

**Viewing One-Year Retention as a Continuum: the Use of Dichotomous Logistic Regression,  
Ordered Logit and Multinomial Logit**

Paper presented at the Association of Institutional Research 1999 annual meeting,  
Seattle, Washington, May 30 – June 2, 1999

Stephen R. Porter, Ph.D.  
Office of Institutional Studies  
2119 Main Administration Building  
University of Maryland, College Park  
College Park, MD 20742

Phone: (301) 405-5608  
Fax: (301) 314-9443  
Email: [sporter@accmail.umd.edu](mailto:sporter@accmail.umd.edu)  
Homepage: <http://www.wam.umd.edu/~srporter>

## **Abstract**

Studies of one-year retention commonly use a binary outcome (retained after one year, yes/no) as the dependent variable. This paper examines the use of alternate specifications of one-year retention that include information about spring semester stopout behavior. Different types of discrete choice models are first described. Next, the multinomial logit model using four outcomes (enrolled first fall semester only; enrolled first fall semester and spring semester; enrolled first fall semester, stopped out and returned the following fall; and enrolled all three semesters) is compared with the traditional dichotomous logit approach.

The pseudo R-square indicates the multinomial model fits the data better, but for practical reasons it has very poor predictive power. The multinomial model does yield interesting results for the impact of independent variables on retention: for example, increases in high school grade point average impacts the probability of a student being in the fall-spring outcome more than the probability of falling in the first fall semester only outcome.

## **Introduction**

My motivation for this analysis was to examine alternate ways of analyzing one-year retention while using recent cohorts of data as well as retaining information about spring semester stopout behavior. Survival analysis is generally superior than logistic regression because it can take into account the semester-to-semester behavior of students over time (for more information see DesJardins et. al. 1997, Ronco 1996, Willett and Singer 1991). However, survival analysis has two drawbacks. First, data is generally needed over several time periods. This means that cohort data from several years ago must be used; unfortunately, most administrators are interested in analyses based on the most recent available data. Second, many institutional researchers are unfamiliar with survival analysis and would have to learn an entirely new statistical technique in order to employ it efficiently.

Other formulations of the logistic regression approach offer one solution to this problem. These models allow the researcher to take into account spring stopout behavior while also using the most recent data available. In addition, many institutional researchers are familiar with dichotomous logistic regression. Learning to use other formulations of these models and interpreting their coefficients should not be too difficult for most researchers.

Use of these models allows the researcher to retain information about student behavior that would otherwise be lost. This is important since a general rule of thumb in any statistical analysis is that information should be used and not thrown away whenever possible. The need for this in studies of one-year retention can be seen in Table 1. The 3,589 first-time full-time degree-seeking freshmen who enrolled in the Fall 1995 semester at the University of Maryland have been divided into four groups depending on their registration behavior over the next year:

- enrolled Fall 1995 only
- enrolled Fall 1995 and Spring 1996 but did not enroll in the following fall semester
- enrolled Fall 1995, stopped out Spring 1996 and returned in Fall 1996
- enrolled all three semesters

Dichotomous logistic regression combines these four groups into two groups, as can be seen in the far right column. The purpose of this paper is to see what advantage is gained by keeping these four groups of students separate during an analysis of one-year retention.

## **Discrete Choice Models**

Discrete choice models are a class of maximum likelihood techniques that allow researchers to model behavior where the outcome, or dependent variable, is discrete rather than

**Table 1. One-Year Retention Outcomes**

Group	Distribution		Registered?			Dichotomous logit measure
	N	%	Fall 95	Spring 96	Fall 96	
A	153	4.3	Yes	No	No	0
B	321	8.9	Yes	Yes	No	0
C	33	.9	Yes	No	Yes	1
D	3,082	85.9	Yes	Yes	Yes	1
Total	3,589	100.0				

continuous. *Logistic regression* is used when the dependent variable has only two outcomes. In institutional research the most familiar example would be one-year retention, where the outcome is coded as retained or not retained (e.g. Nora et. al. 1996, St. John et. al. 1996). There are, however, other types of these models that can be used to analyze educational behavior (much of the following discussion is taken from chapter 19 of Greene (1997). Although his textbook is very technical the chapter on discrete choice models has a very clear narrative and is a must-read for anyone working with these techniques).

