

LOGISTIC REGRESSION WITH GROUPED DATA

Let's consider again the beetle data.

The experiment has been run on several beetles for every dose: I can count how many beetles are dead or alive for each level. I obtain the grouped data.

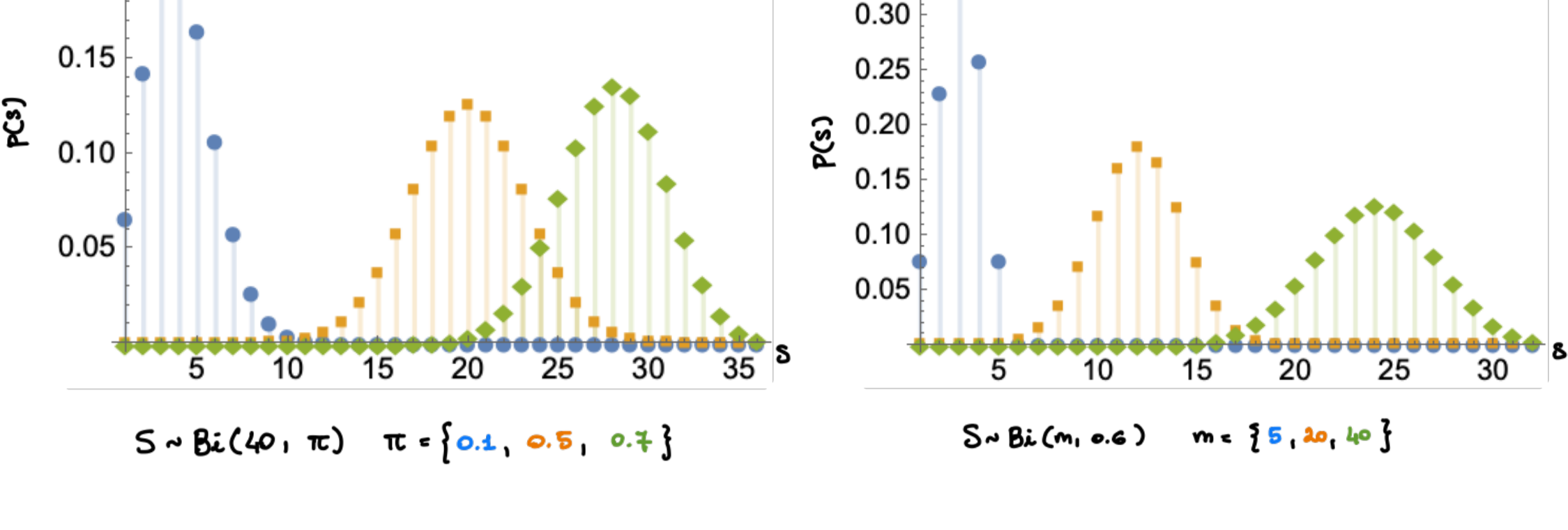
# killed (s)	6	13	...	60
# alive (o)	53	47	...	0
dose (x _i)	1.69	1.714	...	1.93

For the grouped data, an adequate distribution is the BINOMIAL distribution

Recall that

$S \sim \text{Bi}(m, \pi)$

- parameter space: $m \in \{0, 1, 2, \dots\}$ number of trials
 $\pi \in [0, 1]$ success probability
- support: $\mathcal{S} = \{0, 1, \dots, m\}$ number of successes on m trials
- probability mass function: $P_S(s; m, \pi) = P(S=s) = \binom{m}{s} \pi^s (1-\pi)^{m-s}$ with $\binom{m}{s} = \frac{m!}{s!(m-s)!}$
- moments: $E[S] = m\pi$
 $\text{Var}(S) = m\pi(1-\pi)$
- relationship with the Bernoulli distribution: consider a sequence of m independent Bernoulli random variables T_1, \dots, T_m with common success probability π : $T_k \sim \text{Bern}(\pi)$ $k=1, \dots, m$ independent.
Then $S = \sum_{k=1}^m T_k \sim \text{Bi}(m, \pi)$.



How do we define a model for grouped data, eg, in the beetle example?

Assume that in the ungrouped data we observed $T_k \sim \text{Bern}(\pi(x_k))$ $k=1, \dots, N$, with N the total number of beetles that were used in the experiment, and $\pi(x_k)$ the probability of "success" using a dose equal to x_k . However, the experiment was repeated several times for each poison level.

Let's denote with n the number of different levels of poison used in the experiment.

For each dose level x_i ($i=1, \dots, n$), m_i beetles were observed: we can group together the outcome of the experiment for each experimental condition.

