

## Problems in R

### Statistical Models and Procedures

- a) The `mtcars` data set in R contains technical specifications of different car models and their fuel consumption (miles per gallon; *mpg*).

- (a) Estimate the regression model  $M_{hp}$  in which *mpg* is predicted by horse power (*hp*):

$$mpg_i = \beta_0 + \beta_1 hp_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \text{ i.i.d.} \quad (1)$$

Create a scatter plot of *mpg* as a function of *hp*. Add the fitted regression line to the figure via `abline()`. Discuss whether there are any violations of the model assumptions based on the figure you just created and the plot method of the model object. Pay special attention to the “Residuals vs Fitted” portion of the latter plot.

- (b) Now estimate the model  $M_{hp^2}$ , in which the squared horse power is added as a predictor:

$$mpg_i = \beta_0 + \beta_1 hp_i + \beta_2 hp_i^2 + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \text{ i.i.d.} \quad (2)$$

Hint: Use `I(speed^2)` in the model formula. Add the predictions of this model to the figure you created previously; use `predict()` and `lines()` for this. Which of the two models  $M_{hp}$  and  $M_{hp^2}$  would you choose? Conduct an appropriate hypothesis test to aid with your decision.

- (c) The variable *disp* indicates the engine displacement of each car. Estimate the model  $M_{hp\_di}$ :

$$mpg_i = \beta_0 + \beta_1 hp_i + \beta_2 disp_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \text{ i.i.d.} \quad (3)$$

Interpret the Wald tests of the parameters in this model.

Now test the hypothesis that the effect of horse power and engine displacement on *mpg* is equal in size, i.e. in the model above:  $\beta_1 = \beta_2$ .

How can you reconcile the result of this hypothesis test with the Wald tests above?

- b) In a detection experiment, a person is presented with signal and noise trials and asked if he or she recognizes the signal:

	response	
trial	yes	no
noise	89	361
signal	330	120

A measure of sensitivity is  $d'$ , which is defined within signal detection theory (SDT) as the difference between the means of two internal distributions (signal and noise distribution).

- (a) Calculate  $d'$  by using the basic SDT model:

$$d' = z(P(Hit)) - z(P(FalseAlarm))$$

with

$$P(Hit) = P(yes|signal)$$

$$P(FalseAlarm) = P(yes|noise)$$

and  $z(x)$  denoting the quantile function `qnorm()` of the standard normal distribution.

- (b) This basic SDT model may also be written as a probit model:  $\Phi^{-1}(p_{yes}) = \beta_0 + \beta_1 x$ , where  $\Phi^{-1}$  is the quantile function of the standard normal distribution. In this model,  $d' = \beta_1$ .

Estimate the  $d'$  using `glm(..., family = binomial(probit), ...)`.

- c) In a psychophysical experiment, a person is presented with two LEDs, a standard at an intensity of 40 cd/m<sup>2</sup> and a comparison at varying intensities. The task is to decide if the comparison is brighter than the standard. For each comparison intensity  $x$  the person completes 40 trials. The number of positive responses  $y$  are displayed in the following table:

x (cd/m <sup>2</sup> )	37	38	39	40	41	42	43
y (positive)	2	3	10	25	34	36	39

Estimate the parameters  $c$  and  $a$  of the logistic psychometric function

$$p_{pos} = \frac{1}{1 + \exp(-\frac{x - c}{a})}$$

using `glm()`. The model may be written as a GLM: Let  $\text{logit}(p_{pos}) = \beta_0 + \beta_1 x$ , then  $a = 1/\beta_1$  and  $c = -\beta_0/\beta_1$ . Test the goodness of fit of the model. For what intensity of  $x$  is  $p_{pos} = 0.5$  (point of subjective equality)?

Make a plot of the observed proportions of positive responses as a function of comparison intensity. Add the fitted logistic psychometric function to the plot. Hint: Use the `predict()` method to obtain the model predictions. Use `abline()` to indicate the location of the  $c$  parameter estimate.

d) Return to the power-of-the-t-test example on the slides.

- (a) Determine the Type I error rate  $\alpha$  of the test by simulation: Set the effect size to zero and show that “simulated power” is indeed 5%.

Does this simulated  $\alpha$  depend on variability ( $\sigma$ ) or sample size ( $n$ )?

- (b) For the settings on the slides, what sample size is required for the test to have a power of at least 95%?
- (c) Try with a different test: Determine the power to detect a slope of 2 in the regression model

$$y_i = 5 + 2 \cdot x_i + \varepsilon_i; \quad \varepsilon_i \sim N(0, 12^2) \text{ i.i.d.}; \quad i = 1, \dots, 15; \quad x = (1, \dots, 15)$$

e) Go back to the psychophysical data in b) and refit the logistic GLM ( $M_0$ ). Perform a parametric bootstrap of the goodness of fit test:

- (a) Simulate responses from  $M_0$  (check out `?simulate`).
- (b) Fit  $M_0$  to the simulated responses.
- (c) Fit the saturated model ( $M_1$ ) to the simulated responses.
- (d) Compute the likelihood ratio statistic.
- (e) Repeat steps (a)–(d) 1000 times.

The distribution of the simulated likelihood ratio statistics serves as a sampling distribution for the original likelihood ratio statistic computed in b). Calculate the bootstrap p-value from this distribution, which is the proportion of simulated statistics that is larger than the original.