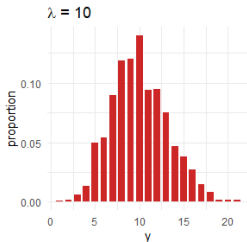
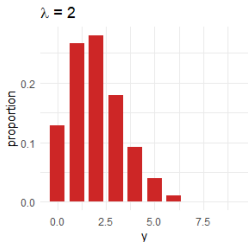
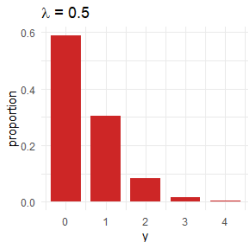


Poisson distribution

We often use the Poisson distribution to model count data. If $Y \sim Poi(\lambda)$ with $\lambda > 0$, then the PMF is given by

$$P(Y = y) = \frac{\lambda^y e^{-\lambda}}{y!}, \quad y = 0, 1, 2, \dots$$

In addition, for Poisson distributed random variables, we have that $E[Y] = \text{var}(Y) = \lambda$. Eventually, we have that $\sum_{i=1}^n y_i \sim Poi(\sum_{i=1}^n \lambda_i)$.



Poisson regression model

Consider n independent observations y_1, \dots, y_n for which we assume a Poisson distribution conditionally on a set of p categorical or numerical covariates x_j , for $j = 1, \dots, p$. The model is given by

$$\ln\left(E[y_i | x_i]\right) = \ln(\lambda_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} = \mathbf{x}_i^T \boldsymbol{\beta}$$

with $i = 1, \dots, n$, with $\mathbf{x}_i^T = (1, x_{i1}, \dots, x_{ip})^T$ and $\boldsymbol{\beta} = (\beta_0, \dots, \beta_p)$.

The **natural link function is the log link**. It ensures that $\lambda_i \geq 0$. It follows that

$$E[y_i | x_i] = \lambda_i = e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}} = e^{\mathbf{x}_i^T \boldsymbol{\beta}}$$

The Poisson GLM is suitable for modeling count data as response variable Y when a set of assumptions are met.

Parameter estimation

The log-likelihood function is given by

$$l(\mathbf{y}, \boldsymbol{\beta}) = \sum_{i=1}^n \left(y_i \mathbf{x}_i^T \boldsymbol{\beta} - e^{\mathbf{x}_i^T \boldsymbol{\beta}} - \ln(y_i!) \right)$$

Differentiating with respect to $\boldsymbol{\beta}$ and setting the new function equal to 0 yields the **Maximum Likelihood equations**

$$\sum_{i=1}^n (y_i - e^{\mathbf{x}_i^T \boldsymbol{\beta}}) x_{ij} = 0$$

with $j = 0, \dots, p$ and $x_{i0} = 1$.

There is **no closed-form solution** for the Maximum Likelihood equations. We therefore have to resort to numerical optimization, for example the Iteratively Weighted Least Squares (IWLS) algorithm or the Newton-Raphson algorithm to obtain estimates of the regression coefficients.

Model assumptions

(i) **Count response:** The response variable is a count (non-negative integers), i.e. the number of times an event occurs in an homogeneous time interval or a given space (e.g. the number of goal scored during a football game). It is suitable for grouped or ungrouped data since the sum of Poisson distributed observations is also Poisson. When the response is a category (a ranking), we should consider a Multinomial GLM instead.

(ii) **Independent events:** The counts, i.e. the events, are assumed to be independent of each other. When this assumption does not hold, we should consider a Generalized Linear Mixed Model (GLMM) instead.

(iii) **Constant variance:** The factors affecting the mean are also affecting the variance. The variance is assumed to be equal to the mean. When this assumption does not hold, we should consider a Quasipoisson GLM for overdispersed (or underdispersed) data or a Negative Binomial GLM instead.

Parameter interpretation

- (i) β_0 represents the change in the log of the mean when all covariates x_j are equal to 0. Thus e^{β_0} represents the change in the mean.
- (ii) β_j , for $j > 0$ represents the change in the log of the mean when x_j increases by one unit and all other covariates are held constant. Thus e^{β_j} represents the change in the mean.

