



Statistische gegevensanalyse

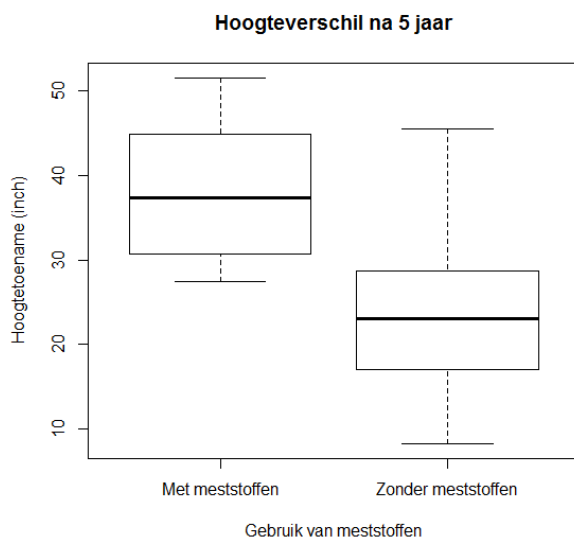
Project

Bart Middag
3^{de} bachelor informatica
Academiejaar 2013-2014

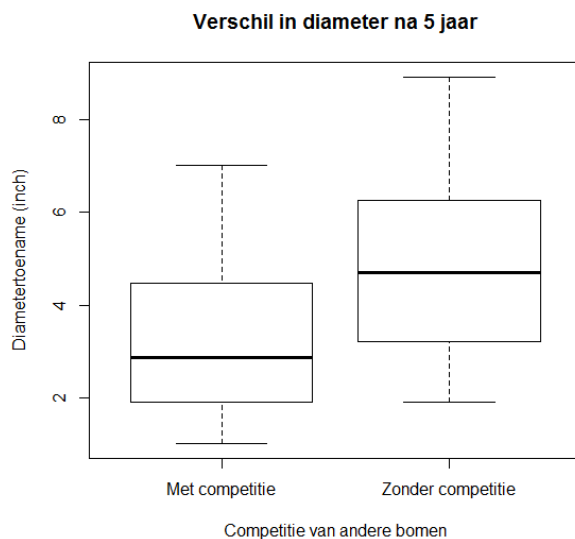
1. DE ZWARTE SPAR

a) Het effect van gebruik van meststoffen op de hoogte:

We bepalen de groei van de bomen met en zonder gebruik van meststoffen. Het resultaat wordt getoond in Figuur 1. In deze steekproef is er dus duidelijk een verschil tussen de groei met gebruik van meststoffen en zonder. Om met zekerheid te kunnen zeggen dat dit geen toeval is, bepalen we aan de hand van een permutatietest dus de kans dat het verschil even groot of groter is. Voor geen enkele van de 10000 permutaties was het verschil zo groot als in de originele steekproef. De kans om een dergelijk verschil te bekomen is dan ook $\frac{1}{10001}$. We kunnen met 95% zekerheid stellen dat bomen bij gebruik van meststoffen gemiddeld tussen 11.60278 en 17.88889 inch hoger zullen worden dan zonder meststoffen.



Figuur 1: Boxplot van de groei van de zwarte spar in verband met het gebruik van meststoffen.



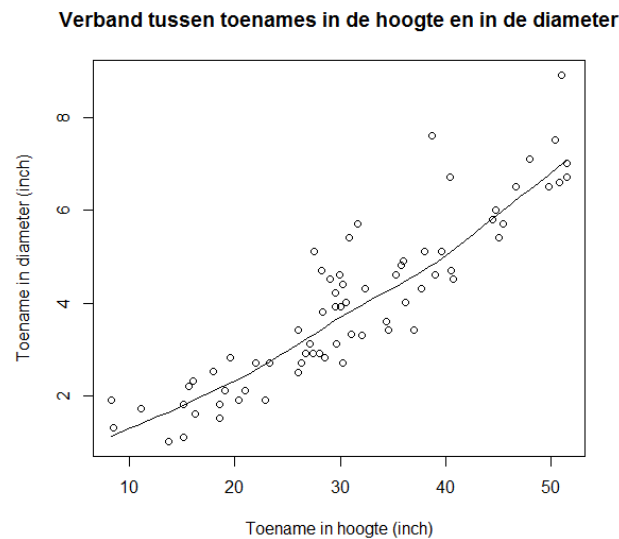
Figuur 2: Boxplot van de groei van de diameter van de zwarte spar in verband met competitie van andere bomen in de omgeving.

