

Module 11: Probit Regression

Rebecca C. Steorts

Agenda

- ▶ Ordinal, numeric, and continuous variables
- ▶ Probit versus linear regression
- ▶ Full conditionals
- ▶ An application to data

Generalized Linear Regression

- ▶ Many datasets include variables whose distributions cannot be represented by the normal, binomial or Poisson distributions we have studied so far.
- ▶ Distributions of common survey variables such as age, education level and income generally cannot be accurately described by any of the above- mentioned sampling models.
- ▶ We extend models to situations where the data are not normal, by expressing non-normal random variables as functions of unobserved, “latent” normally distributed random variables.
- ▶ Multivariate normal and linear regression models then can be applied to the latent data.

Terminology

- ▶ We use the term ordinal to refer to any variable for which there is a logical ordering of the sample space.
- ▶ We use the term numeric to refer to variables that have meaningful numerical scales.
- ▶ We use the term continuous if a variable can have a value that is (roughly) any real number in an interval.

Data

- ▶ The 1994 General Social Survey provides data on variables DEG , $CHILD$ and $PDEG$ for a sample of individuals in the United States.
- ▶ DEG_i indicates the highest degree obtained by individual i , $CHILD_i$ is the number of children and $PDEG_i$ is the binary indicator of whether or not either parent obtained a college degree.

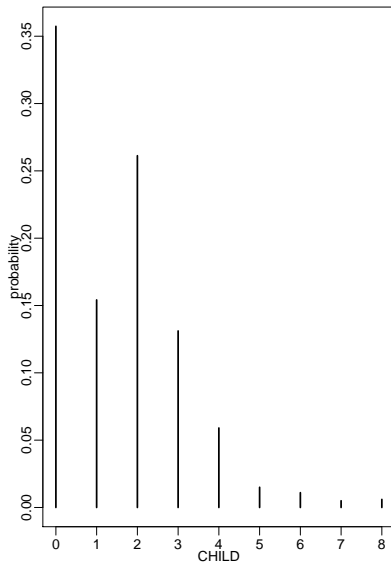
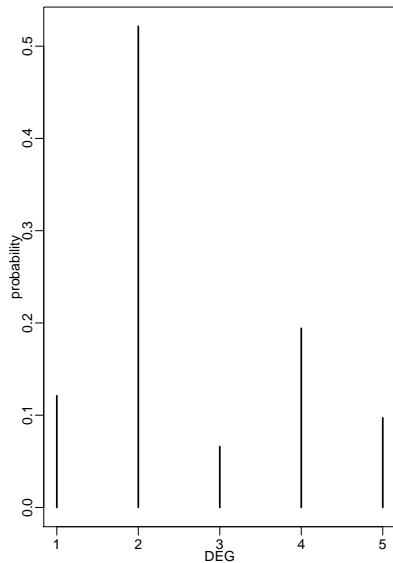
Using these data, we might be tempted to investigate the relationship between the variables with a linear regression model.

Data

```
dat<-read.table("http://lib.stat.cmu.edu/aoas/107/data.txt",header=TRUE)
head(dat)
```

##	INCOME	DEGREE	CHILDREN	PINCOME	PDEGREE	PCHILDREN	AGE
## 1	NA	1	3	3	1	5	59
## 2	11	0	3	NA	0	7	59
## 3	8	1	1	NA	0	9	25
## 4	25	3	2	NA	0	5	55
## 5	100	3	2	4	3	2	56
## 6	40	4	0	NA	4	5	36

Two ordinal variables having non-normal distributions



Probit regression

- ▶ Linear or generalized linear regression models, which assume a numeric scale to the data, may be appropriate for variables like CHILD, height or weight, but are not appropriate for non-numeric ordinal variables like DEG.
- ▶ This idea motivates a modeling technique known as ordered probit regression, in which we relate a variable Y to a vector of predictors x via a regression in terms of a latent variable Z .

