



for Management Research, Policy Analysis, and Planning

AIR Professional File

Building a Student Flow Model

William A. Simpson
Professor

Office of Planning and Budgets
Michigan State University

The purpose of this paper is to introduce the researcher to the usefulness of a certain type of student flow model and to indicate how it may be constructed. Since another article by this author, to appear in the *Journal of Higher Education*, has been devoted entirely to describing the model and how it has been used, this paper focuses primarily on facilitating the development of similar student flow models by indicating the nature of the data needed and identifying useful software packages to assist with the programming.

Conceptual Origins

The original model began as a computer program that followed students through their studies at the university, keeping track of all their declared majors and tabulating them in such a way as to make the data readable and useful. Paul Dressel and this author originally viewed the model as a simple tool for studying the motivations behind student attrition. The concept underlying the model was that nearly every major change represents a student judgment made about the program she or he is both leaving and entering. It was felt that, in many instances, inchoate student dissatisfaction with programs, relevance, instruction, and advising was the proximate cause of a major change; thus, if one could but organize these major changes in a compact and sufficiently clever manner, there was a possibility that some patterns might arise from the details. From this humble beginning, the model quickly evolved into a more sophisticated mechanism replete with many ancillary printouts and options. The rapid refinement was due to a combination of heavy usage, which led to suggestions for further improvements, and a simplicity of concept that made modifications easy to carry out.

The Model

The model began with the premise that an unusually high or low flow of students into or out of a specific program, via changes in major, has something to say about that program—favorable or unfavorable. The task thus set was to tabulate somehow the many major changes made by a student cohort, which in a large university may exceed five thousand. In addition, the display of this data had to be structured so as to make apparent any patterns underlying the data. In effect, the total had to be far more than the sum of its insignificant parts—counts of changes in major.

The original computer model was designed to record all changes in major made by students belonging to a specific cohort. The model then produced two outputs—one displaying all these major changes and a second that tabulated only the first and last major of each student. Each of these methodologies provides a different insight into student preferences, with the relative value of each being dependent upon the nature of the question under investigation. In view of this, the research team decided that the model should be capable of performing both analyses, with the user being given the option of selecting the one desired. The underlying algorithms of the two approaches are nearly identical, and so, from the programming standpoint, this decision was correct.

However, in a paper of this nature some compromises must be made to reduce its length and lessen the complexities of exposition. Because the option that deals with all major changes requires considerably more data and a more sophisticated data base, only the first-to-last methodology will be discussed. This abridgement was made more palatable by the realization that over the last

two years the first-to-last major methodology has proven to be far more useful and that in many situations the two models provided very similar results. The author feels quite confident that if an institution has limited resources and/or a rather sparse student data base, the proper choice for a trial model would be the first-to-last major methodology. Again, the interested reader is referred to the longer paper (W.A. Simpson, *Journal of Higher Education*, in press) for a full discussion of the pros and cons of both methodologies.

If one views a student's change of major as a transaction revealing something useful about both the old major (department) and the new one, it becomes clear that the transaction must be recorded in such a way that both components are discernable. If the change of major is viewed as a flow of students from a former major to a new one, this way of thinking is very suggestive of other situations where a data matrix has proven to be a useful display device. Soon after one begins thinking in terms of a matrix display, it becomes apparent that such a structure meets all the requirements.

The First-to-Last Model. For purposes of explanation, the simple three-department matrix shown as Table 1 is sufficient.

Table 1

Example of Output for Tracking First-to-Last Major

Department of First Major	Department of Last Major			Total
	English	History	Philosophy	
English	3	2	1	6
History	2	4	1	7
Philosophy	0	6	2	8
Total	5	12	4	21

Each entry in the transition matrix can be viewed as the intersection point of a row with a column. Thus, the circled "1" is seen as one student changing from an original major in the English department (row 1) to a final major in philosophy (column 3). Note that a student who makes *no* major change is still captured by this methodology; for example, a student retaining an original choice of a major in English would appear within the cell denoting an English $\xrightarrow{\text{to}}$ English transition.

