

Generalized Linear Models (GLM)

- ★ Regression is the study of the change of the distribution of one variable, y , according to the value of another variable, x .
- ★ Both continuous and discrete variables can be included in regression problems.
- ★ Formally, a relationship between x and y usually means that the expected value of y is different for different values of x .
- ★ At this stage, changes in σ^2 or other aspects of the distribution are not considered.
- ★ When y is continuous, changes in y are smooth (continuous and differentiable), and the model used is:

$$E(y|x) = g(x)$$

for some unknown smooth function $g(\cdot)$.

- ★ A variety of functions can be used to estimate $g(\cdot)$. These functions are called **scatterplot smoother**.

- ★ For a pair (x_i, y_i) , $i = 1, \dots, n$:

$$y_i = g(x_i) + \varepsilon_i$$

where $g(\cdot)$ is called the **regression function**.

- ★ A chemical theory suggests that the temperature at which a certain reaction occurs is 180° .
- ★ The results of 10 independent experiments, contained in `reaction.dat`, are:
179.4, 177.5, 180.6, 178.1, 180.2, 179.3, 180.0, 178.3, 179.1, 181.1.
- ★ In this example, $n = 10$ and

$$E(y_i) = \beta_0 \text{ for } i = 1, \dots, n$$

- ★ More generally, however, different observations are associated with different experimental conditions or different values of one or more **explanatory variables** as the following examples will illustrate.

- ★ To compare two growth hormones, 19 chicken were tested, 11 being assigned to hormone A and 8 to hormone B.
- ★ The following were the gains in weight (grams) over the period of experiment, contained in `hormones.dat`:

Hormone A: 615, 645, 840, 615, 365, 450, 450, 530, 756, 785, 560.

Hormone B: 870, 895, 1035, 700, 560, 756, 785, 980.

- ★ $E(y_{ij}) = \beta_i$ for $j = 1, \dots, n_i$; $i = 1, 2$.
- ★ Thus, β_1 and β_2 are the expectations of y for hormones A and B, respectively.
- ★ Typically, it is interesting in whether or not $\beta_1 = \beta_2$ or in the range of values for the difference $\beta_1 - \beta_2$.

- ★ In an experiment to investigate the effects of a depressant drug, the reaction times of then male rats to a certain stimulus were measured after a specified dose of the drug have been administered to each rat.
- ★ The results, contained in `drug.dat`, were as follows:

Rat number:	1	2	3	4	5	6	7	8	9	10
Dose (mg):	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
Reaction time (secs):	0.32	0.24	0.40	0.52	0.44	0.56	0.64	0.52	0.60	0.80

- ★ A simple model for this data might be a straight line, i.e.

$$E(y_i|x_i) = \beta_0 + \beta_1 x_i \text{ for } i = 1, \dots, n$$

- ★ For simplicity, usually the above formula is written as $E(y_i) = \beta_0 + \beta_1 x_i$, the conditioning on x_i being understood.

- ★ In the study of the polyesterification of fatty acids with glycols, the effect of temperature ($^{\circ}\text{C}$) on the percentage conversion of the esterification process was investigated.
- ★ The following data, contained in `polyesterification.dat`, are the results of an experiment using a catalyst of $4 \cdot 10^{-4}$ mole zinc chloride per 100 grams of fatty acid.

Temperature ($^{\circ}\text{C}$):	175	200	225	250	275	300
Percentage conversion (%):	67	83	93	97	99	92

- ★ This is similar to the previous example, but now a quadratic function of x seems sensible:

$$E(y_i) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 \text{ for } i = 1, \dots, n$$

- ★ The data in the `oxidant.dat` file give levels of an air pollutant, “Oxidant”, together with levels of four meteorological variables recorded on 30 days during one summer.
- ★ Which, if any, of the four independent variables seem to be related to levels of oxidant?
- ★ Again, a simple starting point would be:

$$E(y_i) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 \text{ for } i = 1, \dots, n$$

Model formulae

Expression	Description
$y \sim x$	Simple regression
$y \sim 1 + x$	Explicit intercept
$y \sim -1 + x$	Through the origin
$y \sim x + x^2$	Quadratic regression
$y \sim x1 + x2 + x3$	Multiple regression
$y \sim G + x1 + x2$	Parallel regressions
$y \sim G / (x1 + x2)$	Separate regressions
$\text{sqrt}(y) \sim x + x^2$	Transformed
$y \sim G$	Single classification
$y \sim A + B$	Randomized block
$y \sim B + N * P$	Factorial in blocks
$y \sim x + B + N * P$	with covariate
$y \sim . - x1$	All variables except x1
$y \sim . + A:B$	Add interaction (update)
<code>Nitrogen ~ Times*(River/Site)</code>	More complex design

Common R functions for inference

Expression	Description
<code>coef(obj)</code>	regression coefficients
<code>resid(obj)</code>	residuals
<code>fitted(obj)</code>	fitted values
<code>summary(obj)</code>	analysis summary
<code>predict(obj, newdata = ndat)</code>	predict for newdata
<code>deviance(obj)</code>	residual sum of squares

- ★ All the examples discussed above can be expressed in the form:

$$E(y_i) = \beta_0 + \sum_{j=1}^p \beta_j x_{i,j}$$

- ★ The above formula can be rewritten as

$$\underline{\mu} = E(\underline{y}) = \underline{X}\underline{\beta}$$

where \underline{X} is a $n \times (p + 1)$ matrix, with all ones in the first columns, and $\underline{\beta}$ is the vector of $\beta_0, \beta_1, \dots, \beta_p$.

