

0.1 blogit: Bivariate Logistic Regression for Two Dichotomous Dependent Variables

Use the bivariate logistic regression model if you have two binary dependent variables (Y_1, Y_2), and wish to model them jointly as a function of some explanatory variables. Each pair of dependent variables (Y_{i1}, Y_{i2}) has four potential outcomes, ($Y_{i1} = 1, Y_{i2} = 1$), ($Y_{i1} = 1, Y_{i2} = 0$), ($Y_{i1} = 0, Y_{i2} = 1$), and ($Y_{i1} = 0, Y_{i2} = 0$). The joint probability for each of these four outcomes is modeled with three systematic components: the marginal $\Pr(Y_{i1} = 1)$ and $\Pr(Y_{i2} = 1)$, and the odds ratio ψ , which describes the dependence of one marginal on the other. Each of these systematic components may be modeled as functions of (possibly different) sets of explanatory variables.

Syntax

```
> z.out <- zelig(list(mu1 = Y1 ~ X1 + X2 ,
                    mu2 = Y2 ~ X1 + X3),
                model = "blogit", data = mydata)
> x.out <- setx(z.out)
> s.out <- sim(z.out, x = x.out)
```

Input Values

In every bivariate logit specification, there are three equations which correspond to each dependent variable (Y_1, Y_2), and ψ , the odds ratio. You should provide a list of formulas for each equation or, you may use `cbind()` if the right hand side is the same for both equations

```
> formulae <- list(cbind(Y1, Y2) ~ X1 + X2)
```

which means that all the explanatory variables in equations 1 and 2 (corresponding to Y_1 and Y_2) are included, but only an intercept is estimated (all explanatory variables are omitted) for equation 3 (ψ).

You may use the function `tag()` to constrain variables across equations:

```
> formulae <- list(mu1 = y1 ~ x1 + tag(x3, "x3"), mu2 = y2 ~ x2 +
+   tag(x3, "x3"))
```

where `tag()` is a special function that constrains variables to have the same effect across equations. Thus, the coefficient for `x3` in equation `mu1` is constrained to be equal to the coefficient for `x3` in equation `mu2`.

Examples

1. Basic Example

Load the data and estimate the model:

```
> data(sanction)
```

```
> z.out1 <- zelig(cbind(import, export) ~ coop + cost + target,  
+   model = "blogit", data = sanction)
```

By default, `zelig()` estimates two effect parameters for each explanatory variable in addition to the odds ratio parameter; this formulation is parametrically independent (estimating unconstrained effects for each explanatory variable), but stochastically dependent because the models share an odds ratio.

Generate baseline values for the explanatory variables (with cost set to 1, net gain to sender) and alternative values (with cost set to 4, major loss to sender):

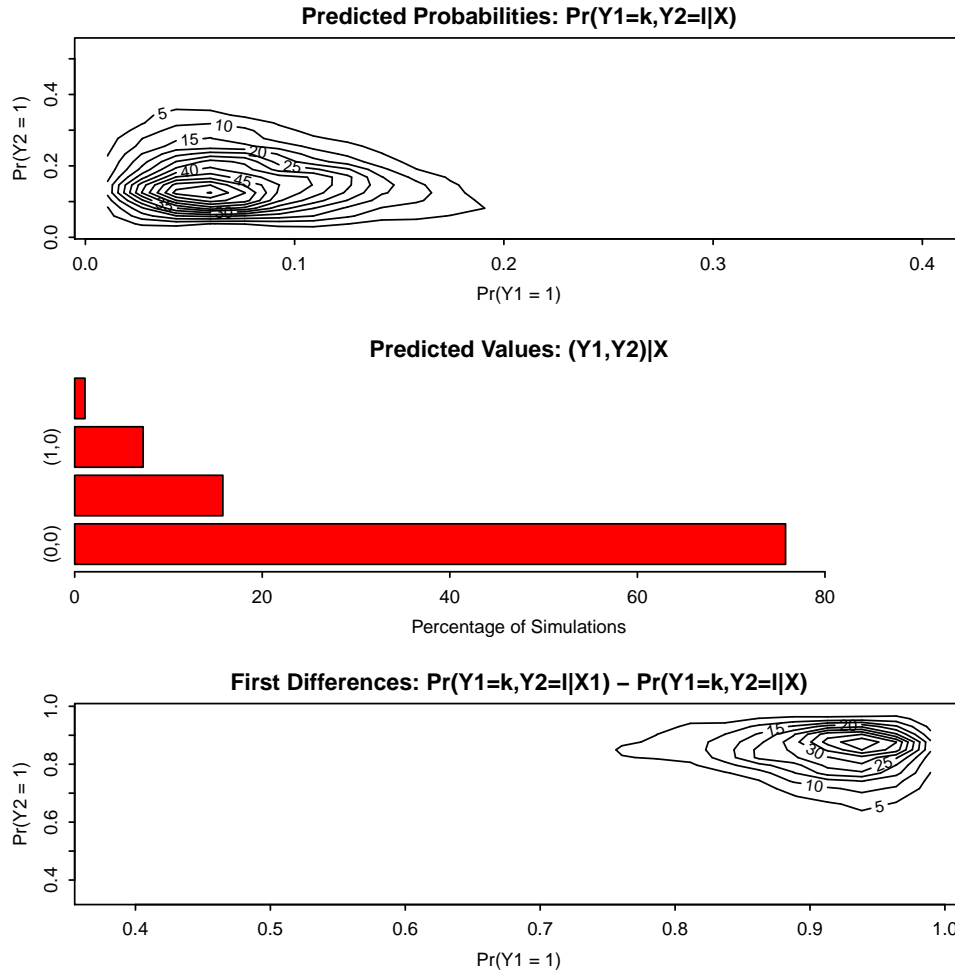
```
> x.low <- setx(z.out1, cost = 1)
```

```
> x.high <- setx(z.out1, cost = 4)
```