Indeed, beetles with a dose = x_i all have the same probability = $\pi(x_i)$.

The exp(dose) is x_i $i=1, \dots, n$.

For each level x_i , we count the number of dead and alive beetles.

$S_i = \sum_{k=1}^{m_i} \mathbb{1}(T_k = 1 | x_k = x_i)$ number of successes at a dose = x_i

m_i is the total number of beetles observed at a dose x_i

$m_i = \sum_{k=1}^{m_i} \mathbb{1}(x_k = x_i)$

Since the T_k $k=1, \dots, N$ were independent, with distribution $\text{Bern}(\pi(x_k)) \rightarrow$ the success probability is common for r.v. with the same x_k . Hence the distribution of S_i is

$S_i \sim \text{Bi}(m_i, \pi_i = \pi(x_i))$ $i=1, \dots, n$ independent with support $\{0, 1, \dots, m_i\}$.

Indeed the grouped data can be expressed as

S	m	x	S	m-S	x
s_1	m_1	x_1	s_1	$m_1 - s_1$	x_1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
s_i	m_i	x_i	s_i	$m_i - s_i$	x_i
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
s_n	m_n	x_n	s_n	$m_n - s_n$	x_n

successes # trials # successes # failures

LOGISTIC REGRESSION (general case: grouped data)

With GLMs we model the MEAN of the random variables.

In this case we have S_1, \dots, S_n independent, $S_i \sim \text{Bi}(m_i, \pi_i)$

$E[S_i] = m_i \pi_i$

If we model directly the S_i , we study (m_i, π_i) . But in a study the quantity of interest is actually $\pi(x_i)$: the success probability at a level x_i (not m_i, π_i , also notice that m_i changes with i).

How do we define a model for π_i ?

Consider a Transformation of the random variables

$Y_i = \frac{S_i}{m_i}$ $i=1, \dots, n$

The expected value is $E[Y_i] = E[\frac{S_i}{m_i}] = \frac{1}{m_i} E[S_i] = \frac{m_i \pi_i}{m_i} = \pi_i$

The mean of Y_i is our parameter of interest π_i .

Support $\mathcal{Y} = \{0, \frac{1}{m_i}, \frac{2}{m_i}, \dots, \frac{m_i-1}{m_i}, 1\}$

what is the distribution of these new r.v.?

$P(Y_i = y_i) = P(\frac{S_i}{m_i} = y_i) = P(S_i = y_i m_i) = \binom{m_i}{m_i y_i} \pi_i^{y_i m_i} (1-\pi_i)^{m_i - y_i m_i} = P_S(m_i y_i; m_i, \pi_i)$
 $S_i \sim \text{Bi}(m_i, \pi_i)$

i.e. $m_i Y_i \sim \text{Bi}(m_i, \pi_i)$ $i=1, \dots, n$ independent.

It is possible to show that the distribution of Y_i is in the exponential family.

$\text{Var}(Y_i) = \text{Var}(\frac{S_i}{m_i}) = \frac{1}{m_i^2} \text{Var}(S_i) = \frac{1}{m_i^2} m_i \pi_i (1-\pi_i) = \frac{\pi_i (1-\pi_i)}{m_i}$ heteroscedastic

