

# Overdispersion and Quasibinomial Logistic Regression

Logistic regression – Stat 230

In this class, you learn about extra-binomial variation (a case of overdispersion) and how we can use quasibinomial logistic regression to adjust our inferences to handle this issue. This topic is not discussed in your book.

## Data overview

*Orobanch*e is a genus of parasitic plants without chlorophyll that grow on the roots of flowering plants. In an experiment, seeds from two varieties, *Orobanch*e aegyptiaca 75 and *Orobanch*e aegyptiaca 73, were brushed with extract from either a cucumber root or a bean root. Researchers recorded the number of seeds that eventually germinated.

```
orobanche <- read.csv("http://aloy.rbind.io/data/orobanche.csv")
```

The variables in the data set are:

- **y**: count of seeds germinated
- **n**: number of seeds
- **variety**: either oa75 or oa73
- **root**: either cucumber or bean

Notice that both explanatory variables are categorical (i.e., factors).

**Task 1.** Does there appear to be a relationship between the proportion of seeds that germinate and variety? What about root? To answer this, I recommend looking at both a summary table and boxplots of the estimated proportion by each variable.

To calculate a summary of counts for each combination of levels, use the commands `f`table() and `x`tabs()

```
f
```

table(xtabs(cbind(y, n) ~ variety + root, data = orobanche))

To create a **prop** column in the data set, you can run the below chunk:

```
orobanche$prop <- orobanche$y/ orobanche$n
```

## Additive model

To begin, let's fit an additive logistic regression model where `variety` and `root` are both included.

```
oro_glm1 <- glm(y/n ~ variety + root, data = orobanche, family = binomial, weights = n)
```

**Task 2.** Perform a deviance goodness-of-fit test. What do you conclude?

**Task 3.** Are there any concerning outliers? Check residual plots to investigate this.

## Adding an interaction term

You should have found evidence of lack of fit and that there were no concerning outliers, so we may be missing a term. The only other thing we can add here is an interaction. The below code chunk updates our first model to include an interaction term.

```
oro_glm2 <- update(oro_glm1, . ~ . + variety:root)
```

**Task 4.** Perform a deviance goodness-of-fit test. What do you conclude?

**Task 5.** Check the residual plots again to see if outliers are a problem.

## Investigating over-dispersion

Since there aren't any more terms we can add and outliers aren't a problem, we suspect over-dispersion.

**Task 6.** Do the *ad hoc* calculation of the estimate of the dispersion parameter. Is it "much" larger than 1?

**Task 7.** The below code to fit the interaction model via quasibinomial (i.e., this fits the quasibinomial model). How do the standard errors and test statistics change? Do any of the variables change in significance?

```
oro_glm3 <- glm(y/n ~ variety * root, data = orobanche, family = quasibinomial,
               weights = n)
```

**Task 8.** In the `summary()` output for `oro_glm3` there is a dispersion parameter given. Multiply the square root of this with the standard error of the interaction term in `oro_glm2`. Compare this to the standard error of the interaction term in `oro_glm3`. What do you notice?

## Comparing models

Let's work with the quasibinomial version now.

**Task 9.** Determine if we can remove the `variety:root` interaction from the model. To do this, fit the reduced quasibinomial logistic regression model and then use the `anova()` command to compare the models, specifying `test = "F"`. (Remember that in the quasibinomial case, the drop in deviance test now uses the F distribution!)