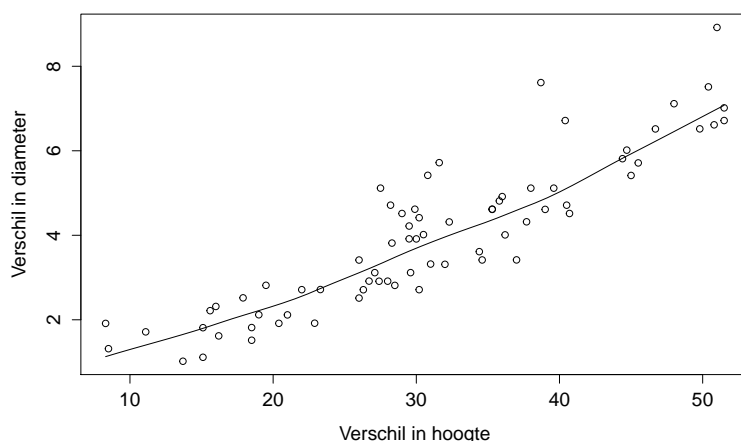
Bij deze opgave gebruiken we het verschil tussen **Diameter5** en **Diameter0**, we gebruiken ook **Competition**. We maken opnieuw histogrammen en beslissen om dezelfde reden terug de mediaan te gebruiken. De staafdiagram die de medianen vergelijkt is te zien in figuur 3a. Geen competitie blijkt een positief effect te hebben op de diameter van de bomen. We voeren ook hier een permutatietoets uit, dit is een tweezijdige aangezien het gevraagde effect niet gespecificeerd is. Om deze reden vermenigvuldigen we onze p-waarde met een factor 2.  $H_0$  stelt dat de groei bij geen of wel competitie gelijk is.  $H_A$  stelt dat de groei bij geen of wel competitie niet gelijk is. De nulhypothese wordt verworpen door een terug lage p-waarde ten voordele van de alternatieve hypothese. Competitie verwijderen is dus positief voor de diameter van de bomen. De permutatieverdeling op basis van  $M = 99999$  is te zien in figuur 3b.



Figuur 3: Barplot en permutatieverdeling van vraag 1b

**C**

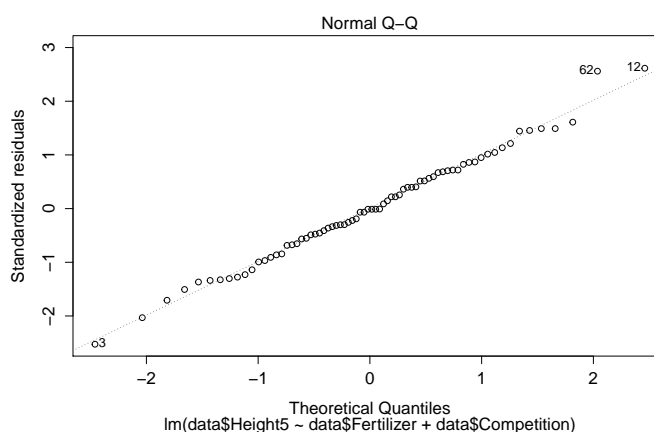
Hier gaan we de associatie tussen de verschillen uit opgaves **a** en **b** bekijken. Aangezien deze beide continue variabelen zijn drukken we de associatie uit aan de hand van de covariantie, deze is 17.82. We kunnen ook gemakkelijk de correlatie berekenen, deze is 0.90. Aan de hand van deze twee waarden kunnen we beslissen dat de associatie groot is. Aangezien we echter niet weten hoe representatief onze steekproef is kunnen we de associatie niet precies bepalen. We kunnen dit ook visualiseren in een scatterplot met een smoother curve, deze is te zien in figuur 4. Deze curve lijkt op een rechte, wat wil zeggen dat er een lineair verband is tussen de toename van de lengte en deze van de diameter. In combinatie met de hoge waarde voor de correlatie kunnen we zeggen dat er een sterk lineair verband is tussen de twee variabelen. Dit is echter data voor onze steekproef en niet voor de populatie.



Figuur 4: Lineair verband tussen toename in

**d**

We stellen hiervoor een lineair regressie model op. Uit de summary van het model zien we dat het niet bemesten een negatieve invloed (schatting van  $-14.772$ ) heeft op de groei en dat het afwezig zijn van competitie een positieve invloed (schatting van  $10.917$ ) heeft.