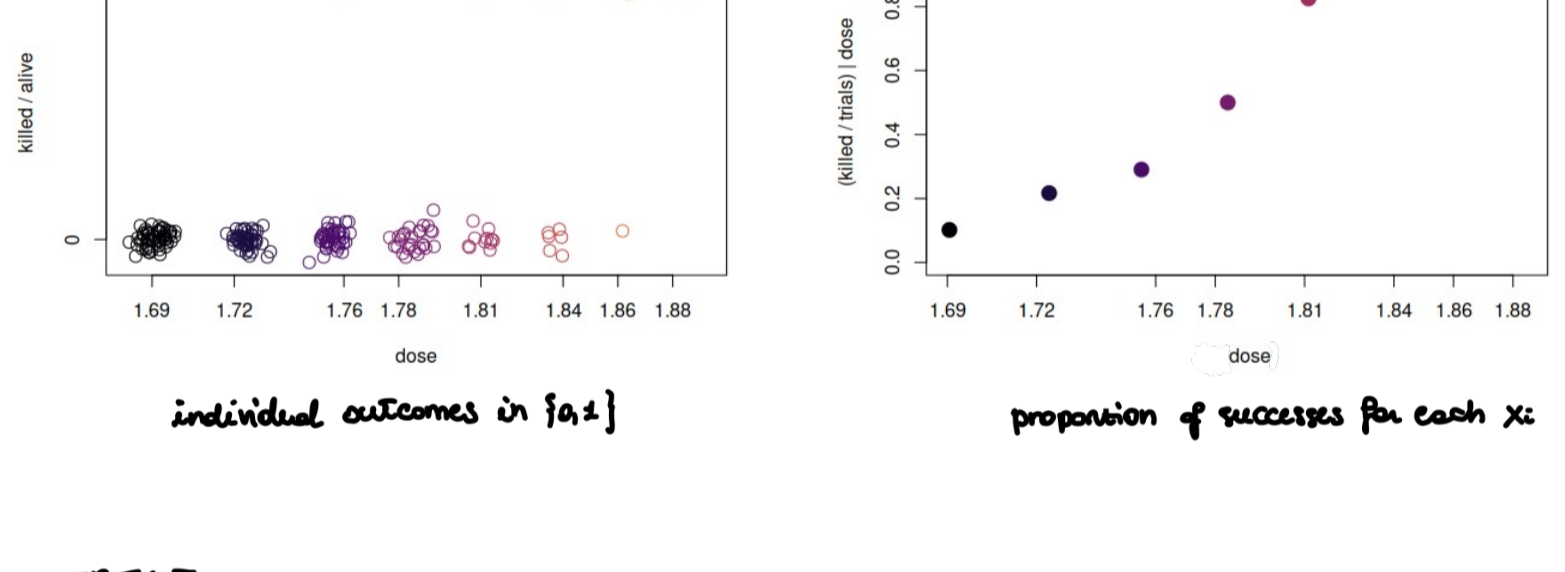
We can fit a glm on these new random variables.

The model is basically the same we have seen for $\{0, 1\}$ data.

The canonical link function is again $g(\pi_i) = \text{log} \frac{\pi_i}{1-\pi_i} = \eta_i = \frac{\beta^T X_i}{\beta}$

The interpretation of the parameters is the same.

Data visualisation with ungrouped and (transformed) grouped data



INFERENCE:

$P_{\beta}(y_i; m_i, \pi_i) = P_{S_i}(m_i y_i; m_i, \pi_i) = \binom{m_i}{m_i y_i} \pi_i^{y_i m_i} (1-\pi_i)^{m_i - y_i m_i}$ with $\pi_i = \frac{e^{\beta^T X_i}}{1 + e^{\beta^T X_i}}$

$L(\beta) = \prod_{i=1}^n P_{\beta}(y_i; m_i, \pi_i) = \prod_{i=1}^n \binom{m_i}{m_i y_i} \pi_i^{y_i m_i} (1-\pi_i)^{m_i - y_i m_i}$
 $\propto \prod_{i=1}^n \pi_i^{y_i m_i} (1-\pi_i)^{m_i - y_i m_i}$

$\ell(\beta) = \sum_{i=1}^n \{ y_i m_i \text{log} \pi_i + m_i (1-y_i) \text{log} (1-\pi_i) \} = \sum_{i=1}^n \{ m_i [y_i \text{log} \pi_i + (1-y_i) \text{log} (1-\pi_i)] \}$
this is what we had in the Bernoulli case

$\ell_x(\beta) = \sum_{i=1}^n \{ m_i x_i (y_i - \frac{e^{\beta^T X_i}}{1 + e^{\beta^T X_i}}) \} = \sum_{i=1}^n x_i (m_i y_i - m_i \frac{e^{\beta^T X_i}}{1 + e^{\beta^T X_i}}) = \sum_{i=1}^n x_i (m_i y_i - m_i \pi_i)$

the likelihood equations are $\sum_{i=1}^n x_i m_i y_i = \sum_{i=1}^n x_i m_i \pi_i$
 $\sum_{i=1}^n x_i m_i y_i = \sum_{i=1}^n x_i m_i \frac{e^{\beta^T X_i}}{1 + e^{\beta^T X_i}}$

$\frac{\partial^2 \ell(\beta)}{\partial \beta_j \partial \beta_k} = - \sum_{i=1}^n m_i x_{ij} x_{ik} \pi_i (1-\pi_i)$

$j(\beta) = - \ell_x(\beta) = X^T U X$ with $U = \text{diag} \{ m_1 \pi_1 (1-\pi_1), \dots, m_n \pi_n (1-\pi_n) \} = U(\beta)$

• INFERENCE about the ESTIMATOR of the REGRESSION COEFFICIENT

$\hat{\beta} \sim N_p(\beta, j(\hat{\beta})^{-1})$

• TEST about NESTED MODELS (about subsets of β)