*Ordered logit* models are used when the dependent variable has more than two discrete outcomes, and these outcomes can be ranked in some fashion. Bond ratings are the common example in economics research, while in the field of education opinion surveys would be another. In this approach we assume that one outcome can be ranked above another, but we know nothing about the distance between outcomes. For example, in an opinion survey there may be three responses such as “very satisfied”, “somewhat satisfied”, and “not satisfied at all”. We know the first response can be ranked above the second in terms of satisfaction, and the second response ranked above the third, but we cannot be sure that the distance between the first and second responses is equal to the distance between the second and third. Multiple regression makes this assumption of common distance, rendering it theoretically unsuitable for such data<sup>1</sup>.

<sup>1</sup> Of course, in practice there may not be much difference between multiple regression and ordered logit for many applications.

Finally, there are two techniques that allow analysis of dependent variables with more than two discrete outcomes than cannot be ranked in any meaningful way. The technique used depends on the data being analyzed. In economics information about choices are very common. For example, analyses of commuter choice behavior will use datasets in which information varies over the commuting choices of bus, car or train. This information may take the form of cost of the commuting choice per mile, or the time of commute for each choice. These models are known as *conditional logit* models; because educational data over choices is uncommon they will not be dealt with here.

The other technique is known as *multinomial logit* and is used when data over individuals is analyzed. Using the commuter example, we may only have data such as income, education and occupation of the individual commuter available (as well as their commute choice). Examples of this technique in the field of education include work by Keil and Partell (1999), Ordovensky (1995) and Weiler (1987, 1989).

I initially thought ordered logit could be used to analyze one-year retention by creating a ranked outcome according to the number of semesters attended by the student. This involved combining categories B and C in the first column of Table 1. Unfortunately this model did not perform as well as either the dichotomous or multinomial models. I decided that analyzing retention as a ranked outcome did not make sense theoretically and abandoned this approach. The remainder of the paper deals with dichotomous and multinomial logistic models.<sup>2</sup>

## **Retention model**

The model estimated is a standard retention model using many of the variables in the retention literature, including demographics, human capital, safety nets, commitment and financial resources:

### *Demographics*

- Age (in years).
- Nonwhite – coded 1 if the student was a minority or international student, 0 otherwise.
- Female – coded 1 if female, 0 otherwise.

---

<sup>2</sup> There is some confusion as to whether a logit or probit distribution should be used when estimating these models. The two distributions are almost exactly alike, and in most applications will yield similar results. They do differ, however, in the tails of the distribution. Analyses with extreme data values may thus yield different results for the two distributions (Green 1997, p. ?).

- MD resident – residency based on tuition status, coded 1 if Maryland state resident, 0 otherwise.

#### *Human capital*

- SATM and SATV – highest math and verbal Scholastic Aptitude Test scores submitted by the student.
- HS GPA - high-school grade point average.
- First gen – taken from the student’s application, coded 1 if student indicated s/he was first in family to attend college, 0 otherwise.

#### *Safety nets*

- Honors – coded 1 if student participated in the university Honors program, 0 otherwise.
- On campus – measures whether the student resided on campus their first semester, coded 1 if so, 0 otherwise.

#### *Commitment*

- App time – number of days between the first day of classes and the date of the student’s application.
- Undecided – coded 1 if the student did not declare a major upon entry, 0 otherwise.

#### *Financial resources*

- Unmet need – the amount of money needed by the student to cover costs of attending the university. Positive amounts indicate need, negative amounts indicate no need. Students who did not apply for financial aid have missing data for this variable; they are assumed to have zero unmet need and are coded 0.

### **Dichotomous logit results**

The estimates for the dichotomous logit model are presented in Table 2. Variables significant at the .05 level are shaded. The binary dependent variable used is the coding in the far right column of Table 1. The results are similar to other studies in this area. Ignoring demographics (which were included purely for control reasons), students with better high school grade point averages are retained at higher rates, as are students participating in the Honors program and students living on campus. Application time, included as a proxy for commitment, is also positively related to being retained after one year. Students who apply far in advance are more likely to be retained than students who apply later. Unmet need is also related to retention. As a student’s financial resources diminish (their unmet need increases) they become less likely to enroll the following fall semester.