Simulate fitted values and first differences:

```
> s.out1 <- sim(z.out1, x = x.low, x1 = x.high)  
> summary(s.out1)
```

```
> plot(s.out1)
```



2. Joint Estimation of a Model with Different Sets of Explanatory Variables

Using sample data `sanction`, estimate the statistical model, with `import` a function of `coop` in the first equation and `export` a function of `cost` and `target` in the second equation:

```
> z.out2 <- zelig(list(import ~ coop, export ~ cost + target),
+   model = "blogit", data = sanction)
> summary(z.out2)
```

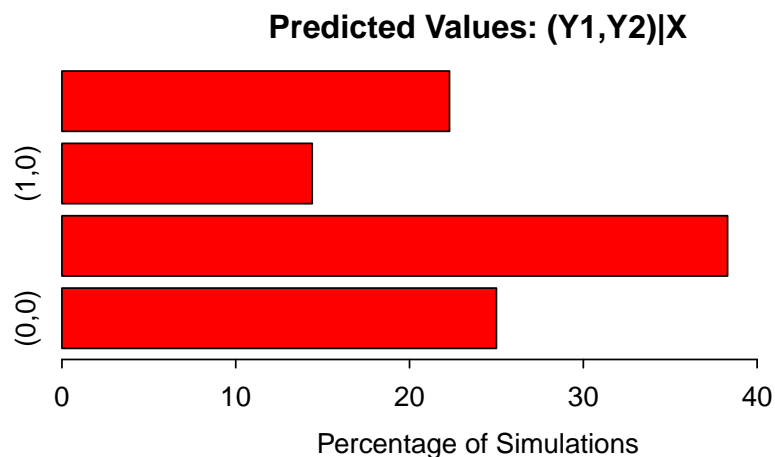
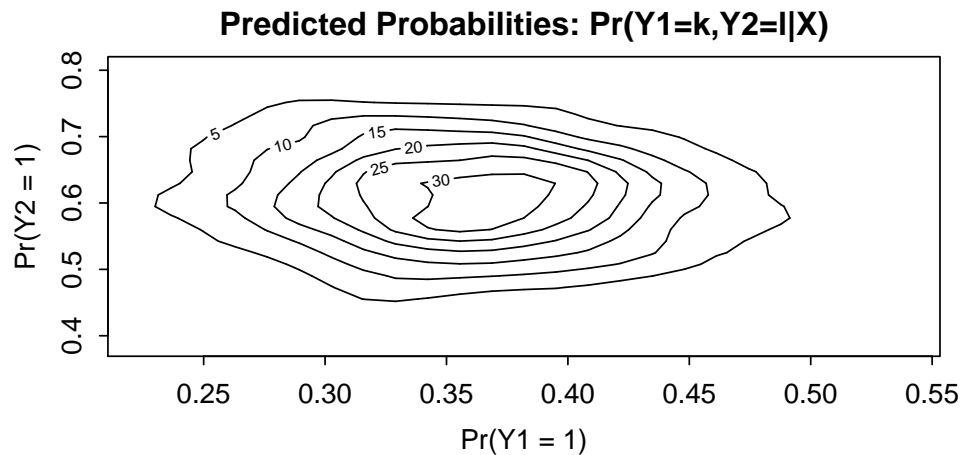
Set the explanatory variables to their means:

```
> x.out2 <- setx(z.out2)
```

Simulate draws from the posterior distribution:

```
> s.out2 <- sim(z.out2, x = x.out2)
> summary(s.out2)
```

```
> plot(s.out2)
```