Probit regression model

The model can be written as

$$Y_i = g(Z_i) \quad (1)$$

$$Z_i = \beta^T x_i + \epsilon_i \quad (2)$$

$$\epsilon_i \stackrel{iid}{\sim} \text{Normal}(0, 1) \quad (3)$$

$$\beta \sim \text{MVN}(0, n(X^T X)^{-1}) \quad (4)$$

g is any non-decreasing function.

Notation

- ▶ $X_{n \times p}$: regression features or covariates (design matrix)
- ▶ $Z_{n \times 1}$: latent variable
- ▶ $\mathbf{y}_{n \times 1}$: response variable (vector)
- ▶ $\beta_{p \times 1}$: vector of regression coefficients

The role of g

$$Y_{n \times 1} = g(Z) \quad (5)$$

$$Z_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1} \quad (6)$$

$$\epsilon_{n \times 1} \stackrel{iid}{\sim} \text{Normal}(0, I) \quad (7)$$

$$\beta_{p \times 1} \sim \text{MVN}(0, n(X^T X)^{-1}) \quad (8)$$

Suppose the sample space for Y takes on K values $\{1, 2, \dots, K\}$, then g can be described with $K - 1$ ordered parameters.

You can think of the values of g_1, \dots, g_{K-1} as thresholds so that moving past z will move y into the next (highest) category.

Full conditional of β

$$p(\beta \mid y, z, g) \propto p(\beta)p(z \mid \beta)$$

Using the MVN conjugacy that we looked at before, $\beta \mid y, z, g$ will be MVN where

$$E[\beta \mid z] = \frac{n}{n+1}(X^T X)^{-1}X^T z$$

$$\text{Var}[\beta \mid z] = \frac{n}{n+1}(X^T X)^{-1}$$

Full conditional of Z

Under the sampling model, the conditional distribution of

$$Z_i \mid \beta \sim \text{Normal}(\beta^T x_i, 1)$$

Given g and observing $Y_i = y_i$, we know that Z_i lies in the interval (g_{i-1}, g_i) .

Let $a = g_{i-1}$, $b = g_i$.

Then

$$p(z_i \mid \beta, z, y, g) \propto \text{dnorm}(z_i, \beta^T x_i, 1) \times I_{a,b}(z_i)$$

This is simply a density of a constrained normal distribution. How to sample? Apply the inverse CDF trick that we have done previously!

Full conditional of g

Suppose the prior distribution is $p(g)$.

Given $Y = y$ and $Z = z$, then we know:

- ▶ g_k must be higher than all z_i 's for which $y_i = k$ and
- ▶ g_k must be lower than all z_i 's for which $y_i = k + 1$

Let $a_k = \max\{z_i : y_i = k\}$ and $b_k = \min\{z_i : y_i = k + 1\}$.

Then the full conditional distribution of g is then proportional to $p(g)$ but constrained to the set $\{g : a_k < g_k < b_k\}$.

Application to data

Some researchers suggest that having children reduces opportunities for educational attainment (Moore and Waite, 1977).

Here we examine this hypothesis in a sample of males in the labor force (meaning not retired, not in school and not in an institution), obtained from the 1994 General Social Survey.

For 959 of the 1,002 survey respondents we have complete data on the variables *DEG*, *CHILD* and *PDEG* described above.

We have the following variables:

- ▶ $Y_i = DEG_i$
- ▶ $x_i = (CHILD_i, PDEG_i, CHILD_i \times PDEG_i)$

Model

Our model specification is the following:

$$Y_{n \times 1} = g(Z) \quad (9)$$

$$Z_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon \quad (10)$$

$$\epsilon_{n \times 1} \stackrel{iid}{\sim} \text{Normal}(0, I) \quad (11)$$

$$\beta_{p \times 1} \sim \text{MVN}(0, n(X^T X)^{-1}) \quad (12)$$

$$p(g) \propto \prod_{k=1}^{K-1} (g_k, 0, 100)$$

Application to Data

```
X<-cbind(ychild,ypdeg,ychild*ypdeg)
y<-ydegr
# replacing missing values with the mean
keep<-(1:length(y))[ !is.na( apply( cbind(X,y),1,mean) ) ]
X<-X[keep,] ; y<-y[keep]
ranks<-match(y,sort(unique(y)))
uranks<-sort(unique(ranks))
n<-dim(X)[1] ; p<-dim(X)[2]
iXX<-solve(t(X)%*%X) ; V<-iXX*(n/(n+1)) ; cholV<-chol(V)
```

