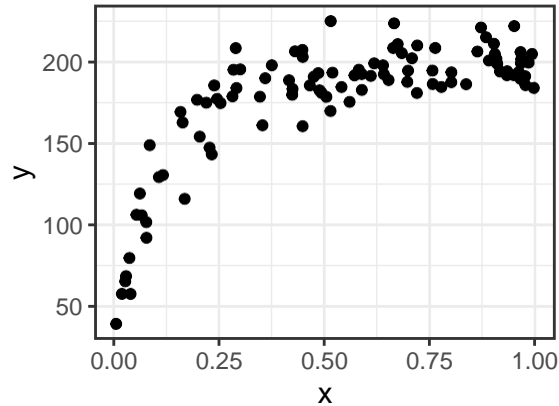


# Statistical Data Modelling With R

## Question 1

The dataframe `nlmodel` contains data on a response variable  $y$  and a single explanatory variable  $x$ . A scatter plot of  $y$  versus  $x$  suggests a strong non-linear relationship:



Suppose for these data we wish to consider the model

$$y_i \sim N \left( \frac{\vartheta_1 x_i}{\vartheta_2 + x_i}, \sigma^2 \right)$$
$$i = 1, 2, \dots, 100, \quad y_i \text{ independent}$$

- (a) [1 mark] Why can't this model be fit using a linear (regression) model?
- (b) [2 marks] Write down the likelihood  $L(\vartheta_1, \vartheta_2, \sigma^2; \mathbf{y}, \mathbf{x})$  and the log-likelihood  $l(\vartheta_1, \vartheta_2, \sigma^2; \mathbf{y}, \mathbf{x})$

- (c) [1 mark] Write an R function `mylike()` which evaluates the negative log-likelihood (i.e.  $-l(\vartheta_1, \vartheta_2, \sigma; \mathbf{y}, \mathbf{x})$ ) for any values of the three parameters
- (d) [5 marks] Use the R function `nlm()` in association with your function `mylike()` to numerically minimise the log-likelihood. Provide some evidence of how you chose sensible starting values. Report the maximum likelihood estimates of the parameters and superimpose a plot of the associated mean relationship on a scatter plot of  $y$  versus  $x$ .
- (e) [5 marks] Report the standard errors for  $\vartheta_1$  and  $\vartheta_2$ , and use those to construct 95% confidence intervals.
- (f) [3 marks] Test the hypothesis that  $\vartheta_2 = 0.08$  at the 5% significance level (not using the confidence interval) and compute the associated p-value of the test.
- (g) [4 marks] Use plug-in prediction to construct and plot 95% prediction intervals.

## Question 2

The dataframe `aids` data relates to the number of quarterly AIDS cases in the UK,  $y_i$ , from January 1983 to March 1994. The variable `cases` is  $y_i$  and `date` is time, symbolised here as  $x_i$ . In this question we consider two competing models to describe the trend in the number of cases. Model 1 is

$$y_i \sim \text{Pois}(\lambda_i) \\ \log(\lambda_i) = \beta_0 + \beta_1 x_i$$

and Model 2 is

$$y_i \sim N(\mu_i, \sigma^2) \\ \log(\mu_i) = \gamma_0 + \gamma_1 x_i$$

- (a) [3 marks] Plot  $y_i$  against  $x_i$  and comment on whether the two proposed models are sensible in terms of the distribution and the relationship of  $x$  with the mean.
- (b) [5 marks] Fit the two models in R. Plot the estimated trends from each model ( $\hat{\lambda}_i$  and  $\hat{\mu}_i$ ) on top of the data with approximate 95% confidence intervals around the mean. Comment on the validity of each model (based on the plot). Obtain the AIC for each model and thus comment on which model is preferable according to this criterion.
- (c) [2 marks] Produce the deviance residuals vs fitted values ( $\hat{\lambda}_i$  and  $\hat{\mu}_i$ ) plot for each model, comment appropriately and thus propose a way that the two models might be extended to improve the fit.
- (d) [4 marks] Implement the proposed extensions to each model, to arrive at a final version for each of them (justified by appropriate hypothesis tests).
- (e) [9 marks] On the basis of your answer to (a) the analogous plots as in (b) and (c) but also on arguments of model fit based on the deviance and the AIC, comment on which (if any) of the two final models in (d) you would choose as the best. Mention at least one reason why either model is not ideal.
- (f) [4 marks] Further extend your final Poisson model to a Negative Binomial model and comment on whether this model is preferable to the other two, on the basis of all the criteria used for comparison so far.

## Question 3

The data frame `titanic` relates to 1309 passengers on the last voyage of the ocean liner ‘Titanic’. The response variable `survived` is a binary variable where the value 1 means the passenger survived the sinking. The data frame also contains predictors relating to passenger class (1st, 2nd, 3rd), gender, age and the fare amount each passenger paid. Passenger names are also available (for interest, rather than for modelling).

[12 marks] Fit a Bernoulli GLM with logistic link of `survived` with `age`, `pclass` and `gender` are predictors, as well as all the associated two-way interactions. Reduce the model if and as appropriate using the AIC in conjunction with the R function `drop1()`, and interpret the final model in terms of parameter estimates and their significance. Perform relevant model checking analysis (model fit and residuals).