

---

# Applied Data Analysis

---

## R-Laboratory 9

---

### Logistic Regression – Baseline-Category Logit Model – Cumulative Link Model

---

Do not use functions from additional R packages (except when it is explicitly stated in the Task, Hint or list of useful packages and functions).

Useful packages and functions:

- |                    |                 |                  |
|--------------------|-----------------|------------------|
| • pROC             | • confint()     | • xtabs()        |
| • pROC::roc()      | • ftable()      | • deviance()     |
| • pROC::plot.roc() | • VGAM          | • df.residuals() |
| • pROC::auc()      | • VGAM::vglm()  | • matplot()      |
| • pchisq()         | • VGAM::step4() |                  |

### Task 26

- Download the file *FieldGoal.csv* from RWTHmoodle and load it as a data frame into your workspace.
- Fit a logistic regression model that predicts `Good.` using `Dist` as explanatory variable and plot the predicted probabilities for a good kick as a function of `Dist`.
- Calculate the percentage of correct classified observations according to the model in (b).
- Test at significance level  $\alpha = 0.05$  the hypotheses

$$H_0 : \beta_1 = 0 \quad \text{versus} \quad H_1 : \beta_1 \neq 0$$

in the model from (b), where  $\beta_1$  is the parameter of the attribute `Dist`. Base your test decision on the Wald's test statistic.

- Compute 90% profile likelihood confidence intervals for the parameters of intercept and `Dist`, as well as for the probability of a successful kick when the distance to the goal is 19, 39 and 64 yards.
- Fit a logistic regression model, that predicts `Good.` using the attributes `Dist`, `Blk.`, `Pressure`, `Roof.type`, `Altitude` and `Field`. Select the model with smallest AIC using a backwards stepwise selection algorithm. Calculate the percentage of correct classified observations according to this model and compare it with the percentage from (e).

- (g) Plot (in a single figure) the ROC curves for the model in (b) and the selected model in (f). Evaluate for both models the area under the curve (AUC).

## Task 27

- (a) Load the dataset `Hoyt` into your workspace and transform it to a 'flat' contingency table, with 4 columns, each corresponding to a level of the nominal attribute `Status`.

*Hint:* You may use the function `ftable` with the argument `col.vars="Status"`.

*Explanation of ftable:* The function `ftable` creates a flat contingency table which means that the usual information contained in a contingency table are re-arranged as a matrix. The rows and columns of this matrix correspond to the combinations of the levels of the involved factors. `ftable` can be applied to a data frame or a contingency table generated with `xtabs` for example, where `xtabs` creates a contingency table from cross-classifying all involved factors which is returned as a table and not re-arranged as a matrix

- (b) Fit a baseline-category logit model for `Status` with explanatory variables `Rank`, `Occupation`, `Sex` and their pairwise interaction terms.

*Hint:* You may use the function `vglm` from the package `VGAM` with the key command

`vglm(formula,family,data).`

You can fit grouped data by placing a matrix as response in `formula`.

- (c) Given the model you fitted in (b), test at significance level  $\alpha = 0.05$ , if there is any influence of the explanatory variables on `Status`.
- (d) Given the model you fitted in (b), test the significance of the effect of `Sex:Rank` on `Status` at significance level  $\alpha = 0.05$ .
- (e) Select the model with smallest AIC and with smallest BIC using a backwards stepwise selection algorithm starting with the model from (b).

## Task 28

- (a) Load the dataset `Vietnam` into your workspace and transform it to a 'flat' contingency table, with 4 columns, each corresponding to a level of the ordinal attribute `response`.
- (b) Compute the baseline-category sample log-odds of `response` at all levels of `sex` and `year` and the cumulative sample log-odds and plot them as a function of `year` (of study) for each gender, with genders indicated by different colors.

*Hint:* You may use the function `matplot`.

- (c) Fit a cumulative logit model for `response` with explanatory variables `sex` and `year`.  
*Hint:* Use the arguments `parallel=FALSE` and `link="logitlink"` within the `family` argument.

- (d) Fit a cumulative probit model for `response` with explanatory variables as in (c).  
*Hint:* Use the arguments `parallel=TRUE` and `link="probitlink"` within the `family` argument.

- (e) For both models in (b) and (c) compute the Pearson's  $X^2$  and the deviance in order to perform a goodness of fit test at significance level  $\alpha = 0.01$ .
- Hint:* You may use the function `predict` with argument `type="response"` to compute the estimated response probabilities for each factor level of **Status**, conditional on the values of the explanatory variables.