Practical example

We will fit a Poisson regression model to a subset of the 'Affairs' dataset.
(after W. H. Greene)

There are $n = 20$ observations and 8 variables in the reduced dataset. The variable 'affairs' is the number of extramarital affairs in the past year and is our response variable. We will include as covariates the variables 'gender', 'age', 'yearsmarried', 'children', 'religiousness', 'education' and 'rating' in our analysis. 'religiousness' ranges from 1 (anti) to 5 (very) and 'rating' is a self rating of the marriage, ranging from 1 (very unhappy) to 5 (very happy).

```
data(Affairs, package = 'AER')
set.seed(1994)
data <- Affairs[sample(nrow(Affairs), size = 20, replace = FALSE), -c(8)]
head(data)
#      affairs gender age yearsmarried children religiousness education rating
# 295         0   male  32          10      yes             4          20      4
# 204         1   male  42          15      yes             4          16      5
# 1584        0 female  37          10      yes             4          16      5
# 1682        7 female  32          15      yes             5          18      4
# 1669        2 female  27           4      no              1          17      1
# 645         0 female  27          10      yes             4          16      3

dim(data)
# [1] 20  8
class(data)
# [1] "data.frame"
```

Fitted Poisson model

```
# Poisson model
poisson.model <- glm(affairs ~ .,
                     family = 'poisson', data = data)
summary(poisson.model)

# Deviance Residuals:
#   Min       1Q   Median       3Q      Max
# -2.3392  -0.7669  -0.4425   0.1047   1.8788
#
# Coefficients:
#              Estimate Std. Error z value Pr(>|z|)
# (Intercept)   0.04201    2.58877   0.016  0.98705
# gendermale    -0.32727    0.60877  -0.538  0.59085
# age           -0.04331    0.05139  -0.843  0.39929
# yearsmarried   0.22417    0.11645   1.925  0.05423
# childrenyes    0.94834    0.67143   1.412  0.15782
# religiousness  0.68438    0.45728   1.497  0.13449
# education     -0.01677    0.11092  -0.151  0.87984
# rating        -1.17513    0.43671  -2.691  0.00713 **
#
# Signif. codes:  0   ***    0.001   **    0.01   *    0.05   .    0.1    1
#
# (Dispersion parameter for poisson family taken to be 1)
#
# Null deviance: 102.924  on 19  degrees of freedom
# Residual deviance:  19.145  on 12  degrees of freedom
# AIC: 60.837
#
# Number of Fisher Scoring iterations: 6
```

Deviance and goodness-of-fit

The **deviance** of the model (also called G-statistic) is given by

$$D_{model} = 2 \sum_{i=1}^n \left(y_i \ln \left(\frac{y_i}{\hat{\lambda}_i} \right) - (y_i - \hat{\lambda}_i) \right)$$

where $\hat{\lambda}_i = e^{\mathbf{x}_i^T \hat{\beta}}$ is the fitted value of λ_i .

The deviance can be used as a goodness-of-fit test. We test H_0 : 'The model is appropriate' versus H_1 : 'The model is not appropriate'. Under H_0 , we have that

$$D_{model} \sim \chi^2_{1-\alpha, n-(p+1)}$$

where $p + 1$ is the number of parameters of the model and $1 - \alpha$ is a quantile of the χ^2 distribution.

```
# p-value of Residual deviance goodness-of-fit test  
1 - pchisq(deviance(poisson.model), df = poisson.model$df.residual)  
# [1] 0.08507918
```

Our model does not fit the data very well. Since our p-value is 0.085, H_0 is just not rejected.

Pearson goodness-of-fit

The **Pearson goodness-of-fit statistic** is given by

$$X^2 = \sum_{i=1}^n \frac{(y_i - \hat{\lambda}_i)^2}{\hat{\lambda}_i}$$

where $\hat{\lambda}_i = e^{\mathbf{x}_i^T \hat{\beta}}$ is the fitted value of λ_i .

We test H_0 : 'The model is appropriate' versus H_1 : 'The model is not appropriate'. Under H_0 , we have that

$$X^2 \sim \chi^2_{1-\alpha, n-(p+1)}$$

where $p + 1$ is the number of parameters of the model and $1 - \alpha$ is a quantile of the χ^2 distribution.

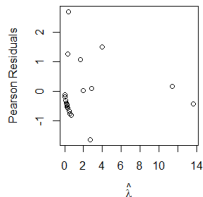
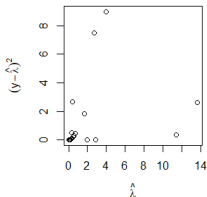
```
# Pearson's goodness-of-fit
Pearson <- sum((data$affairs - poisson.model$fitted.values)^2
               / poisson.model$fitted.values)
1 - pchisq(Pearson, df = poisson.model$df.residual)
# [1] 0.1053663
```

The fit is not much better. Our p-value is 0.1054 and H_0 is not rejected.

