

Applied Statistics With R

Logit and Probit Models for Categorical Response Variables

John Fox

WU Wien
May/June 2006

© 2006 by John Fox

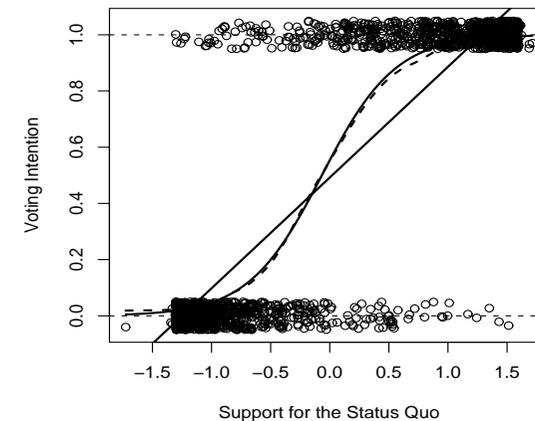
1. Goals:

- To show how models similar to linear models can be developed for qualitative/categorical response variables.
- To introduce logit (and probit) models for dichotomous response variables.
- To introduce similar statistical models for polytomous response variables, including ordered categories.
- To describe how logit models can be applied to contingency tables.

2. Models for Dichotomous Data

- To understand why special models for qualitative data are required, let us begin by examining a representative problem, attempting to apply linear regression to it:
 - In September of 1988, 15 years after the coup of 1973, the people of Chile voted in a plebiscite to decide the future of the military government. A 'yes' vote would represent eight more years of military rule; a 'no' vote would return the country to civilian government. The no side won the plebiscite, by a clear if not overwhelming margin.
 - Six months before the plebiscite, FLACSO/Chile conducted a national survey of 2,700 randomly selected Chilean voters.
 - * Of these individuals, 868 said that they were planning to vote yes, and 889 said that they were planning to vote no.
 - * Of the remainder, 558 said that they were undecided, 187 said that they planned to abstain, and 168 did not answer the question.

- * I will look only at those who expressed a preference.
- The following graph shows voting intention by support for the status-quo (high scores represent general support for the policies of the military regime).
- * The solid straight line is a linear least-squares fit; the solid curved line is a logistic-regression fit; and the broken line is a nonparametric regression fit.
- * Voting intention appears as a dummy variable, coded 1 for yes, 0 for no; the points are jittered in the plot.



- Does it make sense to think of regression as a conditional average when the response variable is dichotomous?
 - * An average between 0 and 1 represents a ‘score’ for the dummy response variable that cannot be realized by any individual.
 - * In the population, the conditional average $E(Y|x_i)$ is the proportion of 1’s among those individuals who share the value x_i for the explanatory variable — the conditional probability π_i of sampling a ‘yes’ in this group:

$$\pi_i = \Pr(Y_i) = \Pr(Y = 1|X = x_i)$$
 and thus,

$$E(Y|x_i) = \pi_i(1) + (1 - \pi_i)(0) = \pi_i$$
- If X is discrete, then in a sample we can calculate the conditional proportion for Y at each value of X .
 - * The collection of these conditional proportions represents the sample nonparametric regression of the dichotomous Y on X .

- * In the present example, X is continuous, but we can nevertheless resort to strategies such as local averaging or local regression, as illustrated in the graph.

2.1 The Linear-Probability Model

- Although non-parametric regression works here, it would be useful to capture the dependency of Y on X as a simple function, particularly when there are several explanatory variables.
- Let us first try linear regression with the usual assumptions:

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$
 where $\varepsilon_i \sim N(0, \sigma_\varepsilon^2)$, and ε_i and ε_j are independent for $i \neq j$.
 - If X is random, then we assume that it is independent of ε .
- Under this model, $E(Y_i) = \alpha + \beta X_i$, and so

$$\pi_i = \alpha + \beta X_i$$
 - For this reason, the linear-regression model applied to a dummy response variable is called the *linear probability model*.
- This model is untenable, but its failure points the way towards more adequate specifications:

- *Non-normality*: Because Y_i can take on only the values of 0 and 1, the error ε_i is dichotomous as well — not normally distributed:
 - * If $Y_i = 1$, which occurs with probability π_i , then

$$\begin{aligned} \varepsilon_i &= 1 - E(Y_i) \\ &= 1 - (\alpha + \beta X_i) \\ &= 1 - \pi_i \end{aligned}$$
 - * Alternatively, if $Y_i = 0$, which occurs with probability $1 - \pi_i$, then

$$\begin{aligned} \varepsilon_i &= 0 - E(Y_i) \\ &= 0 - (\alpha + \beta X_i) \\ &= 0 - \pi_i \\ &= -\pi_i \end{aligned}$$
 - * Because of the central-limit theorem, however, the assumption of normality is not critical to least-squares estimation of the normal-probability model.

– *Non-constant error variance*: If the assumption of linearity holds over the range of the data, then $E(\varepsilon_i) = 0$.

* Using the relations just noted,

$$\begin{aligned} V(\varepsilon_i) &= \pi_i(1 - \pi_i)^2 + (1 - \pi_i)(-\pi_i)^2 \\ &= \pi_i(1 - \pi_i) \end{aligned}$$

* The heteroscedasticity of the errors bodes ill for ordinary-least-squares estimation of the linear probability model, but only if the probabilities π_i get close to 0 or 1.

– *Nonlinearity*: Most seriously, the assumption that $E(\varepsilon_i) = 0$ — that is, the assumption of linearity — is only tenable over a limited range of X -values.

* If the range of the X 's is sufficiently broad, then the linear specification cannot confine π to the unit interval $[0, 1]$.

* It makes no sense, of course, to interpret a number outside of the unit interval as a probability.

* This difficulty is illustrated in the plot of the Chilean plebiscite data, in which the least-squares line produces fitted probabilities below 0 at low levels and above 1 at high levels of support for the status-quo.

- Dummy *regressor* variables do not cause comparable difficulties because the linear model makes no distributional assumptions about the regressors.
- Nevertheless, for values of π not too close to 0 or 1, the linear-probability model estimated by least-squares frequently provides results similar to those produced by more generally adequate methods.

