

REGRESSION MODELS

IN

MORBIDITY ESTIMATION

September, 1982

Gordon E. Willmot, M.Math, ASA  
Actuarial Assistant  
Mutual Life of Canada

## Regression Models in Morbidity Estimates

### Abstract

A statistical model is presented which analyzes group health and dental data. The purpose is to determine the cost of these benefits with special reference to the effects of various characteristics of individuals in the group.

The methodology involves a combination of a binary logistic regression and a normal multiple linear regression.

## Regression Models In Morbidity Estimation

We are interested in the group insurance premium for any or all of the following benefits:

- i) drug benefits
- ii) hospital coverage (private or semi-private)
- iii) vision care
- iv) dental insurance
- v) supplementary health insurance (e.g. ambulance, wheelchair, etc.)

In particular, we wish to determine the effects of an individual's characteristics on the cost of supplying these coverages.

We denote by

$$x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$$

a given set of covariates which may include

- i) age
- ii) sex
- iii) occupation
- iv) location of residence
- v) marital status

among others.

For example, we might have

$$\begin{aligned} x_{i1} &= 1 && \text{if } 18 \leq \text{age} < 35 \\ &= 0 && \text{otherwise} \\ x_{i2} &= 1 && \text{if } 35 \leq \text{age} < 50 \\ &= 0 && \text{otherwise} \end{aligned}$$

$$\begin{aligned}
 x_{i3} &= 1 && \text{if } 50 \leq \text{age} < 65 \\
 &= 0 && \text{otherwise}
 \end{aligned}$$

$$\begin{aligned}
 x_{i4} &= 1 && \text{if sex is male} \\
 &= 0 && \text{if sex is female}
 \end{aligned}$$

etc.

If we denote by  $X$  the total claim amount for an individual for the premium paying period then we are interested in the random variable

$$X \mid x_i$$

or, more specifically, the first couple of moments of its associated distribution.

We assume we have the total amount of claims for each individual in the study for a period of time at least as great as the premium payment period (to eliminate the effects of seasonality on the claims). Furthermore, assume that all claim amounts are strictly positive so that  $X=0$  implies no claims occurred.

Then,

$$\Pr \{X=x \mid x_i\} = \Pr\{\geq 1 \text{ claim} \mid x_i\} \cdot \Pr\{X=x \mid x_i, \geq 1 \text{ claim}\}$$

We will estimate the two quantities on the right hand side of this equation separately.

Let us divide the data into  $n$  groups such that all covariates in a given group are equal (or roughly so). Suppose in group  $i$  there

are  $n_i$  individuals,  $d_i$  of which have at least 1 claim. Let us assume that  $\underline{x}_i$  is the covariate vector of the  $i$ th group.

### Part 1 - Logistic Regression

Let

$$p_i = \Pr \{ \geq 1 \text{ claim} \mid \underline{x}_i \}$$

Then it is reasonable to assume that  $d_i$  is an observation from a binomial variate with parameters  $n_i$  and  $p_i$ .

Let

$$p_i = \frac{e^{\beta_1 x_{i1} + \dots + \beta_p x_{ip}}}{1 + e^{\beta_1 x_{i1} + \dots + \beta_p x_{ip}}} = \frac{e^{\underline{x}_i \underline{\beta}}}{1 + e^{\underline{x}_i \underline{\beta}}} = \frac{e^{\theta_i}}{1 + e^{\theta_i}}$$

where  $\underline{\beta}^T = \{ \beta_1, \beta_2, \dots, \beta_p \}$ ,  $\theta_i = \underline{x}_i \underline{\beta}$

Then we are interested in making inferences about  $\underline{\beta}$ .