3. Joint Estimation of a Parametrically and Stochastically Dependent Model

Using the sample data `sanction` The bivariate model is parametrically dependent if Y_1 and Y_2 share some or all explanatory variables, *and* the effects of the shared explanatory variables are jointly estimated. For example,

```
> z.out3 <- zelig(list(import ~ tag(coop, "coop") + tag(cost, "cost") +
+   tag(target, "target"), export ~ tag(coop, "coop") + tag(cost,
+   "cost") + tag(target, "target")), model = "blogit", data = sanction)
> summary(z.out3)
```

Note that this model only returns one parameter estimate for each of `coop`, `cost`, and `target`. Contrast this to Example 1 which returns two parameter estimates for each of the explanatory variables.

Set values for the explanatory variables:

```
> x.out3 <- setx(z.out3, cost = 1:4)
```

Draw simulated expected values:

```
> s.out3 <- sim(z.out3, x = x.out3)
> summary(s.out3)
```

Model

For each observation, define two binary dependent variables, Y_1 and Y_2 , each of which take the value of either 0 or 1 (in the following, we suppress the observation index). We model the joint outcome (Y_1, Y_2) using a marginal probability for each dependent variable, and the odds ratio, which parameterizes the relationship between the two dependent variables. Define Y_{rs} such that it is equal to 1 when $Y_1 = r$ and $Y_2 = s$ and is 0 otherwise, where r and s take a value of either 0 or 1. Then, the model is defined as follows,

- The *stochastic component* is

$$\begin{aligned} Y_{11} &\sim \text{Bernoulli}(y_{11} \mid \pi_{11}) \\ Y_{10} &\sim \text{Bernoulli}(y_{10} \mid \pi_{10}) \\ Y_{01} &\sim \text{Bernoulli}(y_{01} \mid \pi_{01}) \end{aligned}$$

where $\pi_{rs} = \Pr(Y_1 = r, Y_2 = s)$ is the joint probability, and $\pi_{00} = 1 - \pi_{11} - \pi_{10} - \pi_{01}$.

- The *systematic components* model the marginal probabilities, $\pi_j = \Pr(Y_j = 1)$, as well as the odds ratio. The odds ratio is defined as $\psi = \pi_{00}\pi_{01}/\pi_{10}\pi_{11}$ and describes the relationship between the two outcomes. Thus, for each observation we have

$$\begin{aligned} \pi_j &= \frac{1}{1 + \exp(-x_j\beta_j)} \quad \text{for } j = 1, 2, \\ \psi &= \exp(x_3\beta_3). \end{aligned}$$

Quantities of Interest

- The expected values (`qi$ev`) for the bivariate logit model are the predicted joint probabilities. Simulations of β_1 , β_2 , and β_3 (drawn from their sampling distributions) are substituted into the systematic components (π_1, π_2, ψ) to find simulations of the predicted joint probabilities:

$$\begin{aligned} \pi_{11} &= \begin{cases} \frac{1}{2}(\psi - 1)^{-1} - a - \sqrt{a^2 + b} & \text{for } \psi \neq 1 \\ \pi_1\pi_2 & \text{for } \psi = 1 \end{cases}, \\ \pi_{10} &= \pi_1 - \pi_{11}, \\ \pi_{01} &= \pi_2 - \pi_{11}, \\ \pi_{00} &= 1 - \pi_{10} - \pi_{01} - \pi_{11}, \end{aligned}$$

where $a = 1 + (\pi_1 + \pi_2)(\psi - 1)$, $b = -4\psi(\psi - 1)\pi_1\pi_2$, and the joint probabilities for each observation must sum to one. For n simulations, the expected values form an $n \times 4$ matrix for each observation in \mathbf{x} .

- The predicted values (`qi$pr`) are draws from the multinomial distribution given the expected joint probabilities.
- The first differences (`qi$fd`) for each of the predicted joint probabilities are given by

$$\text{FD}_{rs} = \Pr(Y_1 = r, Y_2 = s \mid x_1) - \Pr(Y_1 = r, Y_2 = s \mid x).$$

- The risk ratio (`qi$rr`) for each of the predicted joint probabilities are given by

$$\text{RR}_{rs} = \frac{\Pr(Y_1 = r, Y_2 = s \mid x_1)}{\Pr(Y_1 = r, Y_2 = s \mid x)}$$

- In conditional prediction models, the average expected treatment effect (`att.ev`) for the treatment group is

$$\frac{1}{\sum_{i=1}^n t_i} \sum_{i:t_i=1}^n \{Y_{ij}(t_i = 1) - E[Y_{ij}(t_i = 0)]\} \text{ for } j = 1, 2,$$

where t_i is a binary explanatory variable defining the treatment ($t_i = 1$) and control ($t_i = 0$) groups. Variation in the simulations are due to uncertainty in simulating $E[Y_{ij}(t_i = 0)]$, the counterfactual expected value of Y_{ij} for observations in the treatment group, under the assumption that everything stays the same except that the treatment indicator is switched to $t_i = 0$.