2.2 Transformations of π : Logit and Probit Models

- To insure that π stays between 0 and 1, we require a positive monotone (i.e., non-decreasing) function that maps the "linear predictor" $\eta = \alpha + \beta X$ into the unit interval.

– A transformation of this type will retain the fundamentally linear structure of the model while avoiding probabilities below 0 or above 1.

– Any cumulative probability distribution function meets this requirement:

$$\pi_i = P(\eta_i) = P(\alpha + \beta X_i)$$

where the CDF $P(\cdot)$ is selected in advance, and α and β are then parameters to be estimated.

– An *a priori* reasonable $P(\cdot)$ should be both smooth and symmetric, and should approach $\pi = 0$ and $\pi = 1$ as asymptotes.

– Moreover, it is advantageous if $P(\cdot)$ is strictly increasing, permitting us to rewrite the model as

$$P^{-1}(\pi_i) = \eta_i = \alpha + \beta X_i$$

where $P^{-1}(\cdot)$ is the inverse of the CDF $P(\cdot)$.

* Thus, we have a linear model for a transformation of π , or — equivalently — a nonlinear model for π itself.

- The transformation $P(\cdot)$ is often chosen as the CDF of the unit-normal distribution

$$\Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{1}{2}Z^2} dZ$$

or, even more commonly, of the *logistic distribution*

$$\Lambda(z) = \frac{1}{1 + e^{-z}}$$

where $\pi \simeq 3.141$ and $e \simeq 2.718$ are the familiar constants.

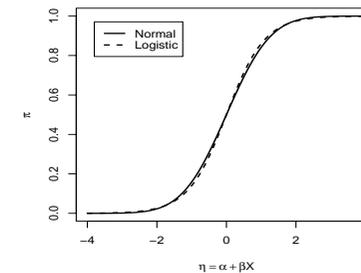
– Using the normal distribution $\Phi(\cdot)$ yields the *linear probit model*:

$$\begin{aligned}\pi_i &= \Phi(\alpha + \beta X_i) \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\alpha + \beta X_i} e^{-\frac{1}{2}Z^2} dZ\end{aligned}$$

– Using the logistic distribution $\Lambda(\cdot)$ produces the *linear logistic-regression* or *linear logit model*:

$$\begin{aligned}\pi_i &= \Lambda(\alpha + \beta X_i) \\ &= \frac{1}{1 + e^{-(\alpha + \beta X_i)}}\end{aligned}$$

– Once their variances are equated, the logit and probit transformations are very similar:



– Both functions are nearly linear between about $\pi = .2$ and $\pi = .8$. This is why the linear probability model produces results similar to the logit and probit models, except for extreme values of π_i .

• Despite their similarity, there are two practical advantages of the logit model:

1. *Simplicity*: The equation of the logistic CDF is very simple, while the normal CDF involves an unevaluated integral.
 - This difference is trivial for dichotomous data, but for polytomous data, where we will require the *multivariate* logistic or normal distribution, the disadvantage of the probit model is more acute.
2. *Interpretability*: The inverse linearizing transformation for the logit model, $\Lambda^{-1}(\pi)$, is directly interpretable as a *log-odds*, while the inverse transformation $\Phi^{-1}(\pi)$ does not have a direct interpretation.
 - Rearranging the equation for the logit model,

$$\frac{\pi_i}{1 - \pi_i} = e^{\alpha + \beta X_i}$$
 - The ratio $\pi_i/(1 - \pi_i)$ is the *odds* that $Y_i = 1$, an expression of relative chances familiar to gamblers.

– Taking the log of both sides of this equation,

$$\log_e \frac{\pi_i}{1 - \pi_i} = \alpha + \beta X_i$$

– The inverse transformation $\Lambda^{-1}(\pi) = \log_e[\pi/(1 - \pi)]$, called the *logit* of π , is therefore the log of the odds that Y is 1 rather than 0.

- The logit is symmetric around 0, and unbounded both above and below, making the logit a good candidate for the response-variable side of a linear model:

Probability π	Odds $\frac{\pi}{1-\pi}$	Logit $\log_e \frac{\pi}{1-\pi}$
.01	1/99 = 0.0101	-4.60
.05	5/95 = 0.0526	-2.94
.10	1/9 = 0.1111	-2.20
.30	3/7 = 0.4286	-0.85
.50	5/5 = 1	0.00
.70	7/3 = 2.333	0.85
.90	9/1 = 9	2.20
.95	95/5 = 19	2.94
.99	99/1 = 99	4.60

- The logit model is also a multiplicative model for the odds:

$$\frac{\pi_i}{1-\pi_i} = e^{\alpha+\beta X_i} = e^\alpha e^{\beta X_i} \\ = e^\alpha (e^\beta)^{X_i}$$

- * So, increasing X by 1 changes the logit by β and multiplies the odds by e^β .
- * For example, if $\beta = 2$, then increasing X by 1 increases the odds by a factor of $e^2 \simeq 2.718^2 = 7.389$.
- Still another way of understanding the parameter β in the logit model is to consider the slope of the relationship between π and X .

- * Since this relationship is nonlinear, the slope is not constant; the slope is $\beta\pi(1-\pi)$, and hence is at a maximum when $\pi = 1/2$, where the slope is $\beta/4$:

π	$\beta\pi(1-\pi)$
.01	$\beta \times .0099$
.05	$\beta \times .0475$
.10	$\beta \times .09$
.20	$\beta \times .16$
.50	$\beta \times .25$
.80	$\beta \times .16$
.90	$\beta \times .09$
.95	$\beta \times .0475$
.99	$\beta \times .0099$

- * The slope does not change very much between $\pi = .2$ and $\pi = .8$, reflecting the near linearity of the logistic curve in this range.

- The least-squares line fit to the Chilean plebescite data has the equation
 - This line is a poor summary of the data.
- The logistic-regression model, fit by the method of maximum-likelihood, has the equation

$$\hat{\pi}_{\text{yes}} = 0.492 + 0.394 \times \text{Status-Quo} \\ \log_e \frac{\hat{\pi}_{\text{yes}}}{\hat{\pi}_{\text{no}}} = 0.215 + 3.21 \times \text{Status-Quo}$$

- The logit model produces a much more adequate summary of the data, one that is very close to the nonparametric regression.
- Increasing support for the status-quo by one unit multiplies the odds of voting yes by $e^{3.21} = 24.8$.
- Put alternatively, the slope of the relationship between the fitted probability of voting yes and support for the status-quo at $\hat{\pi}_{\text{yes}} = .5$ is $3.21/4 = 0.80$.