b) Het effect van de omgeving op de diameter:

Analoog met de vorige opgave, bepalen we de diameter van de bomen met en zonder competitie van andere bomen in de omgeving. Het resultaat in Figuur 2 vertoont een duidelijk negatief effect van bomen in de omgeving. Aan de hand van een t-test bekomen we dat we met 95% zekerheid kunnen stellen dat bomen met competitie van andere bomen in de omgeving tussen 0.81824 en 2.33575 inch minder dik zullen worden dan bomen zonder competitie.

c) Het verband tussen de toename van de hoogte en van de diameter:

We bepalen het verband tussen de toename van de hoogte en van de diameter aan de hand van Pearsons product-moment correlatiecoëfficiënt. Deze wordt geschat op 0.90208. Er is dus bijna een perfect lineair verband tussen de toename van de hoogte en van de diameter. We kunnen met 95% zekerheid stellen dat dit ligt tussen 0.84753 en 0.93777. Dit verband wordt duidelijk getoond op Figuur 3.



Figuur 3: Scatterplot van het verband tussen toename in de hoogte en toename in de diameter van de zwarte spar.

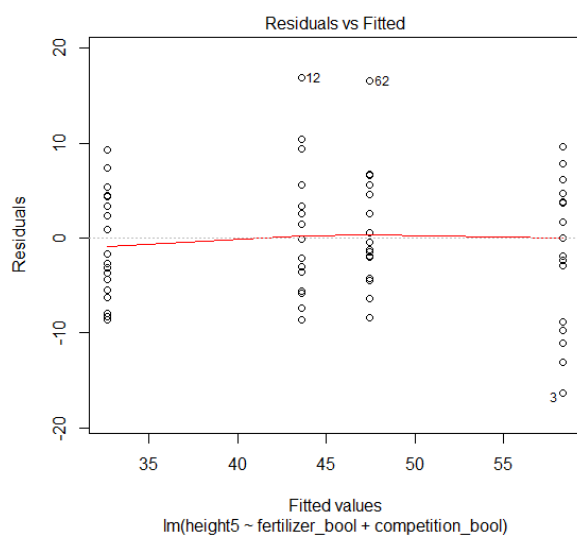
d) Model van de hoogte na 5 jaar a.d.h.v. meststofgebruik en competitie:

Na toepassing van lineaire regressie, bekomen we het volgende model:

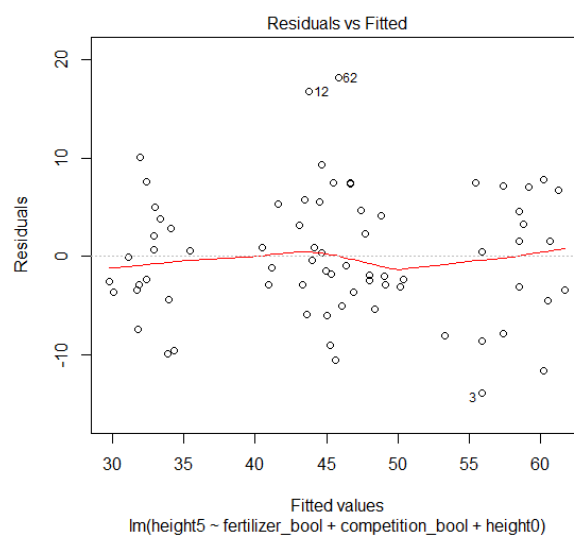
$$w = 43.578 + 14.772 * x - 10.917 * y$$

Hierbij staat x voor het gebruik van meststoffen en y voor de aanwezigheid van andere bomen in de omgeving. De p -waarde die bij de nulhypothese hoort die stelt dat er geen invloed is van meststoffen en/of competitie, is kleiner dan $2.2 * 10^{-16}$.