Figuur 5: Symmetrische verdeling van het model

We zien op de Q-Q plot (figuur 5) dat er bijna geen afwijking van normaliteit is. In combinatie met de determinatie coëfficiënt, die 0.66 is, kunnen we zeggen dat het model dus redelijk kwaliteitsvol is. De p-waarde is extreem laag, waardoor we kunnen beslissen dat de schattingen zeker goed zijn.

e

Hierbij trekken we gewoon **Height0** van **Height5** af. We bekommen terug dat niet bemesten negatief is en dat de afwezigheid van competitie positief is voor de groei van de bomen. De determinatie coëfficiënt is 0.67, net iets hoger dus het model is net iets beter.

f

We zien een schatting voor de gemiddelde hoogtes voor de vier variabelen in de tabel hieronder.

	Bemest	Onbemest
Competitie	33.06	18.35
Geen competitie	43.52	28.81

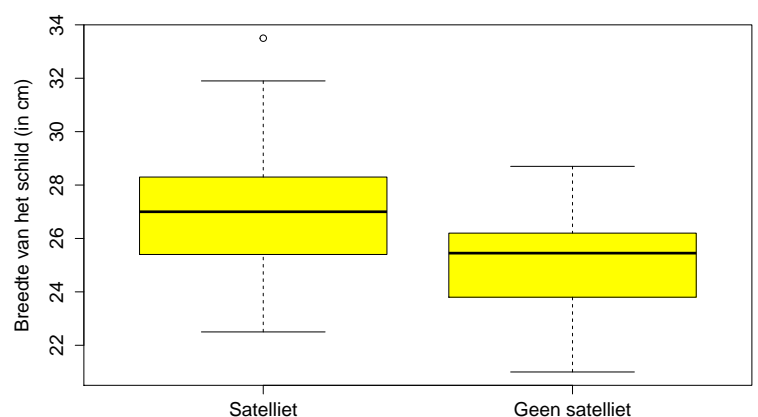
## Opgave 2

a

Bij deze opgave is het belangrijk op te merken dat er gevraagd wordt naar een percentage binnen de populatie, het is dus niet voldoende het percentage uit te rekenen voor de steekproef. Om dit percentage te berekenen maken we gebruik van een 95% betrouwbaarheidsinterval, we doen dit aan de hand van de formules in de cursus. We bekommen dat dit interval  $[0.5676, 0.7092]$  is. We kunnen dus met 95% zekerheid stellen dat 56.76% tot 70.92% van de vrouwelijke degenkrabben in de populatie over minstens één satelliet beschikt tijdens de paringsperiode. Aangezien het percentage aan vrouwelijke krabben met een satelliet in onze steekproef 64.16% is, is dit een representatieve steekproef.

b

De aanwezigheid van satellieten hangt inderdaad af van de grootte van het schild. Dit is duidelijk te zien in de boxplot op figuur 6. Om te weten in hoeverre deze verschilt doen we de t-test van Welch. Deze heeft als  $H_0$  dat beide gemiddeldes (de gemiddelde van de schildgrootte van vrouwtjes met en zonder satellieten) gelijk zijn.  $H_A$  stelt dan dat de gemiddeldes niet gelijk zijn. De p-waarde is  $9.495 \cdot 10^{-9}$ , dit levert een overtuigend bewijs tegen de nulhypothese.



Figuur 6: Verschil in grootte van het schild bij vrouwtjes met en zonder satelliet

Hierdoor kunnen we de nulhypothese verwerpen ten voordele van de alternatieve hypothese. Het betrouwbaarheidsinterval dat we bekomen met de t-test is  $[1.19, 2.33]$ . Vrouwtjes met satellieten hebben dus een schild dat 1.19 tot 2.33 cm groter is dan vrouwtjes zonder satellieten.