★ In the example 1: $\mathbf{X} = (\underline{\mathbf{1}})$ and $\underline{\boldsymbol{\beta}} = (\beta_0)$

★ In the example 2:

$$\mathbf{X} = \begin{pmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{pmatrix}_{n \times 2} \quad \underline{\boldsymbol{\beta}} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}_{2 \times 1}$$

★ In the example 3:

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}_{n \times 2} \quad \underline{\boldsymbol{\beta}} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}_{2 \times 1}$$

★ In the example 4:

$$\mathbf{X} = \begin{pmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix}_{n \times 3} \quad \underline{\boldsymbol{\beta}} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}_{3 \times 1}$$

★ In the example 5:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & x_{1,2} & x_{1,3} & x_{1,4} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n,1} & x_{n,2} & x_{n,3} & x_{n,4} \end{pmatrix}_{n \times 5} \quad \underline{\boldsymbol{\beta}} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix}_{3 \times 1}$$

★ Models of the form $\mu = E(y_i) = \mathbf{X}\boldsymbol{\beta}$ are called **linear models**, because they are linear in the **unknown parameters** comprising the elements of $\boldsymbol{\beta}$.

- ★ The file `hotdogs.dat` contains calories and sodium data for three types of hot dogs: Poultry, Beef, Meat.
- ★ 20 Beef, 17 Meat and 17 Poultry hot dogs. These data may be thought as samples from much larger populations.
- ★ Looking at the density estimates it seems reasonable, at least as a first approximation, to model the distribution as Normal.
- ★ The following model can be adopted (model 1):

$$B_1, \dots, B_{20} \sim \text{i.i.d. } N(\mu_B, \sigma)$$

$$M_1, \dots, M_{17} \sim \text{i.i.d. } N(\mu_M, \sigma)$$

$$P_1, \dots, P_{17} \sim \text{i.i.d. } N(\mu_P, \sigma)$$

- ★ An equivalent formulation (model 2):

$$B_1, \dots, B_{20} \sim \text{i.i.d. } N(\mu, \sigma)$$

$$M_1, \dots, M_{17} \sim \text{i.i.d. } N(\mu + \delta_M, \sigma)$$

$$P_1, \dots, P_{17} \sim \text{i.i.d. } N(\mu + \delta_P, \sigma)$$

- ★ There are three parameters for the population means and one for the standard deviation.

Model 1	Model 2	Interpretation	Approx value
μ_B	μ	MCC of Beef hot dogs	157
μ_M	$\mu + \delta_M$	MCC of Meat hot dogs	159
μ_P	$\mu + \delta_P$	MCC of Poultry hot dogs	119
$\mu_M - \mu_B$	δ_M		2
$\mu_P - \mu_B$	δ_P		-38
σ	σ	Standard deviation	29

MCC: Mean calorie content

- ★ Rewriting the data matrix in a slightly different form reveals some mathematical structure common to many models.
- ★ Let (y_1, \cdot, y_{54}) be the calorie contents. Two new variables are defined:

$$x_{1,i} = \begin{cases} 1 & \text{if the } i\text{-th hot dog is Meat} \\ 0 & \text{otherwise} \end{cases}$$

$$x_{2,i} = \begin{cases} 1 & \text{if the } i\text{-th hot dog is Poultry} \\ 0 & \text{otherwise} \end{cases}$$

- ★ $x_{1,i}$ and $x_{2,i}$ are indicator or dummy variables.
- ★ If there are k populations, then $k - 1$ indicators suffice.

★ Model 2 can be stated as

$$y_i = \mu + \delta_M x_{1,i} + \delta_P x_{2,i} + \varepsilon_i$$

$$\varepsilon_1, \cdot, \varepsilon_{54} \sim \text{i.i.d. } N(0, \sigma)$$

★ Let

$$\begin{aligned} \underline{y} &= (y_1, \cdot, y_{54})^T \\ \underline{\beta} &= (\mu, \delta_M, \delta_P)^T \\ \underline{\varepsilon} &= (\varepsilon_1, \cdot, \varepsilon_{54})^T \end{aligned} \quad \mathbf{X} = \begin{pmatrix} 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{pmatrix}_{54 \times 3}$$

- ★ The general form of the linear model:

$$\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$$

or

$$y_i = \beta_0 + \beta_1 \cdot x_{i1} + \beta_2 \cdot x_{i2} + \cdots + \beta_p \cdot x_{ip} + \varepsilon_i$$

or

$$y_i \sim N(\mu_i, \sigma)$$

with $\mu_i = \beta_0 + \sum_j \beta_j x_{j,i}$ and $\varepsilon_i \sim \text{i.i.d. } N(0, \sigma)$.

- ★ This model has $p + 2$ parameters: $\beta_0, \beta_1, \dots, \beta_p, \sigma$.

- ★ $X\beta = \underline{\mu}$: **systematic** part, **deterministic** part or **signal**.
- ★ $\underline{\varepsilon}$: **random** or **noise** component.
- ★ response = signal + noise
- ★ Models in which the signal is a linear function of parameters (μ , δ_M , δ_P in the example) and response is a linear function of signal and noise, are linear models.