2.3 An Unobserved-Variable Formulation

- An alternative derivation posits an underlying regression for a continuous but unobservable response variable ξ (representing, e.g., the 'propensity' to vote yes), scaled so that

$$Y_i = \begin{cases} 0 & \text{when } \xi_i \leq 0 \\ 1 & \text{when } \xi_i > 0 \end{cases}$$

- That is, when ξ crosses 0, the observed discrete response Y changes from 'no' to 'yes.'
- The latent variable ξ is assumed to be a linear function of the explanatory variable X and the unobservable error variable ε :

$$\xi_i = \alpha + \beta X_i - \varepsilon_i$$

- We want to estimate α and β , but cannot proceed by least-squares regression of ξ on X because the latent response variable is not directly observed.

- Using these equations,

$$\begin{aligned} \pi_i &= \Pr(Y_i = 1) = \Pr(\xi_i > 0) = \Pr(\alpha + \beta X_i - \varepsilon_i > 0) \\ &= \Pr(\varepsilon_i < \alpha + \beta X_i) \end{aligned}$$

- If the errors are independently distributed according to the unit-normal distribution, $\varepsilon_i \sim N(0, 1)$, then

$$\pi_i = \Pr(\varepsilon_i < \alpha + \beta X_i) = \Phi(\alpha + \beta X_i)$$

which is the probit model.

- Alternatively, if the ε_i follow the similar logistic distribution, then we get the logit model

$$\pi_i = \Pr(\varepsilon_i < \alpha + \beta X_i) = \Lambda(\alpha + \beta X_i)$$

- We will return to the unobserved-variable formulation when we consider models for ordinal categorical data.

2.4 Logit and Probit Models for Multiple Regression

- To generalize the logit and probit models to several explanatory variables we require a linear predictor that is a function of several regressors.
- For the logit model,

$$\begin{aligned} \pi_i &= \Lambda(\eta_i) = \Lambda(\alpha + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik}) \\ &= \frac{1}{1 + e^{-(\alpha + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik})}} \end{aligned}$$

or, equivalently,

$$\log_e \frac{\pi_i}{1 - \pi_i} = \alpha + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik}$$

- For the probit model,

$$\pi_i = \Phi(\eta_i) = \Phi(\alpha + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik})$$

- The X 's can be as general as in the general linear model, including, for example:
 - quantitative explanatory variables;
 - transformations of quantitative explanatory variables;

- polynomial regressors formed from quantitative explanatory variables;
- dummy regressors representing qualitative explanatory variables; and
- interaction regressors.

- Interpretation of the partial regression coefficients in the general logit model is similar to the interpretation of the slope in the logit simple-regression model, with the additional provision of holding other explanatory variables in the model constant.

- Expressing the model in terms of odds,

$$\begin{aligned} \frac{\pi_i}{1 - \pi_i} &= e^{(\alpha + \beta_1 X_{i1} + \dots + \beta_k X_{ik})} \\ &= e^\alpha (e^{\beta_1})^{X_{i1}} \dots (e^{\beta_k})^{X_{ik}} \end{aligned}$$

- Thus, e^{β_j} is the multiplicative effect on the odds of increasing X_j by 1, holding the other X 's constant.
- Similarly, $\beta_j/4$ is the slope of the logistic regression surface in the direction of X_j at $\pi = .5$.

- The general linear logit and probit models can be fit to data by the method of maximum likelihood.
- Hypothesis tests and confidence intervals follow from general procedures for statistical inference in maximum-likelihood estimation.
 - For an individual coefficient, it is most convenient to test the hypothesis $H_0: \beta_j = \beta_j^{(0)}$ by calculating the Wald statistic

$$Z_0 = \frac{B_j - \beta_j^{(0)}}{\text{ASE}(B_j)}$$

where $\text{ASE}(B_j)$ is the asymptotic standard error of B_j .

- * The test statistic Z_0 follows an asymptotic unit-normal distribution under the null hypothesis.

- Similarly, an asymptotic $100(1 - \alpha)$ -percent confidence interval for β_j is given by

$$\beta_j = B_j \pm z_{\alpha/2} \text{ASE}(B_j)$$

where $z_{\alpha/2}$ is the value from $Z \sim N(0, 1)$ with a probability of $\alpha/2$ to the right.

- Wald tests for several coefficients can be formulated from the estimated asymptotic variances and covariances of the coefficients.
- Wald tests in logistic regression usually behave reasonably but can sometimes be far off the mark, and so likelihood-ratio tests (and more complicated confidence intervals based on them) should generally be preferred.

- It is also possible to formulate a likelihood-ratio test for the hypothesis that several coefficients are simultaneously zero, $H_0: \beta_1 = \dots = \beta_q = 0$. We proceed, as in least-squares regression, by fitting two models to the data:

- * The full model (model 1)

$$\begin{aligned} \text{logit}(\pi) = & \alpha + \beta_1 X_1 + \dots + \beta_q X_q \\ & + \beta_{q+1} X_{q+1} + \dots + \beta_k X_k \end{aligned}$$

- * and the null model (model 0)

$$\begin{aligned} \text{logit}(\pi) = & \alpha + 0X_1 + \dots + 0X_q \\ & + \beta_{q+1} X_{q+1} + \dots + \beta_k X_k \end{aligned}$$

$$= \alpha + \beta_{q+1} X_{q+1} + \dots + \beta_k X_k$$

- * Each model produces a maximized likelihood: L_1 for the full model, L_0 for the null model.

- * Because the null model is a specialization of the full model, $L_1 \geq L_0$.
- * The generalized likelihood-ratio test statistic for the null hypothesis is

$$G_0^2 = 2(\log_e L_1 - \log_e L_0)$$

- * Under the null hypothesis, this test statistic has an asymptotic chi-square distribution with q degrees of freedom.
- A test of the omnibus null hypothesis $H_0: \beta_1 = \dots = \beta_k = 0$ is obtained by specifying a null model that includes only the constant, $\text{logit}(\pi) = \alpha$.