### c

Voor dit model op te stellen gebruiken we een lineaire en een kwadratische discriminant analyse. We vergelijken dan welke de beste is. We gebruiken ook leave-one-out cross-validation. Aan de hand van de tabellen kunnen we zien dat er bij de lineaire analyse met cross-validation een predictiefout is van 29.48% en bij de kwadratische met cross-validation een van 31.79%. Voor de lineaire zonder cross-validation vinden we een schijnbare predictiefout van 28.90% en voor de kwadratische zonder cross-validation 29.48%. De predictiefouten bij de schattingen met cross-validation liggen hoger dan die zonder, ze overfitten dus de training dataset wat leidt tot een instabieler classificatieregels die minder betrouwbare predicties oplevert. Op basis van de leave-one-out cross-validation schattingen van de predictiefout kunnen we besluiten dat de lineaire discriminant analyse stabielere resultaten oplevert dan de kwadratische.

	FALSE	TRUE
FALSE	27	15
TRUE	35	96

Tabel 1: Misclassificatietabel voor lineaire discriminant analyse

	FALSE	TRUE
FALSE	30	19
TRUE	32	92

Tabel 2: Misclassificatietabel voor lineaire discriminant analyse

	FALSE	TRUE
FALSE	26	15
TRUE	36	96

Tabel 3: Leave-one-out cross-validation misclassificatietabel voor lineaire discriminant analyse

	FALSE	TRUE
FALSE	28	21
TRUE	34	90

Tabel 4: Leave-one-out cross-validation misclassificatietabel voor kwadratische discriminant analyse

## Appendix

Hieronder vindt u de gebruikte R code met extra informatie in commentaar.

```

1  # 1)
2  #####
3  data <- read.csv("ZwarteSpar.csv", header = T)
4  #####
5
6
7  # a)
8  #####
9  # Bereken de groei over de laatste 5 jaar
10 dataF = subset(data, Fertilizer == "F")
11 FHeightDiff = dataF$Height5 - dataF$Height0
12 dataNF = subset(data, Fertilizer == "NF")
13 NFHeightDiff = dataNF$Height5 - dataNF$Height0
14
15 # Aangezien de groei in lengte een continue variabele is maken
16 # we hier een histogram van
17 hist(FHeightDiff, xlab = "Groei (in inch)", ylab = "Frequentie")
18 hist(NFHeightDiff, xlab = "Groei (in inch)", ylab = "Frequentie")
19 # Deze zijn niet symmetrisch en hebben pieken dus is het beter
20 # om hier de mediaan te gebruiken dan het gemiddelde
21
22 medF = median(FHeightDiff)
23 medNF = median(NFHeightDiff)
24
25 barplot(c(medF, medNF),
26         names.arg = c("Bemest", "Onbemest"),
27         ylab = "Hoogte groei over 5 jaar (in inch)", xpd = FALSE)
28 # We zien in de staafdiagram dat bemeste bomen een grotere groei
29 # vertonen dan onbemeste bomen. Het verschil is groot, meer dan
30 # 10 inch. Dit betekent dat de bomen sneller groeien als ze bemest
31 # worden.
32
33 # Om zeker te zijn dat dit niet beïnvloed wordt door een derde
34 # variabele doen we een permutatietoets
35 # Dit is een eenzijdige toets aangezien we enkel geïnteresseerd
36 # zijn in de toename
37
38 M <- 99999
39 meddiffs <- numeric()
40 v <- c(FHeightDiff, NFHeightDiff)
41 for(i in 1:M){
42   Fperm <- sample(1:(length(FHeightDiff) + length(NFHeightDiff)),
43                   length(FHeightDiff))
44   meddiffs[i] <- median(v[Fperm]) - median(v[-Fperm])
45 }
46 hist(meddiffs, prob = T, xlab = "Medianenverschil teststatistiek",
47      ylab = "Relatieve frequentie", main = "Eenzijdige permutatietoets",
48      xlim = c(-10, 15))
49
50 abline(v = medF-medNF, col="blue", lty = 2)
51 # Bereken de p waarde
52 p<-(sum(meddiffs >= (medF-medNF))+1)/(M+1)
53 p
54 # Deze is 10^-5, de steekproef levert dus een overtuigend bewijs tegen

