

Statistics II

Logistic Regression

Çağrı Çöltekin

University of Groningen, Dept of Information Science



May 29, 2013

Exam date & time: **June 21, 10:00–13:00.**
(The same day/time planned at the beginning of the semester.)

More on the exam next week.

Previously in this course...

So far...

- Simple regression** Single (typically) numeric predictor, numeric response.
 - Multiple regression** Multiple predictors, numeric response.
 - ANOVA** Single categorical predictor, numeric response.
 - Factorial ANOVA** Multiple categorical predictor, numeric response.
 - Repeated-measures ANOVA** Like single/factorial ANOVA with dependent observations.
- This week: logistic regression.

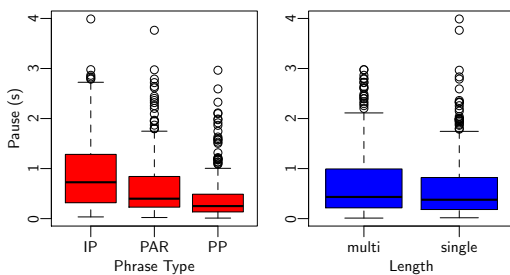
Two-way repeated measures ANOVA: an example

We return to last week's production experiment example.

- ▶ We want to analyze the pauses after three different phrase types: IP, PP and PAR.
- ▶ In addition, we are also interested in effect of length of the phrase. Particularly, whether it is single word or multi word.
- ▶ We have a 3×2 repeated measures ANOVA.
- ▶ We use 5 different examples for each experimental condition, resulting in $3 \times 2 \times 5$ utterances recorded for each of our 30 participants.
- ▶ Note: we need to randomize the order of the utterances for each subject. Why?

Previously in this course...

RM-ANOVA example: the data



Previously in this course...

RM-ANOVA example: assumptions

- ▶ Check each (3×2) group for normality.
- ▶ Check for sphericity:

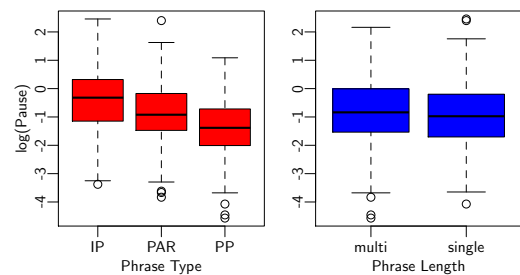

```

'Mauchly's Test for Sphericity'
  Effect      W      p p<.05
 2 phrase.type 0.4153665 0.455
 4 phrase.type:length 0.3779318 0.121
            
```
- ▶ Why don't we have a sphericity test for 'length'?
- ▶ What if sphericity assumption was violated?

Previously in this course...

Previously in this course...

RM-ANOVA example: the data after log transform



Previously in this course...

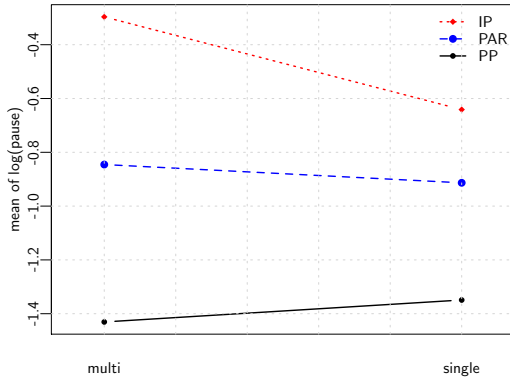
RM-ANOVA example: main ANOVA results

R output:

```

> summary(aov(log(pause) ~ phrase.type*length + Error(subj)))
Error: subj
      Df Sum Sq Mean Sq F value Pr(>F)
Residuals 29  57.27   1.975
Error: Within
      Df Sum Sq Mean Sq F value Pr(>F)
phrase.type  2  153.3   76.67  74.893 <2e-16 ***
length      1    2.9    2.86   2.790 0.0952 .
phrase.type:length  2    9.4    4.68   4.575 0.0106 *
Residuals    865  885.5    1.02
            
```

RM-ANOVA example: interaction



Logistic regression: motivation

- ▶ For all methods discussed in this class, the response variable is numeric.
- ▶ Sometimes we want to analyze/predict categorical responses:
 - ▶ Alive or dead.
 - ▶ Grammatical or ungrammatical.
 - ▶ Correct or incorrect.
 - ▶ Pass or fail.
 - ▶ Win or loose.
 - ▶ Present or absent.