**Table 2. Dichotomous Logit Results**

Variable	Coefficient	Std. Error	z	P> z
Age	-0.04471	0.04516	-0.99	0.32
Nonwhite	0.30656	0.12481	2.46	0.01
Female	0.10406	0.11202	0.93	0.35
MD resident	0.01279	0.11554	0.11	0.91
SATV	-0.00750	0.00717	-1.05	0.30
SATM	0.01269	0.00671	1.89	0.06
HS GPA	1.05400	0.13959	7.55	0.00
Honors	0.61345	0.25516	2.40	0.02
On campus	0.32532	0.13020	2.50	0.01
App time	0.00547	0.00106	5.14	0.00
First gen	-0.08837	0.12239	-0.72	0.47
Undecided	-0.01283	0.10835	-0.12	0.91
Unmet need	-0.00003	0.00001	-2.84	0.01
Constant	-2.79934	1.10962	-2.52	0.01
N	3589			
Log likelihood	-1298.8			
LR	203.8			
Pseudo R <sup>2</sup>	0.0728			

Table 3 presents the results for the same retention model, but the dependent variable used is the coding in the first column of Table 1; in other words, one-year retention is coded as four outcomes rather than two.

At first glance the results appear confusing: instead of one set of coefficients there are three. This results from the nature of the dependent variable. In the binary case the coefficients are usually estimated in the form of measuring the impact of an independent variable on the probability of the yes outcome versus the no outcome. The multinomial case is exactly the same: the coefficients measure the impact of an independent variable on the probability of *one* outcome versus a *base* outcome. Since there are four outcomes and one outcome is treated as the base (or “excluded”) outcome, the result is three sets of coefficients. In Table 3 the excluded outcome is enrolled for all three semesters, so results are given for the three outcomes labeled A, B and C in Table 1. Note that changes in probability remain the same no matter which outcome is excluded; however, the coefficients themselves will change depending on the excluded category.<sup>3</sup>

<sup>3</sup> The probabilities do not change because different formulas are used for different outcomes depending on which outcome is excluded. See Green (1997) p. 875.

**Table 3. Multinomial Logit Results**

Outcome	Variable	Coefficient	Std. Error	z	P> z
Fall only	Age	0.05986	0.05655	1.06	0.29
	Nonwhite	-0.50378	0.21175	-2.38	0.02
	Female	0.02961	0.18589	0.16	0.87
	MD resident	-0.44479	0.19453	-2.29	0.02
	SATV	0.01944	0.01204	1.62	0.11
	SATM	-0.02388	0.01132	-2.11	0.04
	HS GPA	-0.94819	0.23244	-4.08	0.00
	Honors	-1.23987	0.61124	-2.03	0.04
	On campus	-0.75844	0.21533	-3.52	0.00
	App time	-0.00827	0.00168	-4.92	0.00
	First gen	-0.17826	0.20882	-0.85	0.39
	Undecided	0.21796	0.17627	1.24	0.22
	Unmet need	0.00006	0.00002	3.23	0.00
	Constant	2.36594	1.58331	1.49	0.14
Fall - Spring	Age	0.02474	0.05926	0.42	0.68
	Nonwhite	-0.20865	0.14581	-1.43	0.15
	Female	-0.17237	0.13250	-1.30	0.19
	MD resident	0.18730	0.13710	1.37	0.17
	SATV	0.00213	0.00841	0.25	0.80
	SATM	-0.00912	0.00785	-1.16	0.25
	HS GPA	-1.15495	0.16387	-7.05	0.00
	Honors	-0.42117	0.28039	-1.50	0.13
	On campus	-0.10065	0.15408	-0.65	0.51
	App time	-0.00467	0.00128	-3.66	0.00
	First gen	0.20631	0.14192	1.45	0.15
	Undecided	-0.08718	0.12926	-0.67	0.50
	Unmet need	0.00002	0.00001	1.55	0.12
	Constant	2.65279	1.39748	1.90	0.06
Fall - Fall	Age	-0.02981	0.21774	-0.14	0.89
	Nonwhite	0.39741	0.44116	0.90	0.37
	Female	0.04786	0.39356	0.12	0.90
	MD resident	0.33910	0.39615	0.86	0.39
	SATV	0.00389	0.02597	0.15	0.88
	SATM	-0.02842	0.02288	-1.24	0.21
	HS GPA	-1.53521	0.47819	-3.21	0.00
	Honors	-0.37367	1.08721	-0.34	0.73
	On campus	1.10623	0.56992	1.94	0.05
	App time	-0.01353	0.00317	-4.27	0.00
	First gen	-0.11304	0.41191	-0.27	0.78
	Undecided	-0.01122	0.37360	-0.03	0.98
	Unmet need	0.00010	0.00005	2.18	0.03
	Constant	3.99951	4.68092	0.85	0.39
N		3589			
Log likelihood		-1728.3			
LR		306.7			
Pseudo R <sup>2</sup>		0.0815			

