

Biometry. Lecture 21

Alexey Shipunov

Minot State University

April 22, 2015



1 Two-dimensional statistics

- Wordcloud
- Logistic regression
- ANalysis Of VAriation (ANOVA)



```
> setwd("<working folder>")  
or  
"Change dir"  
in menu!
```

On Mac, be sure that startup option is working: `getwd()`
(`getwd()` checks if R is in working folder, `dir()` checks the folder
content)



Two-dimensional statistics

Wordcloud



Wordcloud: fashionable replacement of pie chart

```
> cplants <- read.table("data/cplants.txt", sep="\t", quote="", as.is=T)
> library(wordcloud)
> cfams <- data.frame(table(cplants$V3))
> set.seed(5); wordcloud(cfams[,1], cfams[,2])
```



Two-dimensional statistics

Logistic regression



Numeric influence but categorical response

- What if response is binary?
- It is possible to convert success/failure to the **probability of success** and then apply a **generalized linear model**



Analysis of logistic regression

```
> lo <- read.table("http://ashipunov.info/data/logit.txt")
> head(lo); str(lo)
> lo.logit <- glm(formula=V2 ~ V1, family=binomial,
+ data=lo)
> summary(lo.logit)
```



Visualizing logistic regression

```
> new.points <- seq(min(lo$V1), max(lo$V1), length.out=14)
> predicted.points <- predict(lo.logit,
+ list(V1=new.points), type="response")
> success <- as.numeric(lo$V2) - 1
> plot(success ~ V1, data=lo)
> lines(new.points, predicted.points)
```



Logistic regression example II: poisoning

- Caesar or tomatoes?
- High significance of both terms could be a result of coincidence (people often took these things together)
- If we construct a logistic model and then update it (taking out one of two terms), AIC will show which model is better.



Poisoning analysis

```
> tox <- read.table("http://ashipunov.info/data/tox.txt",
+ h=TRUE)
> tox.logit <- glm(formula=I(2-ILL) ~ CAESAR + TOMATO,
+ family=binomial, data=tox)
> tox.logit2 <- update(tox.logit, . ~ . - TOMATO)
> tox.logit3 <- update(tox.logit, . ~ . - CAESAR)
> tox.logit$aic
> tox.logit2$aic # lowest!
> tox.logit3$aic
> summary(tox.logit2) # highly significant!
```

Caesar!



Two-dimensional statistics

ANalysis Of VAriation (ANOVA)



Categorical influence and numerical response

- What if there are three species in horsetail data? How to compare the diameter of stem of them all?
- Paired comparisons are “temptation” for p-value, and Bonferroni correction is sometimes of no help.
- There is a solution: analysis of variation (ANOVA) and its non-parametric twin, Kruskal-Wallis test.



ANOVA null and alternative hypotheses and assumptions

- Null is that all groups are not different, alternative is that **at least one group is different from all others**.
- All variables **should be normally distributed**. Small deviations from normality are typically accepted but Kruskal-Wallis test is preferable for all “non-normal” data.
- Variances should be at least **similar**. This is possible to overcome with `oneway.test()`
- Actually, t-test is just an ANOVA for only two groups. Wilcoxon test, in turn, is a special case of Kruskal-Wallis test for two groups only.



Introductory example: eight horsetails

```
> eq8 <- read.table("http://ashipunov.info/data/eq8.txt",  
+ h=T)  
> str(eq8); head(eq8)  
> plot(DIA.ST ~ SPECIES, data=eq8)  
> eq8.anova <- lm(DIA.ST ~ SPECIES, data=eq8)  
> anova(eq8.anova)  
# If variables are not normal:  
> kruskal.test(DIA.ST ~ SPECIES, data=eq8)
```



Finishing...

Save your commands!

(`savehistory(<today's date>.r)` or File -> Save as... on Mac)



Summary: most important commands

- `lm()`—estimates the linear regression model and many other models (like ANCOVA)
- `predict()`—predict values with model
- `glm()`—estimates the logistic regression model and many others



For Further Reading



A. Shipunov.

Biometry [Electronic resource].

2012—onwards.

Mode of access:

http://ashipunov.info/shipunov/school/biol_240



A. Shipunov, and others.

Visual statistics. Use R!

Ongoing translation from Russian.