Checking $E[Y] = \text{var}(Y)$ assumption

The variance of y_i is approximated by $(y_i - \hat{\lambda}_i)^2$. From the first graph we can see that the range of the variance differs from the range of the mean. Moreover, from the second graph, we see that the residuals show some kind of pattern. $E[Y] = \text{var}(Y)$ seems not to hold. Let us examine the dispersion of the data and try a Quasipoisson in case of overdispersion.

```
# Checking mean = variance assumption
lambdahat <- fitted(poisson.model)
par(mfrow=c(1,2), pty="s")
plot(lambdahat, (data$affairs - lambdahat)^2,
      xlab=expression(hat(lambda)), ylab=expression((y-hat(lambda))^2))
plot(lambdahat, resid(poisson.model, type="pearson"),
      xlab=expression(hat(lambda)), ylab="Pearson Residuals")
```



Assessing overdispersion

The variance of Y must be somewhat proportional to its mean. We can write

$$\text{var}(Y) = E[Y] = \phi\lambda$$

where ϕ is a **scale parameter of dispersion** and is equal to 1 if the equality $E[Y] = \text{var}(Y)$ holds. If $\phi > 1$, the data are **overdispersed** and if $\phi < 1$, the data are underdispersed. If a Poisson model is fitted under overdispersion of the response, then the standard errors of the estimated coefficients are underestimated. The scale parameter ϕ can be estimated as

$$\hat{\phi} = \frac{\sum_{i=1}^n \frac{(y_i - \hat{\lambda}_i)^2}{\hat{\lambda}_i}}{n - (p + 1)} = \frac{X^2}{n - (p + 1)}$$

```
# Estimated dispersion parameter  
Pearson / poisson.model$df.residual  
# [1] 1.529472
```

The dispersion parameter is roughly equal to 1.53 for our data. Let us try a Quasipoisson regression model.

Fitted Quasipoisson model

The fitted Quasipoisson model yields the following R output. However, the fit seems not to have improved based on the deviance goodness-of-fit test.

```
# Quasipoisson model
quasipoisson.model <- glm(affairs ~ .,
                           family = 'quasipoisson', data = data)
summary(quasipoisson.model)

# Coefficients:
#              Estimate Std. Error t value Pr(>|t|)
# (Intercept)   0.04201    3.20159   0.013  0.9897
# gendermale    -0.32727    0.75287  -0.435  0.6715
# age          -0.04331    0.06355  -0.682  0.5085
# yearsmarried  0.22417    0.14402   1.557  0.1455
# childrenyes   0.94834    0.83037   1.142  0.2757
# religiousness 0.68438    0.56552   1.210  0.2495
# education    -0.01677    0.13718  -0.122  0.9047
# rating       -1.17513    0.54008  -2.176  0.0503
# ---
# Signif. codes:  0      ***      0.001      **      0.01      *      0.05      .      0.1      1
#
# (Dispersion parameter for quasipoisson family taken to be 1.529477)

# p-value of Residual deviance goodness-of-fit test
1 - pchisq(deviance(quasipoisson.model), df = quasipoisson.model$df.residual)
# [1] 0.08507918
```

Variable selection using BIC

Some variables may not be relevant to the model or have low explanatory power. **Stepwise model selection** provides one possible solution to select our covariates based on Akaike Information Criterion (AIC) or **Bayesian Information Criterion (BIC)** reduction (not available for Quasipoisson models).

```
# variable selection using BIC
library(MASS)
stepAIC(poisson.model, direction = 'both', k = log(dim(data)[1]))
# Step: AIC=61.42
# affairs ~ yearsmarried + children + religiousness + rating
#
#
```

	Df	Deviance	AIC
# <none>		20.753	61.423
# + age	1	19.461	63.128
# - children	1	25.501	63.176
# + gender	1	19.879	63.546
# + education	1	20.750	64.417
# - yearsmarried	1	32.187	69.862
# - religiousness	1	32.965	70.640
# - rating	1	57.142	94.817

It appears that the variables 'yearsmarried', 'children', 'religiousness' and 'rating' are the most relevant to our analysis. The next step is to select the best Quasipoisson model between one including all covariates and one for which only those four covariates are incorporated in the model.