For notational simplicity, let

$$\underline{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{pmatrix} \quad n \times 1 \qquad \underline{X} = \begin{pmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \vdots \\ \underline{x}_n \end{pmatrix} \quad n \times p \qquad \underline{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix} \quad p \times 1$$

To get estimates for  $\beta$ , we use likelihood methods.

$$L(\beta) = \prod_{i=1}^n \frac{e^{x_i \beta \cdot d_i}}{(1 + e^{x_i \beta})^{n_i}} \quad (\text{since } d_i \text{ is binomial})$$

$$= \prod_{i=1}^n \frac{e^{\theta_i d_i}}{(1 + e^{\theta_i})^{n_i}}$$

Calculate the score function

$$S(\beta) = \left\{ \frac{\partial}{\partial \beta_i} \log L(\beta) \right\}_{p \times 1}$$

$$= \left\{ \sum_{i=1}^n \left[ d_i x_{ij} - \frac{n_i x_{ij} e^{x_i \beta}}{1 + e^{x_i \beta}} \right] \right\}_{p \times 1}$$

and the information matrix

$$I(\beta) = \left\{ - \frac{\partial^2 \log L(\beta)}{\partial \beta_j \partial \beta_k} \right\}_{p \times p}$$

$$= \left\{ \sum_{i=1}^n \frac{n_i x_{ij} x_{ik} e^{x_i \beta}}{(1 + e^{x_i \beta})^2} \right\}_{p \times p}$$

then the estimates  $\hat{\beta}$  may be calculated iteratively using Newton-Raphson techniques as follows ( $\hat{\beta}_i$  is the value of  $\hat{\beta}$  from the  $i$ th iteration) until  $\hat{\beta}_i$  is sufficiently close to  $\hat{\beta}_{i-1}$ .

$$\hat{\beta}_i = \hat{\beta}_{i-1} + I^{-1}(\hat{\beta}_{i-1}) \cdot S(\hat{\beta}_{i-1})$$

To get a starting value  $\hat{\beta}_0$ , we go back to the likelihood as a function of  $\theta$ . We have

$$L(\theta) = \prod_{i=1}^n \frac{e^{\theta_i d_i}}{(1 + e^{\theta_i})^{n_i}}$$

and so

$$\hat{\theta}_j = \log \frac{d_j}{n_j - d_j}$$

It can be shown using large sample theory techniques that (asymptotically)

$$\hat{\theta}_j \sim N \left\{ \theta_j, \frac{n_j}{d_j (n_j - d_j)} \right\}$$

or 
$$\log \frac{d_j}{n_j - d_j} \sim N \left\{ x_{i1} \beta_1 + \dots + x_{ip} \beta_p, \frac{n_j}{d_j (n_j - d_j)} \right\}$$

which is the usual weighted least squares multiple regression model. Thus, if this were the true result, we would have

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

where

$$V = \left\{ \sqrt{\frac{d_j (n_j - d_j)}{n_j}} x_{jk} \right\} \quad n \times p$$

$$y = \left\{ \sqrt{\frac{d_j (n_j - d_j)}{n_j}} \log \frac{d_j}{n_j - d_j} \right\}$$

This can thus serve as a good starting value for  $\hat{\beta}$  (ie. as  $\hat{\beta}_0$ ). Note that if  $d_j = 0$ , it is recommended to replace

$$\log \frac{d_j}{n_j - d_j} \quad \text{by} \quad \log \frac{0.5}{n_j + 0.5}$$

and if  $d_j = n_j$ , replace

$$\log \frac{d_j}{n_j - d_j} \quad \text{by} \quad \log \frac{n_j + 0.5}{0.5}$$

in the Y matrix above so that  $\hat{\theta}_j \neq \infty$ . This does not affect the final answer.

From large sample theory, it can be shown that  $\hat{\beta}$  is approximately normal with mean  $\hat{\beta}$  and variance covariance matrix  $I^{-1}(\hat{\beta})$ . Thus, the usual tests can be applied to the solutions  $\hat{\beta}$  for inference purposes (ie. to test whether  $\hat{\beta}_i = 0$  so that the corresponding covariates may be dropped).

To test the fit of the model, a one-tailed test involving the chi-squared statistic

$$\chi^2_{(n-p)} = \sum_{i=1}^n \frac{(d_i - n_i \hat{p}_i)^2}{n_i \hat{p}_i (1 - \hat{p}_i)} \quad \text{where}$$

$$\hat{p}_i = \frac{e^{x_i \hat{\beta}}}{1 + e^{x_i \hat{\beta}}}$$

(large values give evidence against the model).