★ The likelihood function is

$$\begin{aligned}l(\beta_0, \dots, \beta_p, \sigma) &= \prod_{i=1}^n p(y_i | \beta_0, \dots, \beta_p, \sigma) = \\&= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{y_i - \mu_i}{\sigma}\right)^2\right) = \\&= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{y_i - (\beta_0 + \sum_j \beta_j x_{j,i})}{\sigma}\right)^2\right) = \\&= (2\pi\sigma^2)^{-\frac{n}{2}} \exp -\frac{1}{2\sigma^2} \sum_i (y_i - (\beta_0 + \sum_j \beta_j x_{j,i}))^2\end{aligned}$$

- ★ The MLE (maximum likelihood estimate) is obtained by take the log of the likelihood expression, differentiate with respect to each parameter in turn, set the derivative equal to 0 and solve.

$$\log l(\beta_0, \dots, \beta_p, \sigma) = c - n \log \sigma - \frac{1}{2\sigma^2} \sum_i (y_i - (\beta_0 + \sum_j \beta_j x_{j,i}))^2$$

$$\begin{cases} \frac{1}{\sigma^2} \sum_i (y_i - (\hat{\beta}_0 + \sum_j \hat{\beta}_j x_{j,i})) = 0 \\ \frac{1}{\sigma^2} \sum_i (y_i - (\hat{\beta}_0 + \sum_j \hat{\beta}_j x_{j,i})) x_{1,i} = 0 \\ \vdots \\ \frac{1}{\sigma^2} \sum_i (y_i - (\hat{\beta}_0 + \sum_j \hat{\beta}_j x_{j,i})) x_{p,i} = 0 \\ -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_i (y_i - (\hat{\beta}_0 + \sum_j \hat{\beta}_j x_{j,i}))^2 = 0 \end{cases}$$

- ★ The first $p + 1$ equations can be multiplied by σ^2 , yielding $p + 1$ linear equations in the $p + 1$ unknown β s.
- ★ Because they are linear, equations can be solved by linear algebra:

$$\underline{\hat{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{y}$$

from which

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{1,i} + \cdots + \hat{\beta}_p x_{p,i}$$

- ★ \hat{y}_i are the **fitted values**, while $r_i = y_i - \hat{y}_i$ are **residuals**, estimates of ε_i .
- ★ The MLE for σ , $\hat{\sigma}$, is found from the last equation that can be rewritten as $-\frac{n}{\hat{\sigma}} + \frac{1}{\hat{\sigma}^3} \sum_i r_i^2 = 0$, from which $\hat{\sigma}^2 = \frac{1}{n} \sum_i r_i^2$.

- ★ The main assumption of linear model is that the model is linear in the parameters.
- ★ Other assumptions concern the error terms $\underline{\varepsilon}$:

① $E(\varepsilon_i) = 0$ $(i = 1, \dots, n)$

② \mathbf{X} is not stochastic

③ $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$ $(i = 1, \dots, n)$

④ $Corr(\varepsilon_i, \varepsilon_{i-j}) = 0$ $(i = 1, \dots, n; 0 < j < i)$

- ★ The residual normality can be checked through graphics and the normality tests, such as **Anderson-Darling test**.
- ★ The Anderson-Darling test compute square mean difference between the empirical cumulative distribution (computed on residuals) and theoretical cumulative distribution of the normal distribution.

$$AD = \sum_{i=1}^n \frac{1-2i}{n} \{ \ln(F_0[Z_{(i)}]) + \ln(1 - F_0[Z_{(n+1-i)}]) \} - n$$

when n is the sample size, $Z_{(i)}$ are observed ordered and standardized residuals and F_0 is the cumulative distribution function of the standard normal distribution.

- ★ The value of AD increase when the difference between residuals distribution and gaussian cumulative distribution also increase.

- ★ If linear model assumptions are satisfied, $\hat{\beta}$ are random variables with a known distribution.

- ★ It is possible:
 - ① to test hypotheses of the form $H_0 : \beta_j = \beta_j^*$ that is to test that β_j (the true, and unknown value of the estimator) is equal to some (hypothesized) value β_j^* ;
 - ② to provide confidence intervals for the regression parameters β_j .

- ★ To test β_j is equal to some value β_j^* , $j = 0, 1, \dots, p$ the t -test statistic can be used:

$$t = \frac{\hat{\beta}_j - \beta_j^*}{\hat{se}(\beta_j)}$$

where $\hat{se}(\beta_j)$ denotes the estimated standard error of the estimator $\hat{\beta}_j$.

- ★ t test statistics follow a **Student's t** distribution with $n - p - 1$ degrees of freedom, where $p + 1$ is the number of estimated parameters.
- ★ The $1 - \alpha$ confidence interval for β_j is given by:

$$[\hat{\beta}_j - t_{n-p-1; \frac{\alpha}{2}} \hat{se}(\beta_j); \hat{\beta}_j + t_{n-p-1; \frac{\alpha}{2}} \hat{se}(\beta_j)]$$

- ★ Linear regression can be used to fit a predictive model to an observed data set of \underline{y} and \mathbf{X} values. After developing such a model, if an additional values of x , \underline{x}_{n+1} , is then given without its accompanying value of y , the fitted model can be used to make a prediction of the value of y_{n+1} .
- ★ It can be shown that the $1 - \alpha$ confidence level for the regression for a given value of \underline{x}_* is limited by:

$$\text{LCL: } f(\underline{x}_*; \hat{\underline{\beta}}) - t_{n-p-1; \frac{\alpha}{2}} \cdot \sqrt{\underline{x}_*^T (\mathbf{X}^T \mathbf{X})^{-1} \underline{x}_*} \cdot \hat{\sigma}_\varepsilon$$

$$\text{UCL: } f(\underline{x}_*; \hat{\underline{\beta}}) + t_{n-p-1; \frac{\alpha}{2}} \cdot \sqrt{\underline{x}_*^T (\mathbf{X}^T \mathbf{X})^{-1} \underline{x}_*} \cdot \hat{\sigma}_\varepsilon$$