- An analog to the multiple-correlation coefficient can also be obtained from the log-likelihood.
 - By comparing $\log_e L_0$ for the model containing only the constant with $\log_e L_1$ for the full model, we can measure the degree to which using the explanatory variables improves the predictability of Y .
 - The quantity $G^2 = -2\log_e L$, called the *deviance* under the model, is a generalization of the residual sum of squares for a linear model.
 - Thus,

$$R^2 = 1 - \frac{G_1^2}{G_0^2}$$

$$= 1 - \frac{\log_e L_1}{\log_e L_0}$$

is analogous to R^2 for a linear model.

- Illustration based on the 1994 wave of the Statistics Canada Survey of Labour and Income Dynamics (the “SLID”): Using data on married women between 20 and 35 ($n = 1935$), I examine how the labor-force participation of these women is related to several explanatory variables (“family income” excludes the woman’s own income, if any):

Variable	Summary
Labor-Force Participation	Yes, 79 percent
Region (R)	Atlantic, 23 percent; Quebec, 13; Ontario, 30; Prairies, 26; BC, 8
Children 0–4 (K04)	Yes, 53 percent
Children 5–9 (K59)	Yes, 44 percent
Children 10–14 (K1014)	Yes, 22 percent
Family Income (I, \$1000s)	5-number summary: 0, 18.6, 26.7, 35.1, 131.1
Education (E, years)	5-number summary: 0, 12, 13, 15, 20

- Allowing for the possibility of interaction between presence of children and each of family income and education in determining women’s labor-force participation, the following models are formulated so that likelihood-ratio tests of terms in the full model can be computed by taking differences in the residual deviances for the models, in conformity with the principle of marginality:

Model	Terms in the Model	Number of Parameters	Residual Deviance
0	C	1	1988.084
1	C, R, K04, K59, K1014, I, E, K04×I, K59×I, K1014×I, K04×E, K59×E, K1014×E	16	1807.376
2	Model 1 – K04×I	15	1807.378
3	Model 1 – K59×I	15	1808.600
4	Model 1 – K1014×I	15	1807.834
5	Model 1 – K04×E	15	1807.407
6	Model 1 – K59×E	15	1807.734
7	Model 1 – K1014×E	15	1807.938
8	Model 1 – R	12	1824.681
9	C, R, K04, K59, K1014, I, E, K59×I, K1014×I, K59×E, K1014×E	14	1807.408

Model	Terms in the Model	Number of Parameters	Residual Deviance
10	Model 9 – K04	13	1866.689
11	C, R, K04, K59, K1014, I, E, K04×I, K1014×I, K04×E, K1014×E	14	1809.268
12	Model 11 – K59	13	1819.273
13	C, R, K04, K59, K1014, I, E, K04×I, K59×I, K04×E, K59×E	14	1808.310
14	Model 13 – K1014	13	1808.548
15	C, R, K04, K59, K1014, I, E, K04×E, K59×E, K1014×E	13	1808.854
16	Model 15 – I	12	1817.995
17	C, R, K04, K59, K1014, I, E, K04×I, K59×I, K1014×I	13	1808.428
18	Model 17 – E	12	1889.223

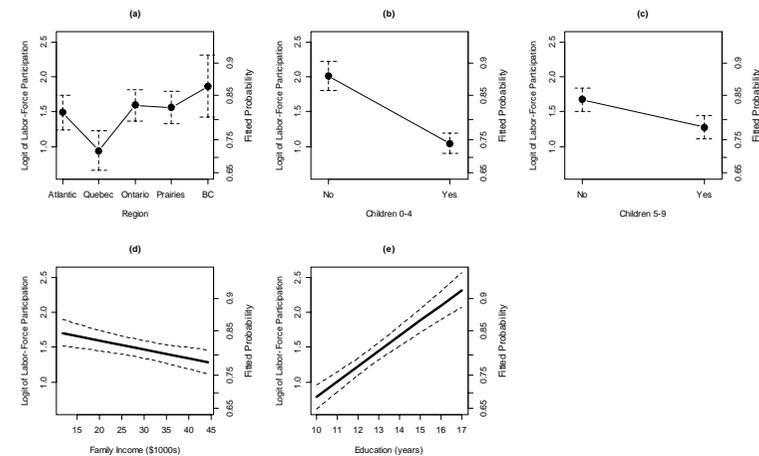
- Likelihood-ratio tests (in a “Type-II” analysis of deviance table):

Term	Models			
	Contrasted	df	G_0^2	p
Region (R)	8-1	4	17.305	.0017
Children 0–4 (K04)	10-9	1	59.281	≪ .0001
Children 5–9 (K59)	12-11	1	10.005	.0016
Children 10–14 (K1014)	14-12	1	0.238	.63
Family Income (I)	16-15	1	9.141	.0025
Education (E)	18-17	1	80.795	≪ .0001
K04×I	2-1	1	0.002	.97
K59×I	3-1	1	1.224	.29
K1014×I	4-1	1	0.458	.50
K04×E	5-1	1	0.031	.86
K59×E	6-1	1	0.358	.55
K1014×E	7-1	1	0.562	.45

- Coefficients for a final model fit to the data:

Coefficient	Estimate (B_j)	Standard Error	e^{B_j}
Constant	-0.3763	0.3398	
Region: Quebec	-0.5469	0.1899	0.579
Region: Ontario	0.1038	0.1670	1.109
Region: Prairies	0.0742	0.1695	1.077
Region: BC	0.3760	0.2577	1.456
Children 0–4	-0.9702	0.1254	0.379
Children 5–9	-0.3971	0.1187	0.672
Family Income (\$1000s)	-0.0127	0.0041	0.987
Education (years)	0.2197	0.0250	1.246
Residual Deviance	1810.444		

- Effect plots for the fitted model (setting other terms to typical values):



3. Models for Polytomous Data

- I will describe three general approaches to modeling polytomous data:
 1. Modeling the polytomy directly as a set of unordered categories, using a generalization of the dichotomous logit model.
 2. Constructing a set of nested dichotomies from the polytomy, fitting an independent logit or probit model to each dichotomy.
 3. Extending the unobserved-variable interpretation of the dichotomous logit and probit models to ordered polytomies.

