

# Chi-squared test and ANOVA

September 30, 2020

## 1. Floral selection of a bumblebee species

By following the foraging activity of bumblebees *Bombus impatiens* on a site, you note the number of bumblebee visits on four genera of plants, as well as the proportion of flowers of each genus on the site.

Genus	Number of visits	Proportions of the flowers of the site
<i>Rubus</i>	8	0.12
<i>Solidago</i>	8	0.24
<i>Trifolium</i>	18	0.33
<i>Vaccinium</i>	11	0.31

The null hypothesis for this study is that *B. impatiens* visits each genus in proportion to its prevalence on the site.

- According to the null hypothesis, what are the expected frequencies for visits to each genus of plant?
- Test the null hypothesis with the `chisq.test` function in R, with a significance level of 5%. If the null hypothesis is rejected, which genus or genera are more or less visited than expected?

## 2. Foraging activity of three bumblebee species

On the same site as the previous exercise, you observe the foraging activity of two other bumble bee species (*B. affinis* and *B. ternarius*). Here is the contingency table showing the number of visits by bumblebee species and plant genus.

	<i>Rubus</i>	<i>Solidago</i>	<i>Trifolium</i>	<i>Vaccinium</i>
<i>B. affinis</i>	10	9	15	8
<i>B. impatiens</i>	8	8	18	11
<i>B. ternarius</i>	20	4	6	5

- What null hypothesis can you test from this table? What is the alternative hypothesis?
- Create a matrix representing this table in R, then test the null hypothesis mentioned in (a) with the `chisq.test` function, with a significance level of 5%.
- Based on the test results in (b), what is the number of degrees of freedom of the  $\chi^2$  distribution? How is this value calculated?
- How can you check the expected frequencies according to the null hypothesis, as well as the residuals?
- If the null hypothesis is rejected, which bumblebee-plant pair has the most positive residual, and which has the most negative residual? How do you interpret these residuals?

### 3. Characteristics of cabbage varieties

The `cabbages` dataset included in the `MASS` package shows the weight in kg (*HeadWt*) and the vitamin C content (*VitC*) of cabbages according to the cultivar (*Cult*) and the planting date. There are 10 replicates for each of the six combinations of cultivar and date.

```
library(MASS)
str(cabbages)
```

```
## 'data.frame': 60 obs. of 4 variables:
## $ Cult : Factor w/ 2 levels "c39","c52": 1 1 1 1 1 1 1 1 1 1 ...
## $ Date : Factor w/ 3 levels "d16","d20","d21": 1 1 1 1 1 1 1 1 1 1 ...
## $ HeadWt: num 2.5 2.2 3.1 4.3 2.5 4.3 3.8 4.3 1.7 3.1 ...
## $ VitC : int 51 55 45 42 53 50 50 52 56 49 ...
```

- Select the subset of data corresponding to cultivar `c52`. For each of the two numeric variables (*HeadWt* and *VitC*), produce a boxplots graph showing the distribution of that variable for each planting date. Before even performing the ANOVA, do you think that the assumptions of the model (especially the equality of variances) will be respected in each case?
- Choose one of the two variables (*HeadWt* or *VitC*) that best fits the ANOVA model based on your result in (a). Perform a one-way ANOVA and determine if the planting date has a significant effect ( $\alpha = 0.05$ ).
- According to Tukey's range test, between which planting dates are there significant differences ( $\alpha = 0.05$ )? What is the estimate of those significant differences?