Logistic regression is an extension of regression for categorical (typically binary) response variables.

Simple regression: a refresher

$$y_i = a + bx_i + e_i$$

- y is the *response* (or outcome, or dependent) variable. The index i represent each unit observation/measurement (sometimes called a 'case').
- x is the *predictor* (or explanatory, or independent) variable.
- a is the intercept.
- b is the slope of the regression line.
- $a + bx$ is the *deterministic* part of the model (we sometimes use \hat{y}).
- e is the *residual*, error, or the variation that is not accounted for by the model. Assumed to be (approximately) normally distributed with 0 mean (e_i are assumed to be i.i.d).

Multiple regression

$$y_i = a + \underbrace{b_1x_{i,1} + b_2x_{i,2} + \dots + b_kx_{i,k}}_{\hat{y}} + e_i$$

- ▶ Multiple regression is a generalization of the simple regression, where we predict the outcome using multiple predictors.
- ▶ **Multicollinearity** causes problems in estimation and interpretation of multiple-regression models.
- ▶ **Model selection** (finding a model that fits the data well, but not more complex than necessary) is important.

RM-ANOVA: a summary and further directions

- ▶ RM-ANOVA is an efficient method for experimentation.
- ▶ ANOVA design can be 'mixed', where one can analyze multiple between- as well as within-subject factors.
- ▶ RM-ANOVA is difficult to please, except in carefully constructed experimental conditions.
- ▶ In our example today, we ignored inter-phrase repeated measures. This is a common problem known as 'language-as-a-fixed-effect fallacy' in psycholinguistics.
- ▶ A newer alternative is to use linear mixed-effect models (or multilevel models).

Logistic regression: some examples

- ▶ survival after a surgery depending on age, length of surgery, ...
- ▶ whether purchase occurs depending on age, income, website characteristics, ...
- ▶ whether speech errors occur depending on alcohol level
- ▶ when linguistic rules apply (final [t] in Dutch) depending on speed of utterance, stress, social group, ...
- ▶ whether one votes to a political party (or not) depending on age, income, ethnicity, ...

Regression assumptions

- independence** cases (observations) should be independent.
- linearity** the relation between 'y vs. x' should be linear.
- normality** residuals should be normally distributed with 0 mean.
- constant variance** residual variance should be constant.

Multiple regression: example

We will try to determine reaction time differences in a lexical decision task in two experimental conditions 'C₁' and 'C₂', and also investigate the effect of age. Here is how our data looks like,

Reaction time (ms)	condition	age
563	C ₁	38
562	C ₂	35
550	C ₁	20
573	C ₂	31
⋮	⋮	⋮

Multiple regression: example (2)

Running multiple regression (R output):

```
Call:
lm(formula = log(rt) ~ age + condition)
Residuals:
    Min       1Q   Median       3Q      Max
-0.0088400 -0.0028946 -0.0004367  0.0024968  0.0099490
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  6.3246836  0.0035730 1770.13 < 2e-16 ***
age           0.0001977  0.0001024   1.93  0.0613 .
conditionC2  0.0167303  0.0014068  11.89 3.33e-14 ***
---
Residual standard error: 0.004374 on 37 degrees of freedom
Multiple R-squared:  0.7927,    Adjusted R-squared:  0.7815
F-statistic: 70.74 on 2 and 37 DF,  p-value: 2.276e-13
```

Note: you should do model diagnostics before interpreting the model.

An example for logistic regression

We go back to our lexical decision task example. In similar experiments typically we also whether the reaction was correct or incorrect.