## Which approach is “better”?

We need some sort of criteria to decide between the two approaches to modeling retention. I believe two criteria are relevant: predictive ability and explanatory power. Predictive ability is simply the ability of the model to correctly predict the outcomes of the dependent variable. Explanatory power, on the other hand, has a different connotation in the context of this paper. Explanatory power simply refers to what the model tells us about student behavior (*not* “what percentage of the variance is explained.”). Are students who live on campus during their first semester more likely to return to the university after a year? Models that can answer these types of questions can be said to have good explanatory power. Obviously explanatory power cannot be measured directly like predictive ability and is more of a judgement call.

The distinction between the two criteria is important because models can have high predictive power and little explanatory power, and vice versa. A simple example makes this clear. Suppose two analysts estimate dichotomous logit models on a dataset where the overall retention rate is 80%. The first analyst uses a typical group of variables such as SAT scores, etc., while the second uses only a constant (or alternatively, a coin flip for every observation, where ‘heads’ assigns a one for that observation, a zero otherwise).

Next, an evaluation committee examines the models to determine which one should be used for policy-making purposes. They discover that the standard retention model correctly predicts student retention outcomes only 45% percent of the time, while the constant model predicts correct outcomes 80% of the time (this follows from the construction of the model, because all students are predicted to be retained and 80% actually are retained). The committee rejects the first model and decides to use the second model for their decision-making because of its superior predictive ability. They ask the second analyst, “What does your model tell us about student behavior?” The answer, of course, is nothing, because the model consists only of a constant. The first model, although a poor predictor of retention, nonetheless can offer interesting information about the impact of various variables on student behavior. This example illustrates the difficulty in relying on predictive power for these types of models, because one can easily develop highly predictive models with little explanatory power.

### *Predictive ability*

From the pseudo R-squares at the bottom of Tables 2 and 3 we can conclude that the multinomial model appears to fit the data better than the dichotomous model.<sup>4</sup> However, if some type of intervention system for at-risk students is under consideration, the real measure of predictive ability is the proportion of outcomes correctly predicted. An institution does not want to waste intervention resources on student who are likely to stay, and they also do not want to miss applying the intervention to those at-risk students who are likely to stop out. Here the multinomial model fails miserably, because the sample used is what Greene (1997, p. 892) terms “unbalanced”. An unbalanced sample has cases that are not evenly distributed across outcomes. This poses a problem because the base probability for an outcome for every individual will be the relative frequency of that outcome. If the relative frequency is very high or low, then only an extraordinary number of regressors could cause the predicted probability of this outcome to shift above or below the predicted probabilities of the other outcomes.

Because of the unbalanced sample, predicting outcomes in the multinomial model is difficult. Like the dichotomous case, a predicted probability for each individual student and each outcome can be derived from the model coefficients<sup>5</sup>. We can use two different decision rules for predicting outcomes based on these probabilities. First, the outcome with the highest predicted probability can be declared the predicted outcome. Unfortunately with this sample every student is predicted to be enrolled all three semesters, because the predicted probability for this outcome is always in the 70%-90% range, much larger than all the other outcomes. Second, we can compare the predicted probability of each outcome with the actual relative frequency for each outcome. For example, if the predicted probability for a student for enrolling in the first fall semester only is 5%, this student is assigned this outcome because 5% is greater than the actual relative frequency (or sample mean) of 4.3%. Unfortunately for most students in the sample *two* outcomes are predicted using this decision rule. That is, one outcome has a reduced probability, and since the sum of the probabilities for the four outcomes must sum to 1, this probability is often shifted to two other outcomes rather than just one. The result is ambiguous predictions for most individuals in the sample. Unfortunately the multinomial approach does not seem useful for actually predicting student outcomes; however, in a more balanced sample the multinomial approach might prove superior to dichotomous logit.