- ★ The confidence interval may be shown as the interval containing the “true” regression line with a given confidence level.
- ★ Instead, the prediction interval will contain, with a given confidence level, the true (and unknown) y_{n+1} value, given \underline{x}_{n+1} values.
- ★ The $1 - \alpha$ confidence level for the prediction interval for a given value of \underline{x}_{n-1} is limited by:

$$\text{LPL: } f(\underline{x}_*; \hat{\underline{\beta}}) - t_{n-p-1; \frac{\alpha}{2}} \cdot \sqrt{\underline{x}_*^T (\mathbf{X}^T \mathbf{X})^{-1} \underline{x}_* + 1} \cdot \hat{\sigma}_\varepsilon$$

$$\text{UPL: } f(\underline{x}_*; \hat{\underline{\beta}}) + t_{n-p-1; \frac{\alpha}{2}} \cdot \sqrt{\underline{x}_*^T (\mathbf{X}^T \mathbf{X})^{-1} \underline{x}_* + 1} \cdot \hat{\sigma}_\varepsilon$$

- ★ Student's t test for independent samples is equivalent to the linear regression of the response variable on the grouping variable, where the grouping variable is recoded to have numerical values, if necessary.
- ★ The t test for independent samples p-value is equal to the slope regression parameter (β_1) p-value.

- ★ The linear regression model says data are normally distributed about the regression line with constant standard deviation. The grouping variable takes on only two values. Therefore, there are only two locations along the regression line where there are data.
“Homoschedastic normally distributed values about the regression line” is equivalent to “two normally distributed populations with equal variances”.
- ★ Thus, the hypothesis of equal means ($H_0 : \mu_A - \mu_B = 0$) is equivalent to the hypothesis that the regression coefficient of x is 0 ($H_0 : \beta_1 = 0$). That is, the population means are equal if and only if the regression line is horizontal.

- ★ More than two samples can be compared using ANOVA. In the ANOVA, the categorical variable is effect coded, which means that each category's mean is compared to the grand mean.
- ★ In the regression, the categorical variable is dummy coded, which means that each category's intercept is compared to the reference group's intercept.
- ★ So an ANOVA reports each mean and a p-value that says at least two are significantly different. A regression reports only one mean (as an intercept), and the differences between that one and all other means, but the p-values evaluate those specific comparisons.
- ★ Thus, the hypothesis of equal means ($H_0 : \mu_1 = \mu_2 = \dots = \mu_p$) of the ANOVA is equivalent to the hypothesis that the regression dummy coefficients $\beta_1, \beta_2, \dots, \beta_p$ are equal to zeros.

- ★ When distributional assumption about \underline{e} is too restrictive and/or some of the independent variables are regarded as random variables:

$$\underline{y} = \underline{X}\underline{\beta} + \underline{Z}\underline{u} + \underline{e}$$

where:

- \underline{Z} : random-effects design matrix; can contain both continuous and dummy matrix, just like \underline{X} ;
- \underline{u} : random-effects parameter.

- ★ \underline{u} and \underline{e} are normally distributed with

$$E \begin{pmatrix} \underline{u} \\ \underline{e} \end{pmatrix} = \begin{pmatrix} \underline{0} \\ \underline{0} \end{pmatrix} \qquad Var \begin{pmatrix} \underline{u} \\ \underline{e} \end{pmatrix} = \begin{pmatrix} \mathbf{G} & \underline{0} \\ \underline{0} & \mathbf{R} \end{pmatrix}$$

- ★ $Var(\underline{y}) = V = \mathbf{ZGZ}^T + \mathbf{R}$
- ★ V can be modeled by defining \mathbf{Z} and specifying covariance structures for \mathbf{G} and \mathbf{R} .

- ★ **Generalized least-squares** (GLS) requires knowledge of V and, therefore, of R and G .
- ★ Lacking of such information, some reasonable estimate of V can be used to minimize

$$(\underline{y} - X\underline{\beta})^T \hat{V}^{-1} (\underline{y} - X\underline{\beta})$$

- ★ A better approach entails using **likelihood-based** method, using the assumption that \underline{u} and \underline{e} are normally distributed.
- ★ Two popular methods are ML and REML (Restricted/residual ML).

★ The log-likelihood functions are:

$$\text{ML: } l(\mathbf{G}, \mathbf{R}) = -\frac{1}{2} \log |\mathbf{V}| - \frac{n}{2} \log \underline{\mathbf{r}}^T \mathbf{V}^{-1} \underline{\mathbf{r}} - \frac{n}{2} \left(1 + \log \left(\frac{2\pi}{n} \right) \right)$$

$$\begin{aligned} \text{REML: } l_R(\mathbf{G}, \mathbf{R}) = & -\frac{1}{2} \log |\mathbf{V}| - \frac{1}{2} \log |\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X}| + \\ & -\frac{n-p}{2} \log(\underline{\mathbf{r}}^T \mathbf{V}^{-1} \underline{\mathbf{r}}) - \frac{n-p}{2} \left(1 + \log \left(\frac{2\pi}{n-p} \right) \right) \end{aligned}$$

where

- $\underline{\mathbf{r}} = \underline{\mathbf{y}} - \mathbf{X}(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \underline{\mathbf{y}}$
- $p = \text{rank}(\mathbf{X})$

