

Logistic regression for case-control studies (Prentice + Pyke, 1979, Biometrika)

Setup: random sample of size n_1 from case population ($Y=1$)
random sample of size n_0 from control population ($Y=0$)

Observations on an explanatory variable X (assume only one X):

$$\begin{aligned} X_{11}, \dots, X_{1n_1} & \text{ iid } \Pr(X=x | Y=1) \\ X_{01}, \dots, X_{0n_0} & \text{ iid } \Pr(X=x | Y=0) . \end{aligned}$$

Assume for simplicity that X has a finite number of possible levels x_0, \dots, x_k .

x_0 = reference level

1) Prospective model vs retrospective model

Prospective model

$$\text{logit } \Pr(Y=1 | X=x) = \alpha + \beta x .$$

Now consider $X=x$ vs $X=x_0$. From above

$$(1) \quad \frac{\Pr(Y=1 | X=x) / \Pr(Y=0 | X=x)}{\Pr(Y=1 | X=x_0) / \Pr(Y=0 | X=x_0)} = e^{\beta(x-x_0)} .$$

This implies

$$(2) \quad \frac{\Pr(X=x|Y=1)/\Pr(X=x_0|Y=1)}{\Pr(X=x|Y=0)/\Pr(X=x_0|Y=0)} = e^{\beta(x-x_0)}$$

(previously showed that $LHS(1) = LHS(2)$).

Define $\delta_k = \log [\Pr(X=x_k|Y=0) / \Pr(X=x_0|Y=0)]$.

(note $\delta_0 \equiv 0$)

Then we can write

$$\Pr(X=x_k|Y=0) = \frac{e^{\delta_k}}{\sum_{l=0}^K e^{\delta_l}}$$

$$\Pr(X=x_k|Y=1) = \frac{e^{\delta_k + \beta x_k}}{\sum_{l=0}^K e^{\delta_l + \beta x_l}}$$

- 2) Approaches to inference :
- a) likelihood inference based
directly on above
 - b) likelihood inference based
on rewritten form of model

3) Rewrite model

Define $q_k = (1-p) \Pr(X=x_k | Y=0) + p \Pr(X=x_k | Y=1)$,

where $p = \frac{n_1}{n}$ ($n = n_0 + n_1$).

Also define $\delta_0 = \log \left[(1-p) / \sum_{l=0}^k e^{\gamma_l} \right]$

$\delta_1 = \log \left[p / \sum_{l=0}^k e^{\gamma_l + \beta x_l} \right]$

Then we can write

$$\Pr(X=x_k | Y=0) = \left(\frac{1}{1-p} \right) e^{\gamma_k + \delta_0}$$

$$= q_k (1-p)^{-1} \left[e^{\gamma_k + \delta_0} / (e^{\gamma_k + \delta_0} + e^{\gamma_k + \delta_1 + \beta x_k}) \right]$$

$$= q_k (1-p)^{-1} \left[1 / (1 + e^{\alpha + \beta x_k}) \right]$$

$$(\alpha = \delta_1 - \delta_0)$$

and $\Pr(X=x_k | Y=1) = q_k p^{-1} \left[e^{\alpha + \beta x_k} / (1 + e^{\alpha + \beta x_k}) \right]$.

Above subject to : $\sum_k q_k = 1$

$$\sum_k \Pr(X=x_k | Y=0) = 1$$

$$\sum_k \Pr(X=x_k | Y=1) = 1$$

} (*)

4) Likelihood-based inference : n_{ik} = # obs in grp i with $X=x_k$

$$L = \left[\prod_{k=0}^K \Pr(X=x_k | Y=0)^{n_{0k}} \right] \times \left[\prod_{k=0}^K \Pr(X=x_k | Y=1)^{n_{1k}} \right]$$

$$= \prod_{k=0}^K g_k^{n_{0k}+n_{1k}} (1-p)^{-n_{0k}} p^{-n_{1k}} \left(\frac{1}{1+e^{\alpha+\beta x_k}} \right)^{n_{0k}} \\ \times \left(\frac{e^{\alpha+\beta x_k}}{1+e^{\alpha+\beta x_k}} \right)^{n_{1k}}$$

$$\mathcal{L} = \log L = \text{const} + \sum_{k=0}^K (n_{0k}+n_{1k}) \log g_k \\ + \sum_{k=0}^K n_{1k} [\alpha + \beta x_k] \\ - \sum_{k=0}^K (n_{0k}+n_{1k}) \log (1+e^{\alpha+\beta x_k}).$$

Ignoring the constraints (*), \mathcal{L} is maximized by:

$$i) \hat{g}_k = \frac{n_{0k}+n_{1k}}{n}$$

ii) solving

$$\frac{\partial \mathcal{L}}{\partial \alpha} = \sum_{k=0}^K n_{1k} - \sum_{k=0}^K (n_{0k}+n_{1k}) \frac{e^{\alpha+\beta x_k}}{1+e^{\alpha+\beta x_k}} \stackrel{\text{set}}{=} 0$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = \sum_{k=0}^K n_{1k} x_k - \sum_{k=0}^K (n_{0k}+n_{1k}) x_k \frac{e^{\alpha+\beta x_k}}{1+e^{\alpha+\beta x_k}} \stackrel{\text{set}}{=} 0.$$

Same equations as prospective logistic regression

Now, in the development just above we ignored the constraints (*), but the solution the results happens to satisfy the constraints automatically. Indeed, we have the estimating equation

$$\sum_{k=0}^K n_{1k} - \sum_{k=0}^K (n_{0k} + n_{1k}) \frac{e^{\hat{\alpha} + \hat{\beta}x_k}}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}} = 0 \tag{1}$$

and we also have $\hat{q}_k = \frac{n_{0k} + n_{1k}}{n}$.

Combining these we get

$$\sum_{k=0}^K n_{1k} - \sum_{k=0}^K n \hat{q}_k \frac{e^{\hat{\alpha} + \hat{\beta}x_k}}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}} = 0$$

$$\Leftrightarrow n \sum_{k=0}^K \hat{q}_k \frac{e^{\hat{\alpha} + \hat{\beta}x_k}}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}} = n_1$$

$$\Leftrightarrow \sum_{k=0}^K p^{-1} \hat{q}_k \frac{e^{\hat{\alpha} + \hat{\beta}x_k}}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}} = 1$$

$$\Leftrightarrow \sum_{k=0}^K \hat{Pr}(X=x_k | Y=1) = 1$$

Further, going back to (1) and using the relation

$$\frac{1}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}} = 1 - \frac{e^{\hat{\alpha} + \hat{\beta}x_k}}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}}$$

we get

$$\sum_{k=0}^K n_{1k} - \sum_{k=0}^K (n_{0k} + n_{1k}) + \sum_{k=0}^K (n_{0k} + n_{1k}) \frac{1}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}} = 0$$

$$\Leftrightarrow \sum_{k=0}^K n \hat{q}_k \frac{1}{1 + e^{\hat{\alpha} + \hat{\beta}x_k}} = n_0$$

$$\Leftrightarrow \sum_{k=0}^K \hat{Pr}(X=x_k | Y=0) = 1$$

It turns out that the analysis based on the retrospective likelihood leads not only to the same estimates as an analysis based on the prospective likelihood, but also the asymptotic $Cov(\hat{\beta})$ is the same under the two analyses. The proof of this latter fact is more complex; see Prentice and Pyke (1979, *Biometrika*) for details.