$\beta = \begin{bmatrix} \beta^{(0)} \\ \beta^{(1)} \end{bmatrix} \in \mathbb{R}^p$
 $\beta = \begin{bmatrix} \beta^{(0)} \\ \beta^{(1)} \end{bmatrix} \in \mathbb{R}^p$

$H_0: \beta^{(1)} = 0$
 $H_1: \beta^{(1)} \neq 0$

likelihood ratio test:

$W = 2 \{ \hat{\ell}(\text{model}) - \hat{\ell}(\text{restricted}) \} \sim \chi^2_{p-p_0}$ under H_0

$W^{obs} = 2 \left\{ \sum_{i=1}^n \{ m_i [y_i \text{log} \hat{\pi}_i + (1-y_i) \text{log} (1-\hat{\pi}_i)] - m_i [y_i \text{log} \pi_0 + (1-y_i) \text{log} (1-\pi_0)] \} \right\}$
 $= 2 \left\{ \sum_{i=1}^n \left\{ m_i \left[y_i \text{log} \frac{\hat{\pi}_i}{\pi_0} + (1-y_i) \text{log} \frac{(1-\hat{\pi}_i)}{(1-\pi_0)} \right] \right\} \right\}$

• TEST about the OVERALL SIGNIFICANCE

$H_0: \beta_2 = \dots = \beta_p = 0$
 $H_1: H_0$

as usual, $W \sim \chi^2_{p-1}$ under H_0

under H_0 we have a common $\pi = \pi^0$ for all $i=1, \dots, n$

what is the estimate in this case?

$\ell(\beta) = \sum_{i=1}^n \{ y_i m_i \text{log} \pi + m_i (1-y_i) \text{log} (1-\pi) \}$

$\ell_x(\beta) = \sum_{i=1}^n \{ y_i m_i \frac{x_i}{\pi} - \frac{m_i (1-y_i)}{1-\pi} \}$

$\ell_x(\beta) = 0 \Rightarrow \sum y_i m_i - \pi \sum \frac{y_i m_i}{\pi} - \pi \sum m_i + \pi \sum \frac{m_i (1-y_i)}{1-\pi} = 0$
 $\frac{\pi}{1-\pi} = \frac{\sum y_i m_i}{\sum m_i} = \frac{\# \text{ successes}}{\# \text{ trials}} = \frac{\sum y_i}{N} = \bar{y}$

• DEVIANCES

I set a separate parameter π_i Y_i : saturated model

$\ell(\pi_i) = m_i y_i \text{log} \pi_i + m_i (1-y_i) \text{log} (1-\pi_i)$

$\ell_x(\pi_i) = \frac{m_i y_i}{\pi_i} - \frac{m_i (1-y_i)}{1-\pi_i}$

$\ell_x(\pi_i) = 0 \Rightarrow m_i y_i - \frac{y_i m_i}{\pi_i} - \pi_i m_i + \frac{m_i (1-y_i)}{1-\pi_i} = 0$
 $\pi_i = y_i$

$\ell(\hat{\pi}_i) = m_i y_i \text{log} y_i + m_i (1-y_i) \text{log} (1-y_i)$ ← now it is not always 0!
 $y_i \in]0, 1[$, $\frac{1}{m_i}, \frac{2}{m_i}, \dots, 1$

So we have now a proper deviance

$D = 2 \{ \hat{\ell}(\text{saturated}) - \hat{\ell}(\text{model}) \}$
 $= 2 \left\{ \sum_{i=1}^n \left[m_i y_i \text{log} y_i + m_i (1-y_i) \text{log} (1-y_i) - m_i y_i \text{log} \hat{\pi}_i - m_i (1-y_i) \text{log} (1-\hat{\pi}_i) \right] \right\}$
 $= 2 \left\{ \sum_{i=1}^n m_i \left[y_i \text{log} \frac{y_i}{\hat{\pi}_i} + (1-y_i) \text{log} \frac{(1-y_i)}{(1-\hat{\pi}_i)} \right] \right\}$

With grouped data, if all m_i are large, we have an approximate distribution for D.

In this case, the deviance can be used as a measure of the goodness of fit of the model:

small values indicate a good fit.

$D < n-p$ is generally ok (where n is the number of different configurations of the covariates, e.g. the number of levels of $x_i = \text{exp}(\text{dose}_i)$)

• RESIDUALS

Similarly to the Poisson regression case, we can carry out the model checking through graphical analysis of the residuals.

Pearson's residuals: $e_i = \frac{y_i - \hat{\pi}_i}{\sqrt{V(\hat{\pi}_i)}}$ $i=1, \dots, n$ with $V(\hat{\pi}_i) = \frac{1}{m_i} \hat{\pi}_i (1-\hat{\pi}_i)$