## Part 2 - Multiple Linear Regression

In order to consider

$$X \mid \underline{x}_i \geq 1 \text{ claim}$$

we consider only the  $d_i$  individuals who had claims. Let  $T_{ij}$  be the total claim amount for person  $j$ , group  $i$ .

Our model is

$$f(T_{ij}) \sim N \{ \tau_1 x_{i1} + \dots + \tau_p x_{ip}, \sigma^2 \}$$

where  $f$  is chosen so that the data appears to be roughly normal.

Appropriate choices of  $f$  include the following

$$f(x) = x$$

$$f(x) = \log x$$

$$f(x) = \sqrt{x}$$

among others.

Let

$$Y_{ij} = f(t_{ij}), \quad \bar{y}_i = \frac{d_j}{\sum_{j=1}^{d_j}} Y_{ij}/d_j$$

Then

$$\bar{y}_i \sim N \{ \tau_1 x_{i1} + \dots + \tau_p x_{ip}, \sigma^2/d_i \}$$

The least squares estimates of  $\underline{\tau}$  are then

$$\hat{\underline{\tau}} = (V^T V)^{-1} V^T Z$$

where

$$V = \{ \sqrt{d_i} x_{ij} \}_{n \times p}, \quad Z = \{ \sqrt{d_i} \bar{y}_i \}_{n \times 1}$$

and

$$\hat{\tau} \sim N \{ \tau, \sigma^2 (V^T V)^{-1} \}$$

An estimate of  $\sigma^2$  is

$$S_1^2 = \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} (y_{ij} - \bar{y}_i)^2}{(\sum_{i=1}^n d_i) - n}$$

To test the fit of the model, compute the F statistic

$$F_{n-p, (\sum_{i=1}^n d_i) - n} = \frac{\sum_{i=1}^n d_i (\bar{y}_i - \hat{\tau}_1 x_{i1} - \dots - \hat{\tau}_p x_{ip})^2 \div (n-p)}{S_1^2}$$

This is the ratio of the "lack of fit" sum of squares to the "pure error" sum of squares. A one-tailed test is appropriate. If the model is good, a better estimate of  $\sigma^2$  is

$$S_2^2 = \frac{\sum_{i=1}^n \sum_{j=1}^{d_i} (y_{ij} - \hat{\tau}_1 x_{i1} - \dots - \hat{\tau}_p x_{ip})^2}{(\sum_{i=1}^n d_i) - p}$$

This problem could have been formulated as a weighted least squares problem with identical results. The use of the mean rather than the original data is primarily to cut down on the computing time.

### Part 3 The Distribution

We may now compute the distribution of  $X|x_i$ .

The density for  $x > 0$  is

$$g_i(x) = \hat{p}_i \cdot \phi \left\{ \frac{f(x) - x_i \hat{f}}{\hat{\sigma}} \right\} \cdot \frac{f'(x)}{\hat{\sigma}}$$

where

$$\phi(x) = \frac{e^{-\frac{x^2}{2}}}{\sqrt{2}}$$

Also

$$g_i(0) = \Pr\{X = 0|x_i\} = 1 - \int_{0+}^{\infty} g_i(x) dx$$

and

$$E\{X|x_i\} = \int_{0+}^{\infty} x g_i(x) dx$$

For a deductible of amount  $c$ , we are interest in the random variable

$$W|x_i = \max\{X-c, 0\} | x_i$$

$$\text{and } \Pr\{W=0|x_i\} = \Pr\{X_i < c | x_i\}$$

and the density for  $x > 0$  is

$$g_i(x+c)$$

The expected value is

$$\int_c^{\infty} (x-c) g_i(x) dx$$

I wish to thank Dr. R. J. Mackay and Dr. J. G. Kalbfleisch of the University of Waterloo for their extremely helpful advice and support in this venture.

### References

- 1) Theoretical Statistics, Cox and Hinkley, 1974, Chapman and Hall, London
- 2) Applied Regression Analysis, Draper and Smith, 1966, John Wiley and Sons, New York
- 3) Discrete Data Analysis Lecture Notes, J. G. Kalbfleisch, University of Waterloo, Waterloo