Model selection using Crossvalidation

We will select the best model in terms of predictions using **leave-one-out Crossvalidation (LOOCV)**. The model with the lowest Root Mean Squared Error (RMSE) will be preferred.

```
# Leave-one-out crossvalidation (LOOCV)
pred.cv.mod_1 <- pred.cv.mod_2 <- numeric(dim(data)[1])
for(i in 1:dim(data)[1]) {
  mod_1 = glm(affairs ~ .,
              family = 'quasipoisson', data = data, subset = -i)
  mod_2 = glm(affairs ~ children + yearsmarried + religiousness + rating,
              family = 'quasipoisson', data = data, subset = -i)
  pred.cv.mod_1[i] = predict.glm(mod_1, data[i,], type = 'response')
  pred.cv.mod_2[i] = predict.glm(mod_2, data[i,], type = 'response')
}

error.mod_1 = (1/dim(data)[1]) * sum((data$affairs - pred.cv.mod_1)^2)
error.mod_2 = (1/dim(data)[1]) * sum((data$affairs - pred.cv.mod_2)^2)

# Root Mean Squared Error (RMSE)
sqrt(c(error.mod_1, error.mod_2))
# [1] 59196.342 337.297
```

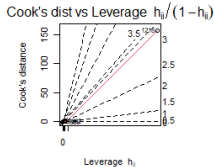
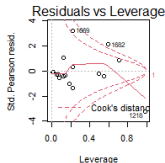
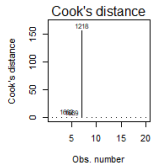
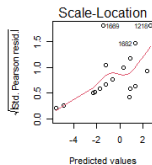
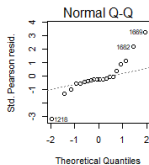
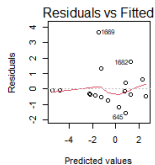
Clearly, the model with four covariates yields better predictions than the complete model and should be preferred. However, the RMSE remains relatively large indicating potential outliers in the dataset.

Diagnostic plots

```
# Diagnostic plots
```

```
par(mfrow = c(2,3))
```

```
plot(quasipoisson.model.2, which = 1:6)
```



Based on the **Cook's distance**, the observation 1218 appears to be atypical and have a strong influence on the parameter estimates as well as on the predictions. This observation should be removed.

Final model

```
# Final model
quasipoisson.model.3 = glm(affairs ~ children + yearsmarried + religiousness + rating,
                           family = 'quasipoisson', data = data2, maxit = 100)
summary(quasipoisson.model.3)

# Coefficients:
#               Estimate Std. Error t value Pr(>|t|)
# (Intercept)  -0.14080    0.70522  -0.200  0.844619
# childrenyes   -2.74742    1.10230  -2.492  0.025841 *
# yearsmarried  0.30447    0.07101   4.288  0.000751 ***
# religiousness 1.64490    0.39165   4.200  0.000891 ***
# rating       -2.03423    0.39565  -5.141  0.000150 ***
# ---
# Signif. codes:  0      ***      0.001      **      0.01      *      0.05      .      0.1      1

# p-value of Residual deviance goodness-of-fit test
1 - pchisq(deviance(quasipoisson.model.3), df = quasipoisson.model.3$df.residual)
# [1] 0.667648

# Pearson's goodness-of-fit
Pearson <- sum((data2$affairs - quasipoisson.model.3$fitted.values)^2
               / quasipoisson.model.3$fitted.values)
1 - pchisq(Pearson, df = quasipoisson.model.3$df.residual)
# [1] 0.7263845
```

Once the outlier has been removed, the fit is much better and the standard errors are much lower compared to the parameter estimates. This is our best model.

Conclusions

- (i) The problems of overdispersion, covariate selection and influence of outliers have been addressed. Our final Quasipoisson model is a good fit for the data. About 86% of the deviance is explained by the model.
- (ii) The level of religiousness and the number of years of marriage seem to be positively related to the average number of affairs, whereas having children and a happy self rated marriage seem to be negatively related to the average number of affairs. Caution however since the dataset only contains 19 observations.
- (iii) If an individual has one child or more, the change in the mean response given all other covariates held constant is $e^{-2.75} \approx 0.064$, hence a decrease of 93.6% of the average number of affairs in the past year.
- (iv) For one more year of marriage, the change in the mean response given all other covariates held constant is $e^{0.304} \approx 1.36$, hence an increase of 36% of the average number of affairs in the past year.
- (v) When the self rating of the marriage changes from unhappy to happy, the change in the mean response given all other covariates held constant is $e^{-2.034} \approx 0.13$, hence a decrease of 87% of the average number of affairs in the past year.