3.1 The Polytomous Logit Model

- The dichotomous logit model can be extended to a polytomy by employing the multivariate-logistic distribution. This approach has the advantage of treating the categories of the polytomy in a non-arbitrary, symmetric manner.
- The response variable Y can take on any of m qualitative values, which, for convenience, we number $1, 2, \dots, m$ (using the numbers only as category labels).
 - For example, a married woman can (1) work full-time, (2) work part-time, or (3) not work outside of the home.
- Let π_{ij} denote the probability that the i th observation falls in the j th category of the response variable; that is,

$$\pi_{ij} \equiv \Pr(Y_i = j) \text{ for } j = 1, \dots, m.$$
- We have k regressors, X_1, \dots, X_k , on which the π_{ij} depend.

- More specifically, suppose that this dependence can be modeled using the *multivariate logistic distribution*:

$$\pi_{ij} = \frac{e^{\gamma_{0j} + \gamma_{1j}X_{i1} + \dots + \gamma_{kj}X_{ik}}}{1 + \sum_{l=1}^{m-1} e^{\gamma_{0l} + \gamma_{1l}X_{i1} + \dots + \gamma_{kl}X_{ik}}}$$

for $j = 1, \dots, m - 1$

$$\pi_{im} = 1 - \sum_{l=1}^{m-1} \pi_{il}$$

- There is one set of parameters, $\gamma_{0j}, \gamma_{1j}, \dots, \gamma_{kj}$, for each response-variable category but the last; category m functions as a type of baseline.
- The use of a baseline category is one way of avoiding redundant parameters because of the restriction that $\sum_{j=1}^m \pi_{ij} = 1$.

- Some algebraic manipulation of the model produces

$$\log_e \frac{\pi_{ij}}{\pi_{im}} = \gamma_{0j} + \gamma_{1j}X_{i1} + \dots + \gamma_{kj}X_{ik}$$

for $j = 1, \dots, m - 1$

- The regression coefficients affect the log-odds of membership in category j versus the baseline category.
- It is also possible to form the log-odds of membership in *any* pair of categories j and j' :

$$\begin{aligned} \log_e \frac{\pi_{ij}}{\pi_{ij'}} &= \log_e \left(\frac{\pi_{ij}}{\pi_{im}} \bigg/ \frac{\pi_{ij'}}{\pi_{im}} \right) \\ &= \log_e \frac{\pi_{ij}}{\pi_{im}} - \log_e \frac{\pi_{ij'}}{\pi_{im}} \\ &= (\gamma_{0j} - \gamma_{0j'}) + (\gamma_{1j} - \gamma_{1j'})X_{i1} \\ &\quad + \dots + (\gamma_{kj} - \gamma_{kj'})X_{ik} \end{aligned}$$

- * The regression coefficients for the logit between any pair of categories are the differences between corresponding coefficients.

- Now suppose that the model is specialized to a dichotomous response variable. Then, $m = 2$, and

$$\begin{aligned}\log_e \frac{\pi_{i1}}{\pi_{i2}} &= \log_e \frac{\pi_{i1}}{1 - \pi_{i1}} \\ &= \gamma_{01} + \gamma_{11}X_{i1} + \dots + \gamma_{k1}X_{ik}\end{aligned}$$

- Applied to a dichotomy, the polytomous logit model is identical to the dichotomous logit model.

- Example adapted from work by Andersen, Heath, and Sinnott on the 2001 British election:
 - Central issue: the potential interaction between respondents' political knowledge and political attitudes in determining vote.
 - The response variable, vote, has three categories: Labour, Conservative, and Liberal Democrat.
 - There are several explanatory variables:
 - * Attitude toward European integration, an 11-point scale, with high scores representing a negative attitude (so-called "Euro-sceptism").
 - * Knowledge of the platforms of the three parties on the issue of European integration, with integer scores ranging from 0 through 3. (Labour and the Liberal Democrats supported European integration, the Conservatives were opposed.)
 - * Other variables included in the model primarily as "controls"—age, gender, perceptions of national and household economic conditions, and ratings of the three party leaders.

- Estimates:

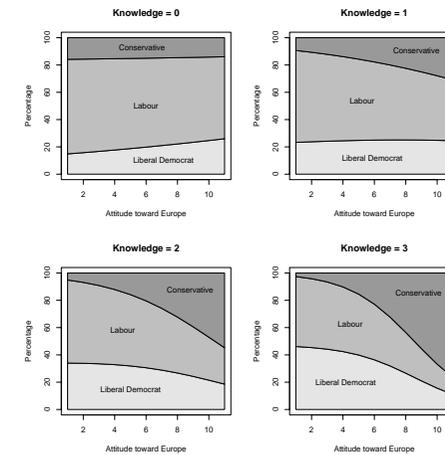
Coefficient	Labour/Lib Dem	
	Estimate	SE
Constant	-0.155	0.612
Age	-0.005	0.005
Gender (male)	0.021	0.144
Perception of Economy	0.377	0.091
Perception of Household Economic Position	0.171	0.082
Evaluation of Blair (Labour leader)	0.546	0.071
Evaluation of Hague (Conservative leader)	-0.088	0.064
Evaluation of Kennedy (Liberal Democrat leader)	-0.416	0.072
Attitude Toward European Integration	-0.070	0.040
Political Knowledge	-0.502	0.155
Europe × Knowledge	0.024	0.021

Coefficient	Cons/Lib Dem	
	Estimate	SE
Constant	0.718	0.734
Age	0.015	0.006
Gender (male)	-0.091	0.178
Perception of Economy	-0.145	0.110
Perception of Household Economic Position	-0.008	0.101
Evaluation of Blair (Labour leader)	-0.278	0.079
Evaluation of Hague (Conservative leader)	0.781	0.079
Evaluation of Kennedy (Liberal Democrat leader)	-0.656	0.086
Attitude Toward European Integration	-0.068	0.049
Political Knowledge	-1.160	0.219
Europe × Knowledge	0.183	0.028