```

```

55 # de nulhypothese van een gelijke hoogte groei
56 #####
57
58 # b)
59 #####
60 # Bereken de groei over de laatste 5 jaar
61 dataC = subset(data, Competition == "C")
62 CDiagDiff = dataC$Diameter5 - dataC$Diameter0
63 dataNC = subset(data, Competition == "NC")
64 NCDiagDiff = dataNC$Diameter5 - dataNC$Diameter0
65
66 # Aangezien de groei in lengte een continue variabele is maken
67 # we hier een histogram van
68 hist(CDiagDiff, xlab = "Diameter groei (in inch)", ylab = "Frequentie")
69 hist(NCDiagDiff, xlab = "Diameter groei (in inch)", ylab = "Frequentie")
70 # We zien dat de piek veel verkleint eens er geen competitie is wat de
71 # verdeling consistentere maakt.
72
73 medC = median(CDiagDiff)
74 medNC = median(NCDiagDiff)
75
76 barplot(c(medNC, medC),
77         names.arg = c("Geen competitie", "Competitie"),
78         ylab = "Diameter groei over 5 jaar (in inch)", xpd = FALSE)
79 # We zien dat er zonder competitie een veel groter groei in diameter
80 # is, meer dan 1 inch. Dit wil zeggen dat bomen zonder competitie
81 # sneller groeien in diameter.
82
83 # Om zeker te zijn dat dit niet beïnvloed wordt door een derde
84 # variabele doen we een permutatietoets
85 # Dit is een tweezijdige toets aangezien er niet gespecificeerd
86 # is wat het effect moet zijn
87 M <- 99999
88 meddiffs <- double()
89 v <- c(NCDiagDiff, CDiagDiff)
90 for(i in 1:M){
91   Cperm <- sample(1:(length(CDiagDiff) + length(NCDiagDiff)), length(CDiagDiff))
92   meddiffs[i] <- median(v[Cperm]) - median(v[-Cperm])
93 }
94 hist(meddiffs, prob = T, xlab = "Medianenverschil teststatistiek",
95      ylab = "Relatieve frequentie", main = "Tweezijdige permutatietoets")
96
97 abline(v = medNC-medC, col = 'blue', lty = 2)
98 # Bereken de p waarde, vermenigvuldigd met 2 aangezien het
99 # een tweezijdige permutatietoets is
100 p <- 2 * ((sum(meddiffs >= (medNC-medC))+1)/(M+1))
101 p
102
103 # Hoewel de p-waarde deze keer iets hoger ligt is dit nog altijd
104 # een lage waarde. Dit zorgt er dus voor dat de steekproef een
105 # overtuigend bewijs vormt tegen de nulhypothese van een gelijke
106 # toename in diameter
107 #####
108
109 # c)
110 #####
111 # Bereken de toenames
112 heightDiff = data$Height5 - data$Height0
113 diagDiff = data$Diameter5 - data$Diameter0

```

```

114 # Dit zijn continue variabelen dus gebruiken we covariantie om de
115 # associatie te modelleren
116 cov(heightDiff, diagDiff)
117 # Bereken toch de correlatie om makkelijker te vergelijken
118 cor(heightDiff, diagDiff)
119 # Plot dit in een scatterplot met een smoother
120 scatter.smooth(heightDiff, diagDiff, xlab = "Verschil in hoogte",
121               ylab = "Verschil in diameter")
122 # We zien dat de smoother lijkt op een rechte, er is dus een lineair
123 # verband tussen de toename van de lengte en deze van de diameter.
124 # Aangezien de correlatie 0.90 is, is dit ook een sterk verband
125 #####
126
127 # d)
128 #####
129 # Stel het model op, Height5 hangt af van Fertilizer en Competition
130 model <- lm(data$Height5 ~ data$Fertilizer + data$Competition)
131 # QQplot p179
132 plot(model)
133 summary(model)
134 # We zien dat niet bemesten een negatief effect heeft en dat de
135 # afwezigheid van competitie een positief effect heeft op de groei.
136 # We zien ook dat de p-waarde extreem laag is (< 2.2e-16),
137 # de schattingen zijn dus zeker goed
138 # De determinatie coefficient is 0.66 dus ook het model is redelijk
139 # goed.
140 #####
141
142 # e)
143 #####
144 # Breng de hoogte bij start in rekening
145 model <- lm(data$Height5 - data$Height0 ~ data$Fertilizer + data$Competition)
146 plot(model)
147 summary(model)
148 # De determinatie coefficient is hier net iets hoger wat dit model
149 # beter maakt.
150 #####
151
152 # f)
153 #####
154 # Print de inschatting voor de gemiddeldes voor de vier combinaties
155 print(paste("Fertilized + Competition: ", model$coeff %*% c(1,F,F)))
156 print(paste("Fertilized + No competition: ", model$coeff %*% c(1,F,T)))
157 print(paste("Not fertilized + Competition: ", model$coef %*% c(1,T,F)))
158 print(paste("Not fertilized + No competition: ", model$coeff %*% c(1,T,T)))
159 #####
160
161
162
163
164 # 2)
165 #####
166 # Clear the environment
167 closeAllConnections()
168 rm(list=ls())
169 # Read the new data
170 data <- read.csv("Krabben.csv", header = T)
171 #####
172