Our data actually looks like this:

Reaction time (ms)	condition	age	correct
563	C ₁	38	1
562	C ₂	35	1
550	C ₁	20	1
573	C ₂	31	0
⋮	⋮	⋮	⋮

Our aim is to predict probability of correct decision.

Transforming the response variable

- ▶ Instead of predicting the probability, p , we can predict **odds**.
- ▶ Odds in favor of an event (e.g., FC Groningen winning against Ajax) is defined based on p , the probability of the event occurring, as

$$\text{odds} = \frac{p}{1-p}$$

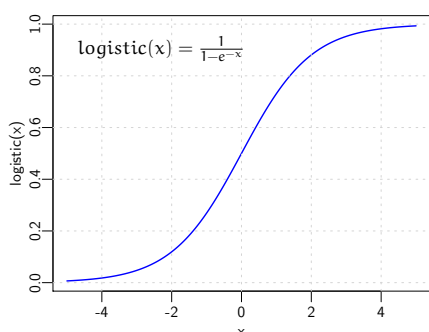
- ▶ odds are bounded between 0 and infinity, $0 \leq \text{odds} \leq \infty$.
- ▶ further if we take the logarithm of odds,

$$\log \frac{p}{1-p} = \text{logit}(p)$$

Bounds of logit is $-\infty \leq \text{logit}(p) \leq +\infty$

But, what about the name 'logistic'?

Logistic function is the inverse of the logit function:



Logistic Regression

- ▶ Logistic regression is a special case of regression, where a binary outcome is predicted based on a number of predictors.
- ▶ The main trick is to predict probability of an event, rather than the outcome directly. This allows converting a categorical variable to a numeric variable (probabilities).
- ▶ One needs to overcome difficulties due to regression assumptions.
 - ▶ Probabilities are strictly bounded between 0 and 1.
 - ▶ Error (residuals) with binary (or proportion) response is not normally distributed.

Trying ordinary regression

Model the outcome directly:

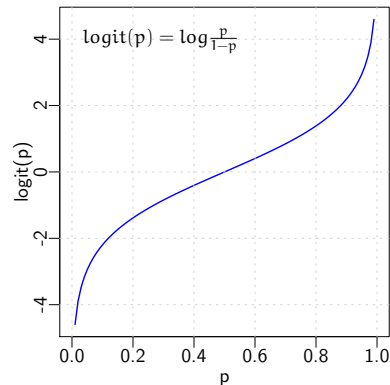
```
lm(formula = correct ~ age + condition)
Residuals:
    Min       1Q   Median       3Q      Max
-0.81335 -0.07136  0.04161  0.17828  0.36306
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.461957  0.243223   1.899  0.0653 .
age           0.016037  0.006973   2.300  0.0272 *
conditionc2 -0.209907  0.095766  -2.192  0.0348 *
...
```

$$P(\text{correct}) = 0.462 + 0.016 \times \text{age} - 0.210 \times \text{condition}$$

Estimated probability of correct response to condition 1 from a 40-year old is: $0.462 + 0.016 \times 40 - 0.210 \times 0 = 1.102$.

Response variable (probability) is not bounded between 0 and 1.

The logit function



Trying ordinary regression with logit transform

```
lm(formula = logit(correct) ~ age + condition)
Residuals:
    Min       1Q   Median       3Q      Max
-5.9595 -0.5229  0.3049  1.3063  2.6602
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.27875  1.78213  -0.156  0.8766
age           0.11751  0.05109   2.300  0.0272 *
conditionc2 -1.53802  0.70169  -2.192  0.0348 *
Residual standard error: 2.182 on 37 degrees of freedom
Multiple R-squared:  0.2501,    Adjusted R-squared:  0.2095
F-statistic: 6.169 on 2 and 37 DF,  p-value: 0.004873
```

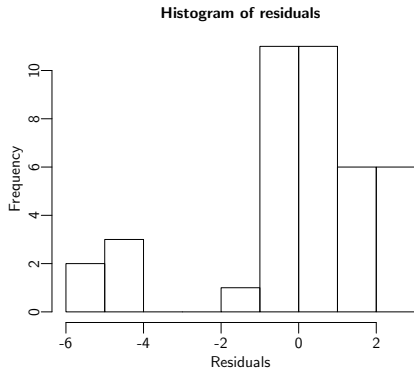
Now our estimate is:

$$\text{logit}(P(\text{correct})) = -0.279 + 0.118 \times \text{age} - 1.538 \times \text{condition}$$

The model is more correct, but interpretation is difficult.

