

```

library(readxl)
data <- read_excel("2020 - assignment 3.xlsx", sheet = "habspecies", na =
"NA")
head(data)
str(data)
colnames(data)

par(mfrow=c(1,3))
plot(data$E, data$N,col=factor(data$minkp),pch=16, xlab= "East", ylab =
"North",main="mink")
plot(data$E, data$N,col=factor(data$martenp),pch=16, xlab= "East", ylab =
"North",main="marten")
plot(data$E, data$N,col=factor(data$otterp),pch=16, xlab= "East", ylab =
"North",main="otterp")

#install packages

install.packages("car")
install.packages("psych")
install.packages("carData")
install.packages("mvtnorm")

library(car)
library(carData)
library(psych)
library(lattice)
library(data.table)
library(plyr)
library(doBy)
library(afex)
library(multcomp)
library(lsmeans)
library(effects)
library(MASS)
library(plot3D)
library(rgl)
library(gmodels)

#Question 1

mytable1=xtabs(~minkp+martenp, data=data)
CrossTable(mytable1, expected=T, chisq=T, fisher = T, format = "SAS")
CrossTable(mytable1, expected=T, chisq=T, format = "SAS")
chisq.test(mytable1, p=c(84.778,27.222,24.222,7.778), rescale.p=T, correct=F)

#Conclusion: In this case, the  $p < 0.05$  would imply that we would reject the
#null hypothesis that both factors are independent, implying a significant
#association between Mink and Marten. Pearson's chi-square value
#( $\chi^2 = 5.955739$ ), which confirms that association between Mink and Marten
#is statistically significant.
#Biological Conclusion: Mink and marten co-occur more than expected by chance.

mytable3=xtabs(~minkp+otterp, data=data)

```

```

CrossTable(mytable3, expected=T, chisq=T, fisher = T, format = "SAS")
CrossTable(mytable3, expected=T, chisq=T, format = "SAS")

#Conclusion: In this case, the  $p > 0.05$  would imply that we would not reject
#the null hypothesis that both factors are independent, implying there is no
#significant association between Mink and Otter. Pearson's chi-square value
#( $\chi^2 = 2.992208$ ), which confirms that association between Mink and Otter
#is not statistically significant.
#Biological Conclusion: Mink and Otter co-occur less than expected by chance.

```

```

mytable5=xtabs(~martenp+otterp, data=data)
CrossTable(mytable5, expected=T, chisq=T, fisher = T, format = "SAS")
CrossTable(mytable5, expected=T, chisq=T, format = "SAS")

```

```

#Conclusion: In this case, the  $p > 0.05$  would imply that we would not able to
#reject the null hypothesis that both factors are independent, implying no
#significant association between Marten and Otter. Pearson's chi-square value
#( $\chi^2 = 2.899$ ), which confirms that association between Marten and Otter
#is not statistically significant.
#Biological Conclusion: Marten and Otter co-occur less than expected by
chance.

```

```

#Question 2

```

```

#created a data set for Figga river, to predict the current of Figga river

```

```

data.Figga <-subset(data, rivertype== "main")
data.FiggaRiver <-subset(data, river== "Figga")

```

```

# this data set is based on Current of Figga river, Figga river is the main
#type river and rest of the rivers are side rivers

```

```

# Model selection and statistical inference without log-transforming BD:
# plot before regression

```

```

#####
#Calculations in R
#####

```

```

data1 = lm(current ~ riverwidth * oceandist *riverdepth*bankheight,
data=data.Figga)
summary (data1)

```

```

# Answer: all variable individually and with interaction are not statistically
significant
# overall model p vaule (0.01725) and F-value (2.242) indicates model is
statistically significant, this model
#has low R squre and adujusted R square is low [Multiple R-squared:0.407,
Adjusted R-squared:0.2254]

```

```

#As all the variabels are statistically insignificant so,we have to choose a
model where some of the variables
#are significant, so we have run some other techniques to determine a better
fit model to predict the current
#of Figga river based on the explanatory variable inluded in the model.

Anova(data1, type = 3)

# Answer:All variable individually and with interation are not statistically
significant

summary(data1)$adj.r.squared
# Answer: [1] 0.225418 Adj R square is bit low that indicated model is not fit
well

#####
#Stepwise Model reduction
#####

#plot before regression

attach(data.Figga)
par(mfrow =c(2,2))
plot(current~riverwidth, main = "current vs riverwidth")
plot(current~oceandist, main = "current vs oceandist")
plot(current~riverdepth, main = "current vs riverdepth")
plot(current~bankheight, main = "current vs bankheight")
boxplot(current~riverdepth:bankheight, main = "current vs joint")

full.model=lm(current ~ riverwidth * oceandist *riverdepth*bankheight,
data=data.Figga)
stepAIC(full.model)

#Start: AIC=-185.76 #Step: AIC=-189.66 #Step: AIC=-191.44
#Step: AIC=-192.54 #Step: AIC=-192.67 #Step: AIC=-194.56
#Step: AIC=-195.3 #Step: AIC=-195.47 #Step: AIC=-195.78

#I generated 9 separaet model (reported above), based on AIC value, AIC=
-195.78 is the
#lowest AIC among all the model I generated. here is my final model reported
below

#linear preduction model

modell=lm(current ~ riverwidth + oceandist + riverdepth +
bankheight + riverwidth:bankheight + riverdepth:bankheight,
data = data.Figga)
summary(modell)
Anova(modell, type = 3)
summary(modell)$adj.r.squared