• Analysis of deviance table:

Source	df	G ₀ ²	p
Age	2	13.87	.0009
Gender	2	0.45	.78
Perception of Economy	2	30.60	≪ .0001
Perception of Household Economic Position	2	5.65	.059
Evaluation of Blair	2	135.37	≪ .0001
Evaluation of Hague	2	166.77	≪ .0001
Evaluation of Kennedy	2	68.88	≪ .0001
Attitude Toward European Integration	2	78.03	≪ .0001
Political Knowledge	2	55.57	≪ .0001
Europe × Knowledge	2	50.80	≪ .0001

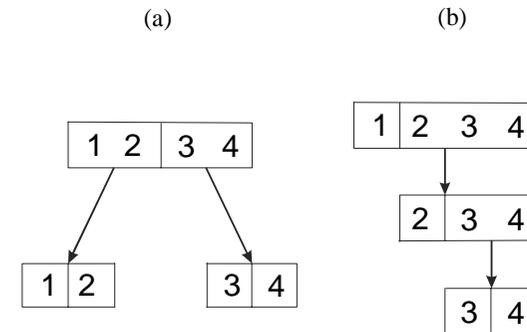
• Effect display for the interaction between attitude and knowledge:



3.2 Nested Dichotomies

- Perhaps the simplest approach to polytomous data is to fit separate models to each of a set of dichotomies derived from the polytomy.
 - These dichotomies are *nested*, making the models statistically independent.
 - Logit models fit to a set of nested dichotomies constitute a model for the polytomy, but are not equivalent to the polytomous logit model previously described.
- A nested set of $m - 1$ dichotomies is produced from an m -category polytomy by successive binary partitions of the categories of the polytomy.

- Two examples for a four-category variable:
 - * In (a), the dichotomies are {1, 2}, {3, 4}, and {1, 2, 3, 4}.
 - * In (b), the nested dichotomies are {1, 2, 3, 4}, {2, 3, 4}, and {3, 4}.



- Because the results of the analysis and their interpretation depend upon the set of nested dichotomies that is selected, this approach to polytomous data is reasonable only when a particular choice of dichotomies is substantively compelling.
- Nested dichotomies are attractive when the categories of the polytomy represent ordered progress through the stages of a process.
 - Imagine that the categories in (b) represent adults' attained level of education: (1) less than high school; (2) high-school graduate; (3) some post-secondary; (4) post-secondary degree.
 - Since individuals normally progress through these categories in sequence, the dichotomy {1, 234} represents the completion of high school; {2, 34} the continuation to post-secondary education, conditional on high-school graduation; and {3, 4} the completion of a degree conditional on undertaking a post-secondary education.

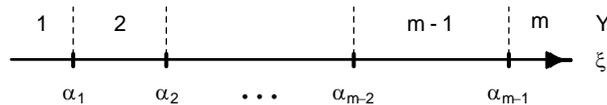
3.3 Ordered Logit and Probit Models

- Imagine that there is a latent variable ξ that is a linear function of the X 's plus a random error:

$$\xi_i = \alpha + \beta_1 X_{i1} + \dots + \beta_k X_{ik} + \varepsilon_i$$
 - Suppose that instead of dividing the range of ξ into two regions to produce a dichotomous response, the range of ξ is dissected by $m - 1$ thresholds into m regions.
 - Denoting the thresholds by $\alpha_1 < \alpha_2 < \dots < \alpha_{m-1}$, and the resulting response by Y , we observe

$$Y_i = \begin{cases} 1 & \text{if } \xi_i \leq \alpha_1 \\ 2 & \text{if } \alpha_1 < \xi_i \leq \alpha_2 \\ \cdot & \\ \cdot & \\ m - 1 & \text{if } \alpha_{m-2} < \xi_i \leq \alpha_{m-1} \\ m & \text{if } \alpha_{m-1} < \xi_i \end{cases}$$

- The thresholds, regions, and corresponding values of ξ and Y are represented graphically as follows:



- Using the model for the latent variable, along with category thresholds, we can determine the cumulative probability distribution of Y :

$$\begin{aligned} \Pr(Y_i \leq j) &= \Pr(\xi_i \leq \alpha_j) \\ &= \Pr(\alpha + \beta_1 X_{i1} + \dots + \beta_k X_{ik} + \varepsilon_i \leq \alpha_j) \\ &= \Pr(\varepsilon_i \leq \alpha_j - \alpha - \beta_1 X_{i1} - \dots - \beta_k X_{ik}) \end{aligned}$$

- If the errors ε_i are independently distributed according to the standard normal distribution, then we obtain the ordered probit model.
- If the errors follow the similar logistic distribution, then we get the ordered logit model:

$$\begin{aligned} \text{logit}[\Pr(Y_i \leq j)] &= \log_e \frac{\Pr(Y_i \leq j)}{\Pr(Y_i > j)} \\ &= \alpha_j - \alpha - \beta_1 X_{i1} - \dots - \beta_k X_{ik} \end{aligned}$$

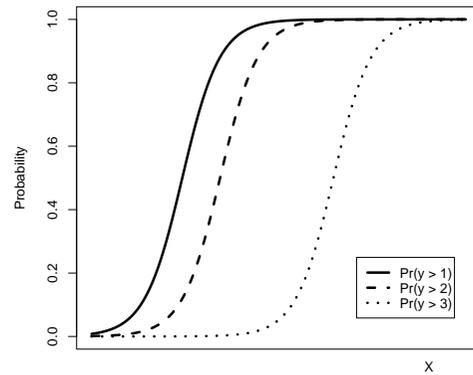
- Equivalently,

$$\begin{aligned} \text{logit}[\Pr(Y_i > j)] &= \log_e \frac{\Pr(Y_i > j)}{\Pr(Y_i \leq j)} \\ &= (\alpha - \alpha_j) + \beta_1 X_{i1} + \dots + \beta_k X_{ik} \end{aligned}$$

for $j = 1, 2, \dots, m - 1$.

- The logits in this model are for cumulative categories — at each point contrasting categories above category j with category j and below.
- The slopes for each of these regression equations are identical; the equations differ only in their intercepts.