Logistic regression: are we there yet?

Let's look at the residual distribution:



Logistic regression

Logistic regression is similar to (multiple) regression, differences are:

- ▶ Categorical response variable.
- ▶ Non-normal errors.
- ▶ Relationship is non-linear (linear after the transformation).
- ▶ Estimation is (typically) done using **maximum likelihood estimation**. Least-squares regression is not applicable.

Logistic regression: example

We return to the same example, this time fitting a correct model:

```
glm(formula = correct ~ age + condition, family = binomial)
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-2.18599  0.00002  0.00005  0.32773  1.09422
Coefficients:
(Intercept)  14.7647  3620.9079  0.004  0.9967
age           0.1904  0.1073  1.775  0.0759 .
condition2   -19.1351  3620.9069 -0.005  0.9958 .
(Dispersion parameter for binomial family taken to be 1)
Null deviance: 30.142 on 39 degrees of freedom
Residual deviance: 17.708 on 37 degrees of freedom
AIC: 23.708
Number of Fisher Scoring iterations: 19
```

Estimated equation is:

$$\text{logit}(P(\text{correct})) = 14.7647 + 0.1904 \times \text{age} - 19.1351 \times \text{condition}$$

Logistic regression: inference and model fit

- ▶ In logistic regression the significance testing for coefficients are carried out using z-statistics (also known as 'Wald statistic' or 'Wald's z').
- ▶ We do not have r^2 as measure of fit (approximations exists, but none have all the properties of r^2).
- ▶ $-2LL$ (-2 times Log Likelihood, smaller better) is the measure of fit in maximum likelihood estimation.
- ▶ Log likelihood is not bounded like r^2 , comparisons are valid only on the same data set.
- ▶ $-2LL$ is approximately χ^2 distributed, a model can be tested against null model (similar to F-test for least-squares regression).

Regression for binary response: problems so far

- ▶ Binary data is strictly bounded, probabilities smaller than 0 and larger than 1 are meaningless.
 - ▶ We solve this by transforming the response variable using logit function (other transformations are also possible).
- ▶ When response is binary, the residuals are not normally distributed with 0 mean.
 - ▶ We estimate the coefficients without assuming normally distributed error.
 - ▶ The cost of is also giving up the familiar (and uniquely determined) least-squares regression. Estimation is done with more complicated procedures (typically maximum likelihood).
 - ▶ More precisely, logistic regression is an instance of **generalized linear models** (GLMs).

Maximum likelihood estimation

$$\text{logit}(p_i) = \alpha + b_1 x_{1,i} + \dots + b_k x_{k,i} + e_i$$

- ▶ Maximum likelihood estimation tries to find the set of model parameters, or coefficients, α, b_1, \dots, b_k , which make the data most likely (or minimizes the error).
- ▶ MLE is an iterative search for the optimum parameter values. There is no exact solution.
- ▶ In some cases, MLE may fail to find a solution.
- ▶ If errors are normally distributed, MLE is equivalent to least-squares estimation.

Logistic regression: interpreting the coefficients

$$\log \frac{P(\text{correct})}{1 - P(\text{correct})} = 14.7647 + 0.1904 \times \text{age} - 19.1351 \times \text{condition}$$

Slope of a variable indicates the change in log odds for unit-change in the predictor (while other predictors kept constant). For example, one unit change in 'age', say from 40 to 41, means

$$\begin{aligned} \log \frac{P(\text{correct}_{41})}{1 - P(\text{correct}_{41})} - \log \frac{P(\text{correct}_{40})}{1 - P(\text{correct}_{40})} &= 0.1904 \\ e^{\log \frac{P(\text{correct}_{41})}{1 - P(\text{correct}_{41})} - \log \frac{P(\text{correct}_{40})}{1 - P(\text{correct}_{40})}} &= e^{0.1904} \\ \frac{\text{odds}(\text{correct}_{41})}{\text{odds}(\text{correct}_{40})} &= 1.209733 \end{aligned}$$

When age is increased form 40 to 41, odds of correct response increase 1.2 times. Note: the change is non-linear.