Zoals we onderzocht hebben in opgave a en b, hebben beide variabelen dus invloed op de hoogte. De determinatiecoëfficiënt van dit model is 0.6586. De standaardfout van het residu is 6.611. Op Figuur 4 zien we de residuen tegenover de gefitte waarden.



Figuur 4: Residuen tegenover de gefitte waarden van een model dat de beginhoogte niet in rekening brengt.



Figuur 5: Residuen tegenover de gefitte waarden van een model dat de beginhoogte in rekening brengt.

e) Model van de hoogte na 5 jaar a.d.h.v. meststofgebruik, competitie en beginhoogte:

Na toepassing van lineaire regressie, bekomen we het volgende model:

$$w = 30.8296 + 14.7195 * x - 10.5235 * y + 0.8631 * z$$

Hierbij staat x voor het gebruik van meststoffen, y voor de aanwezigheid van andere bomen in de omgeving en z voor de beginhoogte. De p -waarde die bij de nulhypothese hoort die stelt dat er geen invloed is van meststoffen en/of competitie en/of de beginhoogte, is kleiner dan $2.2 * 10^{-16}$. Alle variabelen hebben invloed op de hoogte na 5 jaar. De determinatiecoëfficiënt van dit model is 0.6786. De standaardfout van het residu is 6.414. In vergelijking met het vorige model ligt de fout lager en de determinatiecoëfficiënt hoger: dit model is dus beter. Op Figuur 5 zien we de residuen tegenover de gefitte waarden.

f) Schattingen voor de gemiddelde hoogte van de spar na 5 jaar:

Aan de hand van het laatste model voeren we de schattingen uit. De resultaten zijn weergegeven in Tabel 1.

Geen meststoffen	Geen competitie	31.69274 inch
Meststoffen	Geen competitie	46.41222 inch
Geen meststoffen	Competitie	21.16929 inch
Meststoffen	Competitie	35.88876 inch

Tabel 1: Schattingen voor de gemiddelde hoogte van de zwarte spar na 5 jaar

2. DEGENKRABBen

a) Het percentage vrouwelijke degenkrabben met minstens één satelliet:

Op basis van de steekproef schatten we dat het percentage vrouwelijke degenkrabben met minstens één satelliet gelijk is aan 64.16185%.

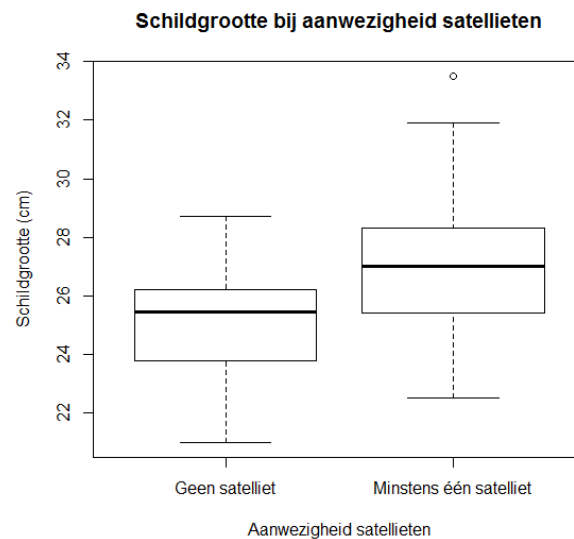
b) Het effect van de schildgrootte op de aanwezigheid van satellieten:

Analoog met de opgaves 1a en 1b, bepalen we de schildgrootte van de krabben met en zonder de aanwezigheid van satellieten. Aan de hand van Figuur 6 zien we dat er meer satellieten aanwezig zullen zijn bij een krab met een groter schild. Met een t-test bekomen we dat we met 95% zekerheid kunnen stellen dat krabben met satellieten een schild zullen hebben dat gemiddeld tussen 1.18848 cm en 2.33227 cm groter is dan het schild van krabben zonder satellieten.

c) Model voor de aanwezigheid van satellieten:

Om dit model op te stellen, gebruiken we lineaire en kwadratische discriminantanalyse. Het schijnbaar foutenpercentage bij lineaire discriminantanalyse is 33.07692% en dit bij kwadratische discriminantanalyse is 29.23077%. We beschikken over relatief weinig data, dus gebruiken we kruisvalidatie gebruikt om deze percentages te berekenen. De bijhorende misclassificatietabellen zijn weergegeven in Tabel 2 en Tabel 3. We zien dat

kwadratische discriminantanalyse een kleiner foutenpercentage oplevert, dus we gebruiken dit model. We zien dat dit model een predictiefout van 37.2093% heeft.