```

```

# My final model is statistically significant based on P value (0.0005351) and
# F statistic (4.749),
# both the value indicates model is overall statistically significant. In my
# final model, I observed all the
# variable statistically insignificant at 5% level of significance except
# "riverdepth:bankheight"
# (interaction variable) which is significant at 5% level of significance. It
# means river depth and bank height
# significantly affect the current of Figga river. According to data Figga
# river is the main river and rest of the
# rivers are side rivers. in the final model river weight, ocean distance and
# bank height got the negative
# coefficient which means Current of Figga River negatively associated with
# river weight, ocean distance
# and bank height. only riverwidth*bankheight interaction has positive
# coefficient which means these two
# jointly positively related to current of Figga river though they are
# statistically insignificant at
# 5% level of significance.

```

```

vif(modell1)

```

```

# variance inflation factors here river width, ocean distance and river depth
# are <5, which seems good
# but bankheight, riverwidth:bankheight riverdepth:bankheight are >5 that's
# doesn't seem good. SO I would run
# again the regression based on the variable where VIF <5, I took out the
# variables from the regression whose VIF >5.
# here is the model using the variables river width, ocean distance and river
# depth

```

```

riazmodell1=lm(current ~ riverwidth + oceandist +riverdepth,
               data = data.Figga)
summary(riazmodell1)
Anova(riazmodell1, type = 3)
summary(riazmodell1)$adj.r.squared

```

```

vif (riazmodell1)

```

```

#riverwidth  oceandist  riverdepth
# 1.037952    1.060488    1.027099
# based on VIF test, value of the all these variables (riverwidth, oceandist
# and riverdepth) are less than 5,
# meaning all three seems good to this model, i will use later

```

```

#regression results of the model
summary(riazmodell1)
summary(riazmodell1)$adj.r.squared

```

```

# regression result of this model showed that riverdepth is the only
significant variable in the model, and rest
#two varibales (riverwidth and oceandist) are statistically insignificant.
river depth got the negative coefficent
#meaning current and river depth are negatively corellated,

#Plot after regression
plot(allEffects(riazmodell))

# Model diagnosis, outliers and influential observations for "final"

residualPlots(riazmodell)
#linearity cannot be assumed

spreadLevelPlot(riazmodell)
# the plot suggests that variance increases

ncvTest(riazmodell)
# Based on ncvtest (Chisquare = 1.934064, Df = 1, p = 0.16431).
#Formal test rejects constant variance,based on p value model is statistically
insignificant

studres.riazmodell=rstudent(riazmodell) # studentized residuals
hist(studres.riazmodell,
     probability=T,
     col="lightgrey",
     xlim=c(-6,6),
     breaks=12,
     main="Distribution of Studentized Residuals",
     xlab="Studentized residuals")
xfit=seq(-6,6,length=100)
yfit=dnorm(xfit) # normal fit
lines(xfit, yfit, col="red",lwd=2)

# The distribution looks very good

shapiro.test(residuals(riazmodell))
# Formal test confirming normality since W (0.90294) > 0.9

outlierTest(riazmodell) # observation 60 is an outlier
#No Studentized residuals with Bonferroni p < 0.05
#Largest |rstudent|:
#rstudent unadjusted p-value Bonferroni p
#60 3.18266 0.0023127 0.15033

influenceIndexPlot(riazmodell,vars=c("Studentized","Bonf"))

```

```
influenceIndexPlot(riazmodell,vars="Cook") # but there are no influential
observations, so no problems there
```

```
# Conclusion
```

```
# The model "final" meets all assumptions, so we keep this model.
# Figga rivers current statistically depends on the riverdepth, current and
riverdepth negatively related
```

```
Anova(riazmodell, type = 3)
summary(riazmodell)
```

```
# The significant river depth effect on current of figga river indicates that
the Figga River current
```

```
# significantly effect the river depth of figga river.
```

```
#The rest of the variables individually for example riverwidth and distance
of ocean statistically insignificant.
```

```
#Meaning, these variables do not effect the current of Figga river.
```

```
#No, there is no significant interaction effect in the model
```