- In conditional prediction models, the average predicted treatment effect (`att.pr`) for the treatment group is

$$\frac{1}{\sum_{i=1}^n t_i} \sum_{i:t_i=1}^n \left\{ Y_{ij}(t_i = 1) - \widehat{Y_{ij}(t_i = 0)} \right\} \text{ for } j = 1, 2,$$

where t_i is a binary explanatory variable defining the treatment ($t_i = 1$) and control ($t_i = 0$) groups. Variation in the simulations are due to uncertainty in simulating $\widehat{Y_{ij}(t_i = 0)}$, the counterfactual predicted value of Y_{ij} for observations in the treatment group, under the assumption that everything stays the same except that the treatment indicator is switched to $t_i = 0$.

Output Values

The output of each Zelig command contains useful information which you may view. For example, if you run `z.out <- zelig(y ~ x, model = "blogit", data)`, then you may examine the available information in `z.out` by using `names(z.out)`, see the `coefficients` by using `z.out$coefficients`, and obtain a default summary of information through `summary(z.out)`. Other elements available through the `$` operator are listed below.

- From the `zelig()` output object `z.out`, you may extract:
 - `coefficients`: the named vector of coefficients.
 - `fitted.values`: an $n \times 4$ matrix of the in-sample fitted values.
 - `predictors`: an $n \times 3$ matrix of the linear predictors $x_j\beta_j$.
 - `residuals`: an $n \times 3$ matrix of the residuals.
 - `df.residual`: the residual degrees of freedom.
 - `df.total`: the total degrees of freedom.
 - `rss`: the residual sum of squares.
 - `y`: an $n \times 2$ matrix of the dependent variables.
 - `zelig.data`: the input data frame if `save.data = TRUE`.
- From `summary(z.out)`, you may extract:
 - `coef3`: a table of the coefficients with their associated standard errors and t -statistics.
 - `cov.unscaled`: the variance-covariance matrix.
 - `pearson.resid`: an $n \times 3$ matrix of the Pearson residuals.
- From the `sim()` output object `s.out`, you may extract quantities of interest arranged as arrays indexed by simulation \times quantity \times \mathbf{x} -observation (for more than one \mathbf{x} -observation; otherwise the quantities are matrices). Available quantities are:
 - `qi$ev`: the simulated expected joint probabilities (or expected values) for the specified values of \mathbf{x} .
 - `qi$pr`: the simulated predicted outcomes drawn from a distribution defined by the expected joint probabilities.
 - `qi$fd`: the simulated first difference in the expected joint probabilities for the values specified in \mathbf{x} and $\mathbf{x1}$.
 - `qi$rr`: the simulated risk ratio in the predicted probabilities for given \mathbf{x} and $\mathbf{x1}$.
 - `qi$att.ev`: the simulated average expected treatment effect for the treated from conditional prediction models.
 - `qi$att.pr`: the simulated average predicted treatment effect for the treated from conditional prediction models.

How to Cite

To cite the *blogit* Zelig model:

Kosuke Imai, Gary King, and Oliva Lau. 2007. "blogit: Bivariate Logistic Regression for Dichotomous Dependent Variables" in Kosuke Imai, Gary King, and Olivia Lau, "Zelig: Everyone's Statistical Software," <http://gking.harvard.edu/zelig>

To cite Zelig as a whole, please reference these two sources:

Kosuke Imai, Gary King, and Olivia Lau. 2007. "Zelig: Everyone's Statistical Software," <http://GKing.harvard.edu/zelig>.

Imai, Kosuke, Gary King, and Olivia Lau. (2008). "Toward A Common Framework for Statistical Analysis and Development." *Journal of Computational and Graphical Statistics*, Vol. 17, No. 4 (December), pp. 892-913.

See also

The bivariate logit function is part of the VGAM package by Thomas Yee (Yee and Hastie 2003). In addition, advanced users may wish to refer to `help(vglm)` in the VGAM library. Additional documentation is available at <http://www.stat.auckland.ac.nz/~yee>. Sample data are from Martin (1992)

Bibliography

- Martin, L. (1992), *Coercive Cooperation: Explaining Multilateral Economic Sanctions*, Princeton University Press, please inquire with Lisa Martin before publishing results from these data, as this dataset includes errors that have since been corrected.
- Yee, T. W. and Hastie, T. J. (2003), “Reduced-rank vector generalized linear models,” *Statistical Modelling*, 3, 15–41.