* The logistic regression surfaces are therefore horizontally parallel to each other, as illustrated for $m = 4$ response categories and a single X :

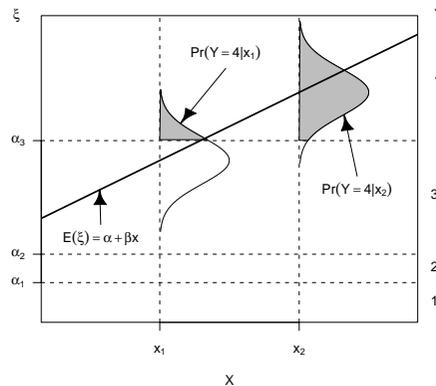


– For a fixed set of X 's, any two different cumulative log-odds — say, at categories j and j' — differ only by the constant $(\alpha_j - \alpha_{j'})$.
 * The odds, therefore, are proportional to one-another, and for this reason, the ordered logit model is called the *proportional-odds model*.

- There are $(k + 1) + (m - 1) = k + m$ parameters to estimate in the proportional-odds model, including the regression coefficients $\alpha, \beta_1, \dots, \beta_k$ and the category thresholds $\alpha_1, \dots, \alpha_{m-1}$.
 - There is an extra parameter in the regression equations, since each equation has its own constant, $-\alpha_j$, along with the common constant α .
 - A simple solution is to set $\alpha = 0$ (and to absorb the negative sign in α_j), producing

$$\text{logit}[\text{Pr}(Y_i > j)] = \alpha_j + \beta_1 X_{i1} + \dots + \beta_k X_{ik}$$

• The following graph illustrates the proportional-odds model for $m = 4$ response categories and a single X :



- Example: Data from the World Values Survey (WVS) of 1995–97.
 - To provide a manageable example, I will restrict attention to four countries: Australia, Sweden, Norway, and the United States.
 - The combined sample size for these four countries is 5381.
 - The response variable in the analysis is the answer to the question, “Do you think that what the government is doing for people in poverty is about the right amount, too much, or too little.”
 - There are several explanatory variables:
 - * gender (a dummy variable coded 1 for *men* and 0 for *women*).
 - * whether or not the respondent belonged to a religion (coded 1 for *yes*, 0 for *no*).
 - * whether or not the respondent had a university degree (coded 1 for *yes* and 0 for *no*).
 - * age (in years, ranging from 18 to 87). Preliminary analysis of the data suggested a roughly linear age effect.

* country (a set of three dummy regressors, with *Australia* as the base-line category).

– Analysis of deviance table for an initial model:

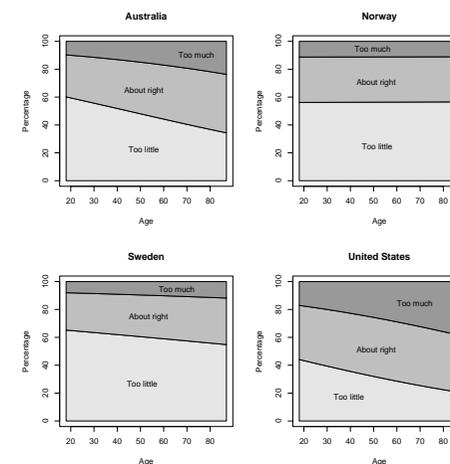
Source	df	G_0^2	p
Country	3	250.881	≪ .0001
Gender	1	10.749	.0010
Religion	1	4.132	.042
Education	1	4.284	.038
Age	1	49.950	≪ .0001
Country×Gender	3	3.049	.38
Country×Religion	3	21.143	< .0001
Country×Education	3	12.861	.0049
Country×Age	3	17.529	.0005

– Estimates for a final model:

Coefficient	Estimate	Standard Error
Gender (Men)	0.1744	0.0532
Country (Norway)	0.1516	0.3355
Country (Sweden)	−1.2237	0.5821
Country (United States)	1.2225	0.3068
Religion (Yes)	−0.0255	0.1120
Education (Degree)	−0.1282	0.1676
Age	0.0153	0.0026

Coefficient	Estimate	Standard Error
Country (Norway)×Religion	0.2456	0.2153
Country (Sweden)×Religion	0.9031	0.5125
Country (United States)×Religion	−0.5706	0.1733
Country (Norway)×Education	0.0524	0.2080
Country (Sweden)×Education	0.6359	0.2141
Country (United States)×Education	0.3103	0.2063
Country (Norway)×Age	−0.0156	0.0044
Country (Sweden)×Age	−0.0090	0.0047
Country (United States)×Age	0.0008	0.0040
Thresholds		
$-\hat{\alpha}_1$ (Too Little About Right)	0.7189	0.1953
$-\hat{\alpha}_2$ (About Right Too Much)	2.5372	0.1986

– Effect display for the age × country interaction:



– Testing the assumption of proportional odds:

<i>Model</i>	<i>Residual Deviance</i>	<i>Number of Parameters</i>
Proportional-Odds Model	10,350.12	18
Cumulative Logits, Unconstrained Slopes	9,961.63	34
Polytomous Logit Model	9,961.26	34

* Likelihood-ratio statistic for testing the assumption of proportional odds:

$$G_0^2 = 10,350.12 - 9,961.63 = 388.49$$

on $34 - 18 = 16$ degrees of freedom.

* This test statistic is highly statistically significant, leading us to reject the proportional-odds assumption for these data.

3.4 Comparison of the Three Approaches

- The three approaches to modeling polytomous data — the polytomous logit model, logit models for nested dichotomies, and the proportional-odds model — address different sets of log-odds, corresponding to different dichotomies constructed from the polytomy.
- Consider, for example, the ordered polytomy {1, 2, 3, 4}:
 - Treating category 1 as the baseline, the coefficients of the polytomous logit model apply directly to the dichotomies {1, 2}, {1, 3}, and {1,4}, and indirectly to any pair of categories.
 - Forming continuation dichotomies (one of several possibilities), the nested-dichotomies approach models {1, 234}, {2, 34}, and {3, 4}.
 - The proportional-odds model applies to the dichotomies {1, 234}, {12, 34}, and {123, 4}, imposing the restriction that only the intercepts of the three regression equations differ.