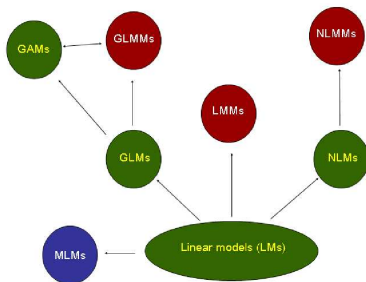
- ★ Obtained \hat{G} and \hat{R} , the mixed model equations can be solved:

$$E \begin{pmatrix} \mathbf{X}^T \hat{R}^{-1} \mathbf{X} & \mathbf{X}^T \hat{R}^{-1} \mathbf{Z} \\ \mathbf{Z}^T \hat{R}^{-1} \mathbf{X} & \mathbf{Z}^T \hat{R}^{-1} \mathbf{Z} + \hat{G}^{-1} \end{pmatrix} \begin{pmatrix} \underline{\hat{\beta}} \\ \underline{\hat{u}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^T \hat{R}^{-1} \underline{y} \\ \mathbf{Z}^T \hat{R}^{-1} \underline{y} \end{pmatrix}$$

- ★ The solutions are:

$$\underline{\hat{\beta}} = (\mathbf{X}^T \hat{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{V}^{-1} \underline{y}$$

$$\underline{\hat{u}} = \hat{G} \mathbf{Z}^T \hat{V}^{-1} (\underline{y} - \mathbf{X} \underline{\hat{\beta}})$$



Model	Acronym	R function
Linear Models	LM	lm, aov
Multivariate LMs	MLM	manova
Generalized LMs	GLM	glm
Linear Mixed Models	LMM	lme, aov
Non-linear Models	NLM	nls
Non-linear Mixed Models	NLMM	nlme
Generalized LMMs	GLMM	glmmPQL
Generalized Additive Ms	GAM	gam

- ★ Collett (1991, p. 75) reports an experiment on the toxicity to the tobacco budworm *Heliothis virescens* of doses of the pyrethroid trans-cypermethrin to which the moths were beginning to show resistance.
- ★ Batches of 20 moths of each sex were exposed for three days to the pyrethroid and the number in each batch that were dead or knocked down was recorded.
- ★ The results were

	Dose (μg)					
Sex	1	2	4	8	16	32
Male	1	4	9	13	18	20
Female	0	2	6	10	12	16

- ★ An experiment recorded the numbers of chromosomal abnormalities observed for various amounts and intensities of gamma radiation.
- ★ The number of cells per measurement varied.
- ★ The results contains, for each one of 27 measurements:
 - the number of cells per measurement;
 - the numbers of chromosomal abnormalities;
 - the amount of gamma radiation;
 - the intensity of gamma radiation.

- ★ In the case of binomial data, response variable may be specified in three different ways:
 - ❶ If the response is a numeric vector it is assumed to hold the data in ratio form, $y_i = s_i/a_i$, in which case the a_i s must be given as a vector of weights using the `weights` argument. If the a_i are all one, the default `weights` suffices.
 - ❷ If the response is a logical vector or a two-level factor it is treated as a 0/1 numeric vector and handled as previously.
 - ❸ If the response is a two-column matrix it is assumed that the first column holds the number of successes, s_i , and the second holds the number of failures, $a_i - s_i$, for each trial. In this case, no `weights` argument is required.

- ★ Generalized Linear Models (**GLM**) are extensions of fixed-effects linear model to cases where standard linear model assumptions are violated.
- ★ Standard linear model:

$$y_i = \beta_0 + \beta_1 \cdot x_{i1} + \beta_2 \cdot x_{i2} + \cdots + \beta_p \cdot x_{ip} + \varepsilon_i$$

The above formula can be written in a matrix format:

$$\underline{\mathbf{y}} = \mathbf{X}\underline{\boldsymbol{\beta}} + \underline{\boldsymbol{\varepsilon}}$$

or, in terms of expected values:

$$E(\underline{\mathbf{y}}) = \mathbf{X}\underline{\boldsymbol{\beta}}$$

Estimation is done by least squares, based on the assumption of normal errors.

- ★ GLM uses a likelihood-based procedure to fit $\mathbf{X}\underline{\beta}$ to a function of $E(\underline{\mathbf{y}})$, suggested by the distribution of data.
- ★ $\underline{\mathbf{y}}$ is described as an additive function of **systematic** and **random** components. The first corresponds to $E(\underline{\mathbf{y}})$, the later to error.

$$\underline{\mathbf{y}} = \underline{\boldsymbol{\mu}} + \underline{\boldsymbol{\varepsilon}}$$

- $Var(\underline{\mathbf{y}}) = Var(\underline{\boldsymbol{\varepsilon}}) = \mathbf{R}$.
- $\underline{\boldsymbol{\eta}} = \mathbf{X}\underline{\boldsymbol{\beta}}$, where $\underline{\boldsymbol{\eta}} = g(\underline{\boldsymbol{\mu}})$.
- $g(\underline{\boldsymbol{\mu}})$ is called the **link function**, because it links the linear model to the mean of $\underline{\mathbf{y}}$.

- ★ Nelder and Wedderburn (1972) showed that the Maximum Likelihood Estimates (MLEs) for $\underline{\beta}$ can be obtained iteratively solving:

$$\mathbf{X}^T \mathbf{W} \mathbf{X} \underline{\beta} = \mathbf{X}^T \mathbf{W} \mathbf{y}^*$$

with

- $\mathbf{W} = \mathbf{D} \mathbf{R}^{-1} \mathbf{D}$
 - $\mathbf{y}^* = \hat{\underline{\eta}} + (\underline{\mathbf{y}} - \underline{\tilde{\mu}}) \mathbf{D}^{-1}$
 - $\mathbf{D} = \partial \underline{\mu} / \partial \underline{\eta}$
 - $\mathbf{R} = \text{Var}(\underline{\varepsilon})$
 - $\underline{\mu} = E(\underline{\mathbf{y}})$
- ★ Estimates of \mathbf{D} and \mathbf{R} are used in place of \mathbf{D} and \mathbf{R} .
 - ★ In the case of standard linear model, $\underline{\eta} = E(\underline{\mathbf{y}}) = \underline{\mu}$ and $\mathbf{D} = \mathbf{I}$.
Thus, $\mathbf{X}^T \mathbf{R}^{-1} \mathbf{X} \underline{\beta} = \mathbf{X}^T \mathbf{R}^{-1} \mathbf{y}$, the normal equation.