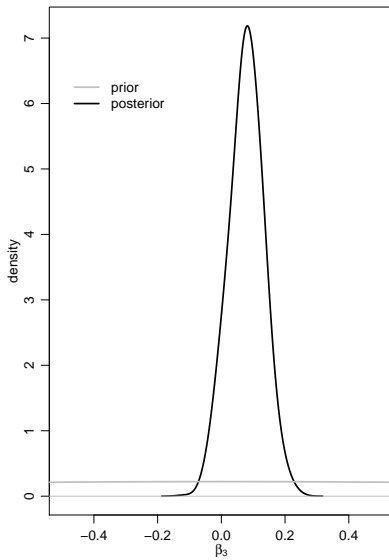
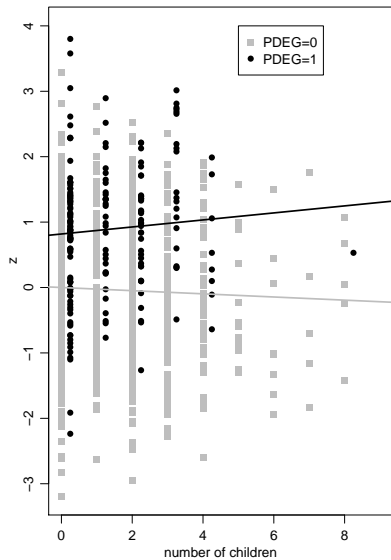
Application to Data

```
###starting values
set.seed(1)
beta<-rep(0,p)
z<-qnorm(rank(y,ties.method="random")/(n+1))
g<-rep(NA,length(uranks)-1)
K<-length(uranks)
BETA<-matrix(NA,1000,p) ; Z<-matrix(NA,1000,n) ; ac<-0
mu<-rep(0,K-1) ; sigma<-rep(100,K-1)
S<-25000
```

Gibbs sampler

```
for(s in 1:S) {  
  #update g  
  for(k in 1:(K-1)){  
    a<-max(z[y==k])  
    b<-min(z[y==k+1])  
    u<-runif(1, pnorm( (a-mu[k])/sigma[k] ),  
            pnorm( (b-mu[k])/sigma[k] ) )  
    g[k]<- mu[k] + sigma[k]*qnorm(u)  
  }  
  
  #update beta  
  E<- V%*( t(X)%*z )  
  beta<- cholV%*rnorm(p) + E  
  
  #update z  
  ez<-X%*beta  
  a<-c(-Inf,g)[ match( y-1, 0:K ) ]  
  b<-c(g,Inf)[y]  
  u<-runif(n, pnorm(a-ez),pnorm(b-ez) )  
  z<- ez + qnorm(u)
```

Plot



Commentary for Left Plot

The posterior mean regression line for people without a college-educated parent ($x_{i,2} = 0$) is

$$E[Z \mid y, x_1, x_2 = 0] = -0.024x_1$$

The posterior mean regression line for people with a college-educated parent $E[Z \mid y, x_1, x_2 = 1] = 0.818 + 0.054x_1$.

In the above figure (left), we see that for people whose parents did not go to college, the number of children they have is indeed weakly negatively associated with their educational outcome. (The opposite is true for people whose parents did go to college).

Commentary for Right Plot

Next we give the posterior distribution of β_3 along with the prior distribution for comparison.

The 95% quantile-based posterior confidence interval for β_3 is $(-0.026, 0.178)$ which contains zero but still represents a reasonable amount of evidence that the slope for the $x_2 = 1$ is larger than the $x_2 = 0$ group.