

Computer Exercise 6 (R)

1. White-throated sparrows occur in two distinct color morphs, referred to as brown and white. It was suspected that females select males of the opposite morph (i.e., white females select brown males and vice versa). This phenomenon is known as negative assortative mating. In 14 mated pairs, the color combinations were as follows:

	White male	Brown male
White female	1	6
Brown female	5	2

Test for assortative mating (negative or otherwise), i.e. test the null hypothesis of no dependence of the male morph on the female morph in a mated pair. Note that there is a problem with small expected numbers, so the ordinary chi-squared test may be unreliable. A good course of action is then to use the so-called Fisher's exact test. You first need to put the data into a matrix (see pp. 54-59 in the R-Book).

R console:

```
tbl <- matrix(c(1, 5, 6, 2), nrow=2)
```

If you like, you may add column and row names by the following commands

```
colnames(tbl) <- c("White male", "Brown male")
```

```
rownames(tbl) <- c("White female", "Brown female")
```

Write `tbl` and press enter to have a look on your data.

Now, perform the Fisher's exact test

```
fisher.test(tbl)
```

You might compare with what you get from

```
chisq.test(tbl)
```

```
chisq.test(tbl, correct=F)
```

Which test do you decide to use and what do you conclude about assortative mating?

2. Read the data in **Foodplnt.txt** (these data are the same as in **Foodplnt.txt** from Exercise 3). The variable `AduDryw` is the dry weight (in mg) of adult green-veined whites (*Pieris napi*), grown on different kinds of host plants. Four types of plants were used, coded as follows: 1 = high quality *Aliaria petiolata*, 2 = low quality *A. petiolata*, 3 = high quality *Berteroa incana*, 4 = low quality *B. incana*. The sexes were coded as: 1 = male, 2 = female. Make sure that the variables `Sex` and `Foodpl` are made into factors (perhaps with suitable labels for the levels).

If you did last week's exercises you know that if one considers just one of the sexes (either males or females), there was an effect of food plant on adult weight. Your job now is to see if male and female adult weight responded differently to different food plants, e.g. if one sex was more sensitive to food plant variation (this was the actual purpose of the experiment from which these data come). Analyze the data with sex and food plant as (fixed effect) factors, for instance

```
R console
fml <- lm(AduDryW ~ Sex + Foodpl + Sex:Foodpl, data=dat)
```

and test for an interaction. You can for instance load the `car` package and do

```
R console
Anova(fml)
```

Then look at an interaction plot, you can write

```
interaction.plot(dat$Foodpl, dat$Sex, dat$AduDryW)
```

What do you conclude? Try to look at the summary output from your model and compare it to the interaction plot. Can you calculate the means for each `Sex:Foodplant` combination using the coefficients?

3. Read the data in **Dietstrs.txt**. These are the "modern life" data used in the lecture. Diet is coded as 1=Regular, 2=Junk, and stress is coded as 1=Low, 2=High. There is an additional column `DietStr` which codes for the treatment combinations: 1=Regular/Low, 2=Regular/High, 3=Junk/Low, 4=Junk/High. See to it that `Diet`, `Stress` and `DietStr` are made into factors. Next, make the two-way analysis with diet and stress shown during the lecture. You can use the commands

```
R console
fml <- lm(WtGain ~ Diet + Stress + Diet:Stress, data=dat)
Anova(fml)
```

Then make a one-way analysis with `DietStr` as factor. Try to devise a contrast for a planned comparison that corresponds to the main effects of Diet and Stress and Diet \times Stress interaction in the two-way analysis of variance. Try

```
cntr <- cbind(c(1, 1, -1, -1), c(1, -1, 1, -1), c(1, -1, -1, 1))
fm2 <- lm(WtGain ~ DietStr, data=dat, contrasts=list(DietStr=cntr))
summary(fm2)
```

You should get the same result as from the two-way analysis of variance, although the summary function gives you t-values, whereas the Anova function gives you the corresponding F-values (the p-values should be identical).