Figuur 6: Boxplot van de schildgrootte in verband met de aanwezigheid van satellieten.

LDA	Geen satelliet	≥ 1 satelliet
Geen satelliet	16.15385%	13.84615%
≥ 1 satelliet	19.23077%	50.76923%

Tabel 2: Misclassificatietabel bij lineaire discriminantanalyse. De echte klassen worden weergegeven in de kolommen, de voorspelde klassen in de rijen.

QDA	Geen satelliet	≥ 1 satelliet
Geen satelliet	19.23077%	13.07692%
≥ 1 satelliet	16.15385%	51.53846%

Tabel 3: Misclassificatietabel bij kwadratische discriminantanalyse. De echte klassen worden weergegeven in de kolommen, de voorspelde klassen in de rijen.

3. APPENDIX: R-CODE

In het verslag werden enkel de resultaten gerapporteerd: de code zelf is iets ingewikkelder en er werden bij enkele opgaves verschillende controles uitgevoerd om te weten welke oplossingsmethode gebruikt moest worden. De werkwijze staat in onderstaande code beschreven. U vindt de code ook terug in het bestand `Project.R`.

```

1 #####
2 # Project: Statistische Gegevensanalyse #
3 # Naam: Bart Middag #
4 # Richting: 3de bachelor informatica #
5 # Academiejaar: 2013-2014 #
6 #####
7
8 # Voorbereiding project
9 nsim <- 10000 # aantal simulaties
10 set.seed(98765) # seed voor reproduceerbare resultaten
11 setwd("C:/BART/UNIF/Statistische Gegevensanalyse/Project")
12
13 #####
14 # OPGAVE 1 #
15 #####
16
17 # Voorbereiding opgave 1
18
19 zwarte_spar <- read.csv("ZwarteSpar.csv", header = T)
20 str(zwarte_spar)
21 summary(zwarte_spar)
22 colnames(zwarte_spar) <- tolower(colnames(zwarte_spar))
23 attach(zwarte_spar)
24
25 # OPGAVE 1A
26
27 # Hoogte met/zonder meststoffen en groei berekenen
28 f_height0 <- height0[fertilizer == "F"]
29 f_height5 <- height5[fertilizer == "F"]
30 nf_height0 <- height0[fertilizer == "NF"]
31 nf_height5 <- height5[fertilizer == "NF"]
32 f_growth <- f_height5 - f_height0
33 nf_growth <- nf_height5 - nf_height0
34 height_growth <- height5 - height0
35
36 # Boxplots
37 boxplot(height_growth~fertilizer, main="Hoogteverschil na 5
jaar", xlab="Gebruik van meststoffen", ylab="Hoogtetoeename
(inch)", xaxt="n")
38 axis(side=1, at=1:2, labels=c("Met meststoffen", "Zonder
meststoffen"))
39
40 # We kijken naar de verdeling van de groei en zien of ze normaal
verdeeld is
41 plot(density(f_growth))
42 plot(density(nf_growth))
43 plot(density(height_growth))
44 qqnorm(f_growth); qqline(f_growth, col = 2)
45 qqnorm(nf_growth); qqline(nf_growth, col = 2)
46 qqnorm(height_growth); qqline(height_growth, col = 2)
47