- ★ Elements for estimating $\underline{\beta}$:
 - **link function**: determines $\underline{\nu}$ and D ;
 - **probability function**: determines $\underline{\mu}$ and R .
- ★ The process for selecting a link function and the structure of the mean and variance can be understood by looking at the probability distribution, or, better said, to the likelihood function.
- ★ Three cases will be considered: the Binomial, the Poisson and the Normal.

- ★ n trials, each trial with success probability of π :

$$f(y_n) = \binom{n}{y_n} \pi^{y_n} (1 - \pi)^{n - y_n}$$

- ★ The log-likelihood function is:

$$l(\pi; y_n) = y_n \log \left(\frac{\pi}{1 - \pi} \right) + n \log(1 - \pi) + \log \binom{n}{y_n}$$

- ★ A sample proportion, $y = y_n/n$, thus has a log-likelihood function of

$$l(\pi; y_n) = ny_n \log \left(\frac{\pi}{1 - \pi} \right) + n \log(1 - \pi) + \log \binom{n}{ny_n}$$

- ★ Mean and variance of a sample proportion are:

$$E(y) = \pi \qquad \text{Var}(y) = \frac{\pi(1 - \pi)}{n}$$

- ★ The probability function of the Poisson distribution is:

$$f(y) = \frac{\lambda^y e^{-\lambda}}{y!}$$

- ★ The log-likelihood function is:

$$l(\lambda; y) = y \log(\lambda) - \lambda - \log(y!)$$

- ★ Mean and variance are:

$$E(y) = \lambda$$

$$Var(y) = \lambda$$

- ★ The probability function of the Normal distribution is:

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y - \mu)^2\right)$$

- ★ The log-likelihood function is:

$$l(\mu, \sigma^2; y) = \left(-\frac{1}{2\sigma^2}\right) (y - \mu)^2 - \frac{1}{2} \log(2\pi\sigma^2)$$

- ★ Mean and variance are:

$$E(y) = \mu$$

$$Var(y) = \sigma^2$$

- ★ The log-likelihood functions for these three distributions have a common form,

$$l(\theta, \phi; y) = \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)$$

where

- θ is the natural parameter,
 - ϕ is the scale parameter.
- ★ θ is a function of the mean, $\theta(\mu)$.
 - ★ It is also possible express the variance as a function of the mean and $a(\phi)$:

$$Var(y) = V(\mu)a(\phi)$$

where $V(\mu)$ is the **variance function**.

	Sample proportion	Poisson	Normal
$E(y)$	π	λ	μ
$\theta(\mu)$	$\log(\pi/(1 - \pi))$	$\log(\lambda)$	μ
$a(\phi)$	$1/n$	1	σ^2
$V(\mu)$	$\pi(1 - \pi)$	λ	1
$Var(y)$	$\pi(1 - \pi)/n$	λ	σ^2

- ★ A distribution whose log-likelihood has the general form

$$l(\theta, \phi; y) = \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)$$

is a member of the **exponential family**.

- ★ The GLM can be applied to data distributed according to the exponential family.
- ★ If y_1, \dots, y_n is a random sample from such a family, the log-likelihood of y_i is

$$l(\theta_i, \phi_i; y_i) = \frac{y_i\theta_i - b(\theta_i)}{a(\phi_i)} + c(y_i, \phi_i)$$

- ★ The joint log-likelihood of y_1, \dots, y_n is

$$l(\underline{\theta}, \underline{\phi}; y_1, \dots, y_n) = \sum_i \frac{y_i\theta_i - b(\theta_i)}{a(\phi_i)} + c(y_i, \phi_i)$$

- ★ Observations are linear in the **natural parameter** θ .
 - For normally distributed data: $\theta = \mu$.
 - For Poisson distributed data: $\theta = \log(\lambda)$.
 - For Binomial distributed data: $\theta = \log(\pi/(1 - \pi))$.

- ★ In these examples, $\theta(\mu)$ is used as a link function, they are called the **canonical link functions**.

- ★ The structure of the variance-covariance matrix of \underline{y} , or errors, can be described in terms of the scale parameter and variance function. Specifically,

$$Var(\underline{\varepsilon}) = \mathbf{R} = \mathbf{R}_{\mu}^{1/2} \mathbf{A} \mathbf{R}_{\mu}^{1/2}$$

where

- \mathbf{R}_{μ} is the diagonal matrix whose i -th diagonal element is $V(\mu_i)$, the variance function for the i -th observation;
- $\mathbf{R}_{\mu}^{1/2}$ is the diagonal matrix of square roots of the corresponding elements of \mathbf{R}_{μ} ;
- \mathbf{A} is the scale parameter matrix.