---

<sup>4</sup> The pseudo R-square is calculated as  $1 - (\log \text{likelihood of the full model} / \log \text{likelihood of a model estimated with only a constant})$  and is bounded from zero to one. It is more formally known as the likelihood ratio index (Greene 1997, p. 891).

<sup>5</sup> Porter (1999) shows how to calculate these and other simulation probabilities for the dichotomous logit case.

### *Explanatory power*

What the model tells us about student behavior is the second criteria by which to judge the two approaches. The impacts of five of the statistically significant variables from the dichotomous logit model on the probability of enrolling in the second fall semester are presented in Table 4 (this probability for the multinomial case is calculated by summing the probability of enrolling the first fall semester, stopping out in spring and enrolling the following fall semester and the probability of enrolling all three semesters). Changes in probability were calculated as follows. The predicted probability of enrolling in the second fall semester was calculated using the sample means for all independent variables except the variable for which the change is calculated. That variable is constrained to the value indicated. The process was repeated using the second value of the independent variable and the difference between the two probabilities was taken. For example, the impact of housing on retention was estimated by calculating the predicted probability with the on campus variable set to zero rather than the sample mean; this was repeated with on campus set to one and the difference taken.

**Table 4. Change in Probability of Being Enrolled in Fall 1996 for Selected Independent Variables, DL and ML Models**

	Dichotomous logit	Multinomial logit
High school GPA changes from 3.0 to 4.0	8.4%	7.7%
Enrolled in Honors program	5.1%	3.9%
Resided on campus	3.4%	2.1%
Applied 12 months instead of 6 months before	11.5%	9.7%
Unmet need changes from \$20,000 to zero	6.8%	5.9%

Note: all probabilities evaluated at the sample means.

The numbers for the two models are similar, as expected. The one major difference is that the multinomial changes are all slightly smaller than the dichotomous logit changes. This is probably due to the fact that the multinomial model slightly overestimates the probability that a student will enroll all three semesters. This type of table is useful for double-checking calculations, because the two models should give similar results.

Table 4 gives the typical results to be expected from a dichotomous logit model. Can the multinomial model improve on this? Table 5 presents similar calculations for the multinomial model from Table 3, with changes calculated for each of the four possible enrollment outcomes. The differences between the four outcomes are quite interesting.

**Table 5. Change in Probability of Four Retention Outcomes for Selected Independent Variables, Multinomial Model**

	Fall only	Fall-Spring	Fall-Fall	Fall-Spring-Fall
High school GPA = 3.0	1.9%	10.1%	0.3%	87.7%
High school GPA = 4.0	0.8%	3.5%	0.1%	95.7%
<i>Difference</i>	<b>-1.1%</b>	<b>-6.7%</b>	<b>-0.2%</b>	<b>8.0%</b>
Not enrolled in Honors program	1.8%	8.4%	0.2%	89.5%
Enrolled in Honors program	0.5%	5.8%	0.2%	93.5%
<i>Difference</i>	<b>-1.3%</b>	<b>-2.7%</b>	<b>-0.1%</b>	<b>4.0%</b>
Resided off campus	2.7%	8.5%	0.1%	88.7%
Resided on campus	1.3%	7.9%	0.3%	90.5%
<i>Difference</i>	<b>-1.4%</b>	<b>-0.7%</b>	<b>0.2%</b>	<b>1.9%</b>
Applied 6 months before	3.4%	12.1%	0.9%	83.6%
Applied 12 months before	0.7%	5.1%	0.1%	94.1%
<i>Difference</i>	<b>-2.7%</b>	<b>-7.0%</b>	<b>-0.8%</b>	<b>10.5%</b>
Unmet need = \$20,000	5.3%	10.4%	1.7%	82.6%
Unmet need is zero	1.6%	8.1%	0.3%	90.0%
<i>Difference</i>	<b>-3.7%</b>	<b>-2.2%</b>	<b>-1.5%</b>	<b>7.4%</b>

Note: all probabilities evaluated at the sample means.