```



```

173 # a)
174 #####
175 # Splits de data
176 dataS = subset(data, satell == "TRUE")
177 dataNS = subset(data, satell == "FALSE")
178 # Bepaal het percent vrouwelijke krabben dat een satelliet heeft
179 # binnen onze steekproef
180 dataSPercent = length(dataS$satell)/length(data$satell)
181 # Bepaal de grenzen van het 95% betrouwbaarheidsinterval
182 noemer = length(data$satell) + 4
183 tellerp = (length(dataS$satell) + 2)/noemer
184 bovengrens = tellerp + 1.96 * sqrt((tellerp*(1-tellerp))/noemer)
185 ondergrens = tellerp - 1.96 * sqrt((tellerp*(1-tellerp))/noemer)
186 ondergrens; bovengrens
187 #####
188
189 # b)
190 #####
191 # Visualiseer de data in histogrammen
192 hist(dataS$width)
193 hist(dataNS$width)
194 # De data is niet symmetrisch dus we gebruiken de mediaan
195
196 barplot(c(median(dataS$width), median(dataNS$width)),
197         names.arg = c("Satelliet", "Geen satelliet"),
198         ylab = "Breedte van het schild (in cm)", xpd = FALSE)
199 # De schildgrootte ligt iets hoger bij de vrouwtjes die over een
200 # satelliet beschikken. Het verschil is echter heel klein, we wensen
201 # meer informatie. Deze kan geleverd worden door een boxplot
202
203 boxplot(dataS$width, dataNS$width, col="yellow",
204         ylab = "Breedte van het schild (in cm)",
205         names = c("Satelliet", "Geen satelliet"))
206 # De schildgrootte heeft dus duidelijk een effect op het al dan niet
207 # beschikken over een satelliet of niet.
208
209 # Nu willen we weten hoezeer de breedte van het schild verschilt
210 # We voeren de t-test van Welch uit
211 t.test(dataS$width, dataNS$width)
212 #####
213
214 # c)
215 #####
216 # Importeer de bibliotheek die nodig is voor discriminant analyses
217 library(MASS)
218 # Voer lineaire en kwadratische discriminant analyses uit
219 lda = lda(satell ~ ., data)
220 qda = qda(satell ~ ., data)
221 # Schatten
222 plda = predict(lda)
223 pqda = predict(qda)
224
225 # Toon de misclassificatietabellen
226 lt = table(plda$class, data$satell)
227 qt = table(pqda$class, data$satell)
228
229 (lt[2] + lt[3]) / sum(lt)
230 # Bereken de kwadratische predictiefout
231 (qt[2] + qt[3]) / sum(qt)

```

```
232
233 # Doe een lineaire discriminant analyse op de satelliet data, met
234 # leave-one-out cross-validation
235 lda = lda(satell ~ ., data, CV=TRUE)
236 # Doe een kwadratische discriminant analyse op de satelliet data, met
237 # leave-one-out cross-validation
238 qda = qda(satell ~ ., data, CV=TRUE)
239
240 # Bereken de leave-one-out cross-validation misclassificatietabel voor
241 # lineaire discriminant analyse
242 lt = table(lda$class, data$satell)
243 # Bereken de leave-one-out cross-validation misclassificatietabel voor
244 # kwadratische discriminant analyse
245 qt = table(qda$class, data$satell)
246 lt
247 qt
248
249 # Bereken de lineaire predictiefout
250 (lt[2] + lt[3]) / sum(lt)
251 # Bereken de kwadratische predictiefout
252 (qt[2] + qt[3]) / sum(qt)
253 #####
```

Oplossingen.R