Logistic regression: where things may go wrong

- ▶ **Overdispersion** is the case when variance in the data is higher than expected.
 - ▶ Logistic regression requires response to follow **binomial** distribution.
 - ▶ Variance of binomial distribution is predictable from it's mean.
- ▶ In case of **complete separation**, i.e., when predictors perfectly separate cases of success and failure, logistic regression parameters cannot be estimated.
- ▶ Logistic regression may also fail to find a solution, especially when data is not evenly distributed.
- ▶ Logistic regression is multiple regression, other problems, such as collinearity is also problems for logistic regression.

Logistic regression: another example

We would like to guess whether a child would develop dyslexia or not based on a test applied to pre-verbal children. Here is a simplified problem:

- ▶ We test children when they are less than 2 years of age.
- ▶ The hypothesis is that the test may be relevant in diagnosing dyslexia.
- ▶ We observe the same children when they are in the school age, and note whether they are diagnosed with dyslexia or not.
- ▶ Our data looks like the following:

Test score	Dyslexia
8.2	0
2.2	1
6.2	1
⋮	⋮

* The research question is from the longitudinal study by Ben Maasen and his colleagues. Data is fake as usual.

Example: interpretation

$$\text{logit}(p(\text{dyslexia})) = 4.4373 - 0.9419 \times \text{score}$$

score	logit = $\log \frac{p}{1-p}$	odds = $\frac{p}{1-p}$	p
1	3.50	32.96	0.97
3	1.61	5.01	0.83
5	-0.27	0.76	0.43
7	-2.16	0.12	0.10
9	-3.10	0.05	0.04

Summary

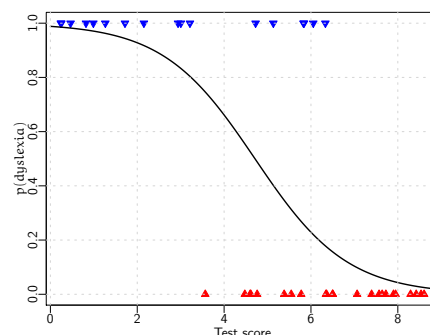
- ▶ Logistic regression is applicable when response is binomial.
- ▶ Two major differences from ordinary least-squares regression:
 - ▶ Response variable is strictly bounded. logit transformation make sure regression output is mapped to range [0, 1].
 - ▶ Errors are not normal. GLM framework allows non-normal errors. However, we use maximum likelihood estimation instead of least squares.
- ▶ Problems to watch out for:
 - ▶ overdispersion
 - ▶ complete separation, or insufficient data for estimation
 - ▶ like in multiple regression: multicollinearity
- ▶ Extensions
 - ▶ Multinomial logistic regression, when there are more than two categories of the response variable.
 - ▶ Generalized linear mixed-effect models when observations are not independent.

Example: the analysis

```
glm(formula = dys ~ score, family = binomial, data = d)
Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.6558  -0.5927  -0.2519   0.3590   1.8567
Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  4.4373    1.5758   2.816  0.00486 **
score       -0.9419    0.2920  -3.225  0.00126 **
---
(Dispersion parameter for binomial family taken to be 1)
Null deviance: 54.548 on 39 degrees of freedom
Residual deviance: 30.337 on 38 degrees of freedom
AIC: 34.337
Number of Fisher Scoring iterations: 5
```

$$\text{logit}(p(\text{dyslexia})) = 4.4373 - 0.9419 \times \text{score}$$

Visualizing the logistic regression



Next week

- ▶ A summary of all methods.
- ▶ When to use which analysis.
- ▶ What to do when assumptions fail.
- ▶ Beyond hypothesis testing: effect sizes.
- ▶ Time for your questions.