For example, Table 4 shows that the impact of a change in high school grade point average from 3.0 to 4.0 is about 8 percentage points; that is, such an increase would decrease the probability of stopping out by 8%. The results in Table 5 show exactly where this occurs. If a hypothetical student experienced such a change, the impact on stopping out is in the Fall 95 – Spring 96 category, rather than the Fall 95 only category. One explanation for this result may be the university’s academic dismissal policy. Students who perform poorly face dismissal, and the minimum grade point average needed to avoid dismissal rises with the number of credits earned. In essence this means that it is very difficult to be academically dismissed after the first semester since so few credits have been earned, but it is much less difficult after two semesters worth of credit have been earned. Given that students with 3.0 grade point averages are more likely to do poorly at the college level and are more likely to face academic dismissal, and that it is difficult to be academically dismissed after the first semester, the impact of a change in high school grade point average is greatest in the Fall 95 – Spring 96 category.

Some of the other changes are more difficult to interpret. It is not clear, for example, why changes in application time (our proxy for commitment) should have differential impact on the first two enrollment outcomes, Fall 95 only and Fall 95 – Spring 96.

## **Conclusion**

The multinomial approach allows the use of the most recent student cohort data while taking into account information about spring stopout behavior. The multinomial model fared poorly in predicted actual student outcomes, but this was due to the unbalanced nature of the sample analyzed in this study. It is likely that institutions with more balanced samples (i.e., more stopout behavior) would have better results with the multinomial approach. The multinomial model, however, does provide a more nuanced view of one-year retention as seen in Table 5.

## References

DesJardins, Stephen L., Dennis A. Ahlburg and Brian P. McCall. "Using Event History Methods to Model the Different Modes of Student Departure from College." Paper presented at the 1997 Association of Institutional Research meeting, Orlando, FL.

Frazis, Harley. 1993. "Selection Bias and the Degree Effect." *Journal of Human Resources*, vol. 28 pp. 538-554.

Greene, William H. 1997. *Econometric Analysis*. New York: MacMillan Publishing Company.

Keil, Jack and Peter J. Partell. 1999. "Is Supplemental Instruction by Graduate Teaching Assistants Beneficial to Undergraduates?" Paper presented at the Association of Institutional Research 1999 annual meeting, Seattle, Washington.

Mason, Paul M. et. al. 1995. "Student Evaluations of Faculty: A New Procedure for Using Aggregate Measures of Performance." *Economics of Education Review*, vol. 14, pp. 403-416.

Nora, Amaury, Alberto Cabrera, Linda Serra Hagedorn and Ernest Pascarella. 1996. "Differential Impacts of Academic and Social Experiences on College-Related Behavioral Outcomes Across Different Ethnic and Gender Groups at Four-Year Institutions." *Research in Higher Education*, vol. 37 pp. 427-451.

Ordozensky, J. Farley. 1995. "Effects of Institutional Attributes on Enrollment Choice: Implications for Postsecondary Vocational Education." *Economics of Education Review*, vol. 14, pp. 335-350.

Porter, Stephen R. 1999. "Understanding and Interpreting Dichotomous Logistic Regression Coefficients." Manuscript, University of Maryland.

Ronco, Sharron L. 1996. "How Enrollment Ends: Analyzing the Correlates of Student Graduation, Transfer and Dropout with a Competing Risks Model." *AIR Professional File*, number 61.

St. John, Edward P., Michael B. Paulsen and Johnny B. Starkey. 1996. "The Nexus Between College Choice and Persistence." *Research in Higher Education*, vol. 37 pp. 175-220.

Weiler, William C. "An Application of the Nested Multinomial Logit Model to Enrollment." *Research in Higher Education*, vol. 27 pp. 273-282.

Weiler, William C. 1989. "A Flexible Approach to Modelling Enrollment Choice Behavior." *Economics of Education Review*, vol. 8 pp. 277-283.

Willett, John B. and Judith D. Singer. 1991. "From Whether to When: New Methods for Studying Student Dropout and Teacher Attrition." *Review of Educational Research*, vol. 61 pp. 407-450.