- Which of these models is most appropriate depends partly on the structure of the data and partly upon our interest in them.

4. Discrete Explanatory Variables and Contingency Tables

- When the explanatory variables — as well as the response variable — are discrete, the joint sample distribution of the variables defines a contingency table of counts.

- An example, drawn from *The American Voter* (Converse et al., 1960), appears below.
 - This table, based on data from a sample survey conducted after the 1956 U.S. presidential election, relates voting turnout in the election to strength of partisan preference, and perceived closeness of the election:

Perceived Closeness	Intensity of Preference	Turnout	
		Voted	Did Not Vote
One-Sided	Weak	91	39
	Medium	121	49
	Strong	64	24
Close	Weak	214	87
	Medium	284	76
	Strong	201	25

- The following table gives the *empirical logit* for the response variable, $\log_e \frac{\text{proportion voting}}{\text{proportion not voting}}$ for each of the six combinations of categories of the explanatory variables:

Perceived Closeness	Intensity of Preference	$\log_e \frac{\text{Voted}}{\text{Did Not Vote}}$
One-Sided	Weak	0.847
	Medium	0.904
	Strong	0.981
Close	Weak	0.900
	Medium	1.318
	Strong	2.084

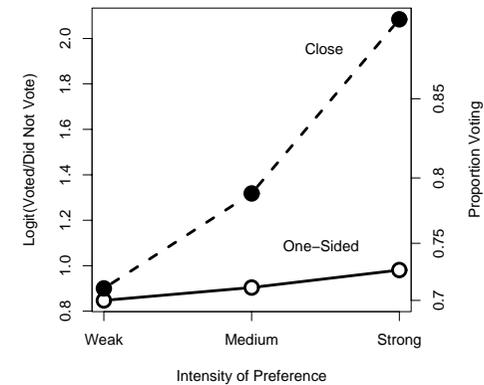
* For example,

$$\begin{aligned} \text{logit}(\text{voted}|\text{one-sided, weak preference}) &= \log_e \frac{91/130}{39/130} \\ &= \log_e \frac{91}{39} \\ &= 0.847 \end{aligned}$$

* Because the conditional proportions voting and not voting share the same denominator, the empirical logit can also be written as

$$\log_e \frac{\text{number voting}}{\text{number not voting}}$$

* Graph of empirical logits:



- Logit models are fully appropriate for tabular data.
 - When, as in the example, the explanatory variables are qualitative or ordinal, it is natural to use logit or probit models that are analogous to analysis-of-variance models.
 - Treating perceived closeness of the election as the ‘row’ explanatory variable and intensity of partisan preference as the ‘column’ explanatory variable, for example, yields the model

$$\text{logit } \pi_{jk} = \mu + \alpha_j + \beta_k + \gamma_{jk}$$

where

- * π_{jk} is the conditional probability of voting in combination of categories j of perceived closeness and k of preference;
- * μ is the general level of turnout in the population;
- * α_j is the main effect on turnout of membership in the j th category of perceived closeness;
- * β_k is the main effect on turnout of membership in the k th category of preference; and

- * γ_{jk} is the interaction effect on turnout of simultaneous membership in categories j of perceived closeness and k of preference.
- Under the usual sigma constraints, this model leads to deviation-coded regressors, as in the analysis of variance.

- Deviances under several models for the *American-Voter* data:

Model	$k + 1$	Deviance G^2
α, β, γ	6	1356.434
α, β	4	1363.552
α, γ	4	1368.042
β, γ	5	1368.554
α	2	1382.658
β	3	1371.838

- An analysis-of-deviance table showing alternative Type-II and Type-III tests for the main effects:

Source	df	G_0^2	p
Perceived Closeness	1		
$\alpha \beta$ (Type II)		8.286	.0040
$\alpha \beta, \gamma$ (Type III)		12.120	.0005
Intensity of Preference	2		
$\beta \alpha$ (Type II)		19.106	<.0001
$\beta \alpha, \gamma$ (Type III)		11.608	.0030
Closeness \times Preference	2		
$\gamma \alpha, \beta$		7.118	.028

- The log-likelihood-ratio statistic for testing

$$H_0: \text{all } \gamma_{jk} = 0$$

for example, is

$$\begin{aligned} G_0^2(\gamma|\alpha, \beta) &= G^2(\alpha, \beta) - G^2(\alpha, \beta, \gamma) \\ &= 1363.552 - 1356.434 \\ &= 7.118 \end{aligned}$$

with $6 - 4 = 2$ degrees of freedom, for which $p = .03$.

5. Summary

- It is problematic to apply least-squares linear regression to a dichotomous response variable:
 - The errors cannot be normally distributed and cannot have constant variance.
 - Even more fundamentally, the linear specification does not confine the probability for the response to the unit interval.
- More adequate specifications transform the linear predictor $\eta_i = \alpha + \beta_1 X_{i1} + \dots + \beta_k X_{ik}$ smoothly to the unit interval, using a cumulative probability distribution function $P(\cdot)$.
 - Two such specifications are the probit and the logit models, which use the normal and logistic CDFs, respectively.

- Although these models are very similar, the logit model is simpler to interpret, since it can be written as a linear model for the log-odds:

$$\log_e \frac{\pi_i}{1 - \pi_i} = \alpha + \beta_1 X_{i1} + \dots + \beta_k X_{ik}$$

- The dichotomous logit model can be fit to data by the method of maximum likelihood.
 - Wald tests and likelihood-ratio tests for the coefficients of the model parallel t -tests and F -tests for the general linear model.
 - The deviance for the model, defined as $G^2 = -2 \times$ the maximized log-likelihood, is analogous to the residual sum of squares for a linear model.

- Several approaches can be taken to modeling polytomous data, including:
 - (a) modeling the polytomy directly using a logit model based on the multivariate logistic distribution;
 - (b) constructing a set of $m - 1$ nested dichotomies to represent the m categories of the polytomy; and
 - (c) fitting the proportional-odds model to a polytomous response variable with ordered categories.
- When all of the variables — explanatory as well as response — are discrete, their joint distribution defines a contingency table of frequency counts.
 - It is natural to employ logit models that are analogous to analysis-of-variance models to analyze contingency tables.