4. Read the data in **Buttfat.txt**. The file contains butterfat percentages for mature and young cows (10 of each) for each of five breeds. The codes are as follows. Breeds: 1=Ayrshire, 2=Canadian, 3=Guernsey, 4=Holstein-Friesian, 5=Jersey; Age classes: 1=Mature, 2=Two-year old (young). See to it that Breed and AgeCls are made into factors, for instance like this

R console

```
dat$Breed <- factor(dat$Breed, levels=c(1, 2, 3, 4, 5),
labels=c("Ayr", "Can", "Guern", "HolFrie", "Jers"))
dat$AgeCls <- factor(dat$AgeCls, levels=c(1, 2),
labels=c("Mature", "Young"))
```

Answer the following questions. Is there an Age class \times Breed interaction? Do young and mature cows differ in butterfat percentage? Do the breeds differ in butterfat percentage? You may need to transform the Buttfat variable to achieve homogeneous variances. A hint: the inverse transformation does the job (to get the best possible Box-Cox transformation, run the command `boxcox(Buttfat ~ Breed * AgeCls, data=dat)`, which will show that a transformation with exponent of about -1.4, instead of -1 as for the inverse, is the best, but the inverse transformation is good enough). Having done this, fit the model

R console

```
fm <- lm(InvButt ~ Breed + AgeCls + Breed:AgeCls, data=dat)
```

and analyze it with the `Anova` function (you need to load the `car` package).

It might also be useful to perform some post-hoc testing concerning which breeds differ from which. The Tukey HSD test might be used by loading the `multcomp` package. This function does not handle interactions so well. So before applying the command below be sure to redefine your model by

dropping the interaction term. Now, assuming you called your reduced fitted model `fm`, give the command:

```
cmp <- glht(fm, linfct = mcp(Breed = "Tukey"))  
  
summary(cmp)  
  
plot(cmp)
```

Which breeds differ significantly from each other and which do not?

5. Read the data in **Jimson.txt**. They consist of the length/width ratio (`LenWid`) of second seedling leaves of two types of Jimsonweed, called globe (coded as 1) and nominal (coded as 2). Three seeds of each type were planted in 16 pots (`Pot` is the pot identification number). You might assign G and N as level labels for the types and you also need to make `Pot` into a factor:

R console

```
dat$Type <- factor(dat$Type, levels=c(1, 2), labels=c("G", "N"))  
  
dat$Pot <- as.factor(dat$Pot)
```

Begin by restricting the dataset to only the first two pots

```
dat1 <- subset(dat, (Pot=="16533") | (Pot=="16534"))
```

The question is now whether the data from these two pots are sufficient to conclude that globe and nominal differ in length/width ratio. Do a two-way analysis of variance on `LenWid` with `Type` and `Pot` as factors (`Pot` is a blocking factor). First treat `Pot` as a fixed effect and then repeat the analysis with `Pot` as a random effect. Which type of analysis do you think is to be preferred? The fixed effect analysis you could do like this

```
fm1 <- lm(LenWid ~ Type + Pot + Type:Pot, data=dat1)  
  
Anova(fm1)
```

and the random effect analysis like this

```
fm2 <- lmer(LenWid ~ Type + (1|Pot) + (1|Type:Pot), data=dat1)  
  
Anova(fm2)
```

To help you think about this question, it might be of interest to look at the entire dataset, and in particular to check for the presence of a `Type×Pot` interaction. The effect that you are actually interested in is if one Type has a different Length-width ratio than the other. If the difference between the types are consistent across all the 16 pots (or say in 13 out of the 16) we would feel quite

convinced that there is indeed a difference. But with only 2 pots, how sure can you be? Your number of observations could be considered to be the two pots, and not the total number of plants grown from them. Using the dataset with all 16 pots the sample size for the observation of variance between pots is increased substantially, and there is hardly any, so the main effect of Type is highly significant. Notably, there are still uncertainties about how reliable the p-values from mixed models based on maximum likelihood algorithms are when the number of observations of the random effects are small.

Now try the same using the full dataset. What happened? It seems that there is very little variation explained by this effect, so is it that important? With the facts at hand (from looking at all 16 pots) we could see that there indeed seems to be very little variance in the difference between type G and N among pots, but you would not have known this with only the two pots.