```

```

48 # We kijken of de groei normaal verdeeld is - de Shapiro-Wilk
    test geeft ons duidelijke resultaten.
49 shapiro.test(f_growth) # Niet normaal verdeeld, de t-test mogen
    we dus niet gebruiken!
50 shapiro.test(nf_growth) # Wel normaal verdeeld
51 shapiro.test(height_growth) # Algemeen normaal verdeeld
52
53 # We bekijken het verschil in gemiddeldes
54 mean_diff <- mean(f_growth) - mean(nf_growth)
55 # We bepalen de kans dat een willekeurige permutatie een
    resultaat geeft >= het gemiddelde in deze situatie.
56 means <- numeric()
57 for(i in 1:nsim) {
58   permutation <- sample(height_growth)
59   means[i] <- mean(permutation[fertilizer == "F"]) -
    mean(permutation[fertilizer == "NF"])
60 }
61 hist(means)
62 means[nsim+1] <- mean_diff
63 means_p <- sum(means >= mean_diff)/(nsim+1)
64 # Op basis van 10000 permutaties is de kans ongeveer 1/10001.
65
66 # We bepalen het betrouwbaarheidsinterval.
67 differences <- numeric()
68 for(i in 1:nsim) {
69   differences[i] <- mean(sample(f_growth, replace = T)) -
    mean(sample(nf_growth, replace = T))
70 }
71 differences[nsim+1] <- mean_diff
72 differences <- sort(differences)
73 interval <- c(differences[0.05*(nsim+1)], differences[(1-
    0.05)*(nsim+1)])
74 cat(paste0("We kunnen met 95% zekerheid stellen dat de extra
    groei zal liggen tussen ", interval[1], " en ", interval[2], "
    inch. "))
75
76
77 # OPGAVE 1B
78
79 # Diameter met/zonder competitie en groei berekenen
80 c_diameter0 <- diameter0[competition == "C"]
81 c_diameter5 <- diameter5[competition == "C"]
82 nc_diameter0 <- diameter0[competition == "NC"]
83 nc_diameter5 <- diameter5[competition == "NC"]
84 c_growth <- c_diameter5 - c_diameter0
85 nc_growth <- nc_diameter5 - nc_diameter0
86 diameter_growth <- diameter5 - diameter0
87
88 # Boxplots
89 boxplot(diameter_growth~competition, main="Verschil in diameter
    na 5 jaar", xlab="Competitie van andere bomen",
    ylab="Diameter toename (inch)", xaxt="n")
90 axis(side=1, at=1:2, labels=c("Met competitie", "Zonder
    competitie"))
91
92 # We kijken naar de verdeling van de groei en zien of ze normaal
    verdeeld is
93 plot(density(c_growth))
94 plot(density(nc_growth))
95 qqnorm(c_growth); qqline(c_growth, col = 2)
96 qqnorm(nc_growth); qqline(nc_growth, col = 2)