Distribution	R_μ	A
Normal	I	I_{σ^2}
Poisson	$\text{diag}(\lambda_i)$	I
Sample proportion	$\text{diag}(\pi_i/(1 - \pi_i))$	$\text{diag}(1/n_i)$

- ★ The **inverse link function** (sometimes referred to as the **mean function**) is defined as $h(\underline{\eta}) = \underline{\mu}$.
- ★ It can be used to predict $\underline{\mu}$ from $\underline{\hat{\beta}}$.
- ★ It is assumed that the relationship between $\underline{\eta}$ and $\underline{\mu}$ is one-to-one, thus $h(\underline{\eta}) = g^{-1}(\underline{\eta})$. For some complex GLMs, this may be false.
- ★ Because $\underline{\eta}$ is estimated by $\mathbf{X}\underline{\hat{\beta}}$, $\hat{\mu} = h(\mathbf{X}\underline{\hat{\beta}})$.

- Normal:

$$\eta = \mu \Rightarrow \mu = h(\mathbf{X}\underline{\hat{\beta}}) = \mathbf{X}\underline{\hat{\beta}}$$

- Poisson:

$$\eta = \log \lambda \Rightarrow \lambda = h(\mathbf{X}\underline{\hat{\beta}}) = \exp \mathbf{X}\underline{\hat{\beta}}$$

- Sample proportion:

$$\eta = \log (\pi / (1 - \pi)) \Rightarrow \lambda = h(\mathbf{X}\underline{\hat{\beta}}) = \exp \mathbf{X}\underline{\hat{\beta}} / (1 + \exp \mathbf{X}\underline{\hat{\beta}})$$

- ★ Deviance is defined as:

$$Dev(\hat{\underline{\mu}}; \underline{y}) = 2 (l(\underline{y}; \underline{y}) - l(\hat{\underline{\mu}}; \underline{y}))$$

where

- $l(y; y)$ is the value of the maximum log-likelihood achievable in a full model; i.e. it is the value of the log-likelihood for which θ is expressed as a function of y ;
 - $l(\hat{\mu}; y)$ is the value of the log-likelihood over β , i.e. it is the value of the log-likelihood for which θ is expressed as a function of $\hat{\mu}$.
- ★ The deviance is a generalization of the Sum of Squares for Error in the Analysis of Variance and the likelihood ratio χ^2 in contingency tables:
 - ★ Deviance can be used to:
 - evaluate **goodness-of-fit**;
 - **test hypotheses**.

- ★ Dev is a function of θ and ϕ .
- ★ When ϕ is known then:
 - $Dev \sim \chi^2(n - p)$ where n is the number of observations and p is the number of parameters in $\underline{\beta}$ (including intercept, if any);
 - therefore, Dev can be used as a χ^2 statistics to test the goodness-of-fit (GOF) of the model.
- ★ When ϕ is unknown, it can be estimated and used to compute the **scaled deviance**:

$$Dev^*(\underline{\hat{\mu}}; \underline{\mathbf{y}}) = Dev(\underline{\hat{\mu}}; \underline{\mathbf{y}}) / \hat{\phi}$$

- ★ If $\underline{\beta}$ is partitioned in $\underline{\beta}_1$ and $\underline{\beta}_2$ then

$$\underline{X}\underline{\beta} = \underline{X}_1\underline{\beta}_1 + \underline{X}_2\underline{\beta}_2$$

- ★ The difference between the deviance of the full model and that of the model fitting $\underline{X}_1\underline{\beta}_1$ can be used as likelihood ratio (LR) to test

$$H_0 : \underline{\beta}_2 = \underline{0}$$

$$LR \sim \chi^2_{p - p_1}$$

where p_1 is the number of parameters in $\underline{\beta}_1$.

- ★ When ϕ is unknown, the scaled difference can be used to test hypotheses.

- ★ It can be shown that

$$Var(\underline{\beta}) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1}$$

then

$$Var(\mathbf{K}^T \underline{\beta}) = \mathbf{K}^T (\mathbf{X}^T \mathbf{W} \mathbf{X}) \mathbf{K}$$

where $\mathbf{K}^T \underline{\beta}$ is an estimable function.

- ★ The **Wald statistic** for $H_0 = \mathbf{K}^T \underline{\beta} = \mathbf{K}^T \underline{\beta}_0$ is

$$(\mathbf{K}^T \underline{\beta} - \mathbf{K}^T \underline{\beta}_0)^T \left[\mathbf{K}^T (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{K} \right] (\mathbf{K}^T \underline{\beta} - \mathbf{K}^T \underline{\beta}_0)$$

- ★ The scalar form of the Wald statistic is

$$\frac{\beta - \beta_0}{Var(\beta - \beta_0)}$$

- ★ The Wald statistic is distributed as a χ^2_ν where $\nu = \text{rank}(\mathbf{K})$.

- ★ Residuals are computed as in the linear model:

$$r_i = y_i - \hat{y}_i$$

- ★ Residual diagnostic test the fulfillment of the assumptions of the generalized linear model:
 - 1 testing the assumption of homogeneity of error variance;
 - 2 testing for outliers.

- ★ The assumption of homogeneity of error variance can be performed with the Levene's test.
- ★ Levene's test is performed by computing:

$$Z_{ij} = |y_{ij} - E(y_{.j})|$$

and then computing an F test on the j groups.

- ★ A nonsignificant result indicates no heteroskedasticity.

★ Outliers can be detected:

- ❶ Look for standardized residuals greater than 3.5 or less than -3.5
- ❷ and look for high Cook's D , greater than $4p/(n - p - 1)$. For instance, if $n = 100$, $p = 5$ then high Cook's D are higher than $4 \cdot 5 / (100 - 5 - 1) = 20/94$.

★ Standardized residuals are computed by

$$\frac{r_i - \hat{\mu}_r}{\hat{\sigma}_r}$$

where $\hat{\mu}_r$ and $\hat{\sigma}_r$ are the mean and the standard deviation of r_i , respectively.

