

Biometry. Lecture 14

Alexey Shipunov

Minot State University

March 23, 2016



1 Two-dimensional statistics

- Hypotheses and tests
- Tests for the independence of two variables



```
> 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

Hypotheses and tests



Hypotheses are cornerstones of science

- The inferential science is based on hypotheses construction and calculation of their probability.
- The simplest approach is to establish null hypothesis and reject it if needed.
- More complicated approach is to consider null and alternative hypotheses together.



Statistical errors

- Type I error is a false alarm: we accept alternative when null is true
- Type II error is a carelessness: we accept null when alternative is true



Level of significance

- The probability to have greater or equal effect when null hypothesis is true is a p-value
- We may ignore this probability if it is too low, in other words, below the level of significance
- The level of significance is a matter of experience and agreement, it could be 0.05, but sometimes also 0.1 and 0.01
- p-value is related with Type I error



Two-dimensional statistics

Tests for the independence of two variables



What is tested?

- Null: difference equal to 0 \approx similar \approx related \approx samples came from same population
- Alternative: difference not equal to 0 \approx different \approx non-related \approx samples came from different populations



p-value

- p-value is the probability to have equal or greater effect in case if null hypothesis is true
- p-value is related with the accepted level of Type I statistical error (false alarm): the bigger is the significance level of p-value, the more false alarms we accept and at the same time, the more sensitive is our research
- In biology, the most common level of significance is 0.05
- As a rule, if p-value is equal or less then 0.05, we reject the null hypothesis, if more – we stay with null hypothesis



Paired and non-paired

- Paired: came from one set of objects (e.g., measurements done at different time)
- Non-paired: do not belong to one set of objects



Artificial example of the paired test

```
> a <- 51:59  
> set.seed(1); t.test(a, (a+rnorm(9)), paired=T)
```

We introduced here a random noise (`rnorm()` function)



Two sample tests for Keller's data

```
> ph <- read.table(  
+ "http://ashipunov.info/data/phaseolus.txt", h=T)  
> Normality <- function(x, p=.05)  
+ {  
+ ifelse(shapiro.test(x)$p.value > p, "NORMAL", "NOT NORMAL")  
+ }  
> sapply(ph, Normality) # all normal!  
> t.test(ph$EXPER.2, ph$EXPER.1, alt="greater", paired=T)  
> with(ph, t.test(CONTROL.2, CONTROL.1, alt="greater",  
+ paired=T))
```

In first two cases, we use “greater” because we **already** know that leaves **grow**.



Homoscedasticity

```
> a <- 51:59
> b <- 1:9
> x <- rep(5, 9)
> t.test(a, b)
> t.test(b, x)
```

Homoscedasticity, similarity of variance (like in a and b but not like in b an x) is an important assumption of all parametric two-dimensional tests. In R, the Welch correction for **non-homogeneity of variance** is by default applied inside `t.test()`



Two main questions

- Normal?
 - Yes: `t.test()`
 - No: `wilcox.test()`
- Paired?
 - Yes: `(..., paired=T)`
 - No: `(...)`



“Classical” sleep data and model formula

```
> str(sleep)
> boxplot(extra ~ group, data=sleep)
> t.test(extra ~ group, data=sleep)
```

sleep is a data in so-called long format, `extra ~ group` is a **model formula** of response ~ factor form.

For t-test, “group” should have exactly 2 levels!



Model formula for leaves data

```
> leaves12 <- stack(ph[,1:2])  
> leaves12  
> t.test(values ~ ind, data=leaves12, paired=T)
```

`stack()` converts from short to long form



Air quality data in May and August

```
> str(airquality)
> sapply(airquality, Normality)
> boxplot(Ozone ~ Month, data=airquality,
+ subset=Month %in% c(5,8))
> wilcox.test(Ozone ~ Month, data=airquality,
+ subset=Month %in% c(5,8))
```

`unstack()` converts from long to short form
`%in%` is a selection operator



Finishing...

Save your commands!

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



Summary: most important commands

- `binom.test()` and `prop.test()` —tests for the equality of proportions
- `t.test()` —paired and non-paired two-sample parametric test
- `wilcox.test()` —paired and non-paired two-sample non-parametric test



For Further Reading



A. Shipunov.

Biometry [Electronic resource].

2012—onwards.

Mode of access:

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



A. Shipunov, and many others.

Visual statistics. Use R!

2016—onwards.

Mode of access: http://ashipunov.info/shipunov/school/biol_240/en/visual_statistics.pdf