```

```
97 # We kijken of de groei normaal verdeeld is - de Shapiro-Wilk
    test geeft ons duidelijke resultaten.
98 shapiro.test(c_growth) # Wel normaal verdeeld
99 shapiro.test(nc_growth) # Wel normaal verdeeld
100
101 # We mogen de t-test dus gebruiken.
102 t.test(c_growth,nc_growth)
103 # Er is dus een negatief effect op de toename in diameter als er
    competitie is.
104
105 # OPGAVE 1C
106
107 # We bepalen de associatie tussen de toenames van de hoogte en
    van de diameter
108 height_diameter <- cor(height_growth, diameter_growth)
109 cor.test(height_growth, diameter_growth)
110
111 # We plotten het verband
112 scatter.smooth(height_growth, diameter_growth, main="Verband
    tussen toenames in de hoogte en in de diameter", xlab="Toename in
    hoogte (inch)", ylab="Toename in diameter (inch)")
113
114 # OPGAVE 1D
115
116 # We zetten dit om naar een logische eenheid voor R, zodat R niet
    NC en NF beschouwt maar F en C.
117 # Als we zouden werken met een model dat NC en NF beschouwt,
    moeten we het effect van F en C inverteren en dat is verwarrend.
118 fertilizer_bool <- as.logical(fertilizer == "F")
119 competition_bool <- as.logical(competition == "C")
120
121 lmfit <- lm(height5~fertilizer_bool+competition_bool)
122 summary(lmfit)
123 plot(lmfit)
124 coef(lmfit)
125
126 # OPGAVE 1E
127
128 lmfit_height0 <-
    lm(height5~fertilizer_bool+competition_bool+height0)
129 summary(lmfit_height0)
130 plot(lmfit_height0)
131 coef(lmfit_height0)
132
133 # OPGAVE 1F
134
135 lmfit_height0$coeff %*% c(1,F,F,T)
136 lmfit_height0$coeff %*% c(1,T,F,T)
137 lmfit_height0$coeff %*% c(1,F,T,T)
138 lmfit_height0$coeff %*% c(1,T,T,T)
139
140 #####
141 # OPGAVE 2 #
142 #####
143
144 # Voorbereiding opgave 2
145
146 krabben <- read.csv("Krabben.csv",header = T)
147 str(krabben)
148 summary(krabben)
```



```
149 # colnames(krabben) <- tolower(colnames(krabben)) # Ze zijn al
    lowercase.
150 attach(krabben)
151 library(MASS)
152
153 # OPGAVE 2A
154
155 satellites_percent <- (length(satell[satell=='TRUE']) /
    length(satell))*100
156 satellites_percent
157
158 # OPGAVE 2B
159
160 # Grootte van het schild in functie van aanwezigheid van
    satellieten berekenen
161 t_width <- width[satell == T]
162 f_width <- width[satell == F]
163
164 # Boxplots
165 boxplot(width~satell, main="Schildgrootte bij aanwezigheid
    satellieten", xlab="Aanwezigheid satellieten",
    ylab="Schildgrootte (cm)", xaxt="n")
166 axis(side=1, at=1:2, labels=c("Geen satelliet", "Minstens één
    satelliet"))
167
168 # We kijken naar de verdeling van beide groepen en zien of ze
    normaal verdeeld zijn
169 plot(density(t_width))
170 plot(density(f_width))
171 qqnorm(t_width); qqline(t_width, col = 2)
172 qqnorm(f_width); qqline(f_width, col = 2)
173
174 # We kijken of de groepen normaal verdeeld zijn - de Shapiro-Wilk
    test geeft ons duidelijke resultaten.
175 shapiro.test(t_width) # Wel normaal verdeeld
176 shapiro.test(f_width) # Wel normaal verdeeld
177
178 # We mogen de t-test dus gebruiken.
179 t.test(t_width,f_width)
180 # Er is dus een duidelijk verband tussen de schildgrootte en de
    aanwezigheid van satellieten.
181
182 # OPGAVE 2C
183
184 krabben <- krabben[,c(1,3,2)] # kolom 2 en 3 switchen
185 index <- sample(c(rep("training", 130), rep("test", 43))) # test:
    training/3
186 krabben_train <- krabben[index == "training",]
187 krabben_test <- krabben[index == "test",]
188
189 # Trainen van modellen
190 lda_train <- lda(satell ~ ., data = krabben_train)
191 qda_train <- qda(satell ~ ., data = krabben_train)
192
193 # Validatie (functie uit practicum 10)
194 K <- 5 #aantal folds
195
196 own.cv <- function(x,y, K = 5,method = lda){
197   f <- method
198   n <- nrow(x)
199
```

```
200 grid <- rep(1:K, n%/%K+1 ) [1:n]
201 id <- sample(grid)
202 preds <- rep(NA,n)
203 for(i in 1:K){
204   f.model <- f( x[id != i,],y[id != i])
205   preds[id == i]<- predict(f.model,newdata = x[id ==
i,])$class
206 }
207 preds-1
208 }
209
210 # Schijnbaar foutenpercentage lineaire discriminantanalyse
211 preds.cvlda<-own.cv(krabben_train[,1:2], krabben_train[,3], K =
K, method = lda)
212 sum(krabben_train$satell != preds.cvlda)/nrow(krabben_train)
#cross-validation error
213 table(preds.cvlda,krabben_train$satell)/130
214
215 # Schijnbaar foutenpercentage kwadratische discriminantanalyse
216 preds.cvqda<-own.cv(krabben_train[,1:2],krabben_train[,3],K =
K,method = qda)
217 sum(krabben_train$satell != preds.cvqda)/nrow(krabben_train)
#cross-validation error
218 table(preds.cvqda,krabben_train$satell)/130
219
220 # QDA > LDA, maar dit hangt af van de seed.
221 # Predictiefout
222 sum(predict(qda_train,newdata = krabben_test)$class !=
krabben_test$satell)/43
```