- ★ Cook's D are computed by

$$\text{Cook's } D_i = \left(\frac{1}{p} \right) \left(\frac{h_{ii}}{1 - h_{ii}} \right) \left(\frac{r_i^2}{\hat{\sigma}^2 \cdot (1 - h_{ii})} \right)$$

where h_{ii} is the i -th diagonal entry of the hat matrix \mathbf{H} .

- ★ Cook and Weisberg (1982) suggested that values of D_i that exceed 50% of the F distribution with p and $n - p$ degrees of freedom are large

- ★ The hat matrix \mathbf{H} transforms \mathbf{Y} into the predicted scores. It is computed by

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

- ★ The diagonals for the hat matrix indicate which values could be outliers or not.
- ★ The diagonals are therefore measures of leverage.
- ★ Leverage is bounded by two limits: $1/n$ and 1. The closer the leverage is to unity, the more leverage the value has.
- ★ The trace of the hat matrix is equal to the number of variables in the model.
- ★ When the leverage is greater than $2p/n$ there is high leverage according to Belsley et al. (1980), cited in Long J.F., Modern Methods in Data Analysis (page 262). For smaller samples, Vellman and Welsch (1981) suggested that $3p/n$ is the criterion.

- ★ The leverage of outliers can be assessed:
 - constructing and analyzing studentized residuals;
 - constructing and analyzing the leverage of the high and low studentized residuals;
 - using Cook's D to help determine how problematic outliers are.

- ★ Studentized residuals are computed by

$$\frac{r_i}{\hat{\sigma}\sqrt{1 - h_{ii}}}$$

- ★ A Generalized Linear Mixed Model (GLMM) is an extension to the GLM in which the linear predictor contains random effects in addition to the usual fixed effects.
- ★ A GLMM use the **quasi-likelihood**, an extension of likelihood methods.
- ★ The quasi-likelihood function is defined as follows.
 - Given n observation y_i ($i = 1, \dots, n$) such that $E(y_i) = \mu_i$, $Var(y_i) \propto Var(\mu_i)$ (known function), $\mu = function(\beta_1, \dots, \beta_p)$ the quasi-likelihood function is defined by

$$\frac{\partial Q(\mu_i, y_i)}{\partial \mu_i} = \frac{y_i - \mu_i}{Var(\mu_i)}$$

- An example of $\mu = function(\beta_1, \dots, \beta_p)$ is given by $\mu_i = h(\eta_i)$, $\eta_i = \sum_j X_{ij}\beta_j$, that is the inverse link function of a GLM
- The log-likelihood function is a special case of $Q(\mu_i, y_i)$.

- ★ In the normal errors mixed model $\underline{y} = \underline{X}\underline{\beta} + \underline{Z}\underline{u} + \underline{e}$
 - the conditional mean of the observations given the random model effects is $E(\underline{y}|\underline{u}) = \underline{\mu} = \underline{X}\underline{\beta} + \underline{Z}\underline{u}$,
 - the conditional variance is $Var(\underline{y}|\underline{u}) = \underline{R} = Var(\underline{e})$.

- ★ Observation can be described as $\underline{y} = \underline{\mu} + \underline{e}$ where $\underline{\mu}$ is the conditional mean, $E(\underline{y}|\underline{u})$, $u \sim MVN(\underline{0}, \underline{G})$.

- ★ In the generalized linear mixed model, the **conditional distribution** of $\underline{y}|\underline{u}$ plays the same role as the distribution of \underline{y} in the fixed effects generalized linear model.
- ★ The **conditional quasi-likelihood** of an observation, $\underline{y}_i|\underline{y}$ is

$$Q(u_i, y_i | u_i) = \left[\frac{y_i \theta_i - b(\theta_i)}{a(\phi_i)} \right]$$

- ★ The **joint quasi-likelihood** of an observation, $\underline{y}_i | \underline{y}$ of the observations, \underline{y} , is the sum of the quasi-likelihood of $\underline{y} | \underline{u}$ and \underline{u} . In matrix terms, the joint quasi-likelihood is

$$Q(\underline{\mu}, \underline{u}; \underline{y}) = \left[\underline{y}^T \underline{A}^{-1} \underline{\theta} - (\underline{b}_{\underline{\theta}}^{1/2})^T \underline{A}^{-1} \underline{b}_{\underline{\theta}}^{1/2} \right] + \frac{1}{2} \underline{u}^T \underline{G}^{-1} \underline{u}$$

where:

- \underline{A} is the matrix of $a(\phi_i)$,
- $\underline{\theta}$ is the vector of $\theta(\mu_i)$,
- $\underline{b}_{\underline{\theta}}$ is the vector of $b(\theta_i)$.

- ★ The strategies for fitting a GLMM to $E(\underline{y})$ are the same as those for a GLM:
 - the form of the conditional quasi-likelihood determines the **variance structure** and contains the natural parameter,
 - the conditional mean function $\theta(\mu)$ is used a **canonical link**.

$$\eta = \underline{X}\underline{\beta} + \underline{Z}\underline{u}$$

$$\eta = g(\mu)$$

- ★ The strategies for fitting a GLMM to $E(\underline{y})$ are the same as those for a GLM:
 - the form of the conditional quasi-likelihood determines the **variance structure** and contains the natural parameter,
 - the conditional mean function $\theta(\mu)$ is used a **canonical link**.

$$\eta = \underline{X}\underline{\beta} + \underline{Z}\underline{u}$$

$$\eta = g(\mu)$$