Table 2

Persistence/Attrition Matrix
for First-Time Freshmen Entering Fall 1978

1978 Initial Major		1984 Final Major		
		Physics	Mathematics	
Physics	Yet enrolled	2	5	7 120 55 182 final distribution of all students entering physics as freshmen
	Graduated	100	20	
	Not enrolled	50	5	
	Total	152	30	
Mathematics	Yet enrolled	1	10	11 202 52 265
	Graduated	2	200	
	Not enrolled	2	50	
	Total	5	260	
Final Major as of 1984	Yet enrolled	3	15	18 322 107 447 final distribution of all students (the standard "attrition" study data)
	Graduated	102	220	
	Not enrolled	52	55	
	Total	157	290	

final distribution of all students electing physics as their final major

Further, a row provides a complete disclosure of all major changes *out* of the associated department; for example, row 2 tells us that two students changed from history to English, four remained in history, and one went from history to philosophy. The row total indicates that seven students out of the total cohort of 21 (the lower right corner total) began in the History Department.

In a similar fashion, each column indicates the nature of the flow of students *into* the associated department. Thus, the second column shows that the History Department gains two and six students who began in English and philosophy, respectively. This, combined with the four students who began and ended in history, makes a total of 12 students who were history majors at the end. In our interpretations of this table, we can equate major changes with students, i.e., every student makes one, and only one, major change from (first major) $\xrightarrow{\text{to}}$ (last major).

The Departmental-Level Persistence Model. By taking the flow-of-majors model one step further, we can obtain a model that produces attrition/persistence reports at the department level. If the rows of the matrix in Table 1 (using the first-to-last major methodology) are further subdivided to indicate the final status of the students involved (i.e., "yet enrolled," "graduated," "not enrolled"), the matrix provides a wealth of persistence information. Table 2 represents such a hypothetical example for a two-department university. Note that the lower right-hand corner block contains the data usually associated with a total university persistence/attrition report. Aside from using the data presented in the matrix to calculate various attrition percentages or indices, the output of this model is interpreted in the same way as that in Table 1. In fact, once this attrition/persistence model is developed, there is no need to preserve the first-to-last major model as a separate option since all the data in the latter is contained in the more detailed output of the persistence model. Development of the more powerful persistence model requires so little additional data and programming effort that most researchers would be well served to set their sights on this version. Hereafter, the generic term "student flow model" will signify both the first-to-last major model and the persistence version.

Output Formats

In producing the major change matrix for a large university, one must come to terms with the fact that a readable 80-department by 80-department matrix is unwieldy. The printout may consist of sixteen pages which, if pasted together, would create a 2' x 4' matrix (even larger in the case of the persistence report). The problem cannot be solved by making a reduced copy because the cell entries become unreadable before a convenient size is reached. Those who use this matrix frequently soon become proficient at turning quickly to the correct page(s), but they cannot expect other users to accommodate to this. It is clear that a potentially useful tool might be disregarded by administrators if the outputs were too much of a nuisance to read.

Recognizing this, a second version of the output was created, consisting of a package of fifteen matrices. In each of these fifteen matrices, most of the transition

data is accumulated at the college level rather than the department level, i.e., major changes are shown as switches from one college or another college. However, within each matrix, one specific college has been expanded into a more detailed department-by-department matrix. For example, a planner, interested in major changes within the Agriculture College, would turn to the matrix shown in Figure 1, which contains the detailed *departmental* data for the Agriculture College and *college-level* aggregations (only) for all major changes occurring outside that college. This compromise facilitated the reduction of the large university matrix to a set of two-page matrices and yet preserved the departmental level of detail within the college of concern.

While attention is being directed toward Figure 1, it may be useful to illustrate how the characteristics, which were explained using Tables 1 and 2, carry over to this far more detailed format. Figure 1 shows the flow of majors across departments within the Agriculture College (the upper left-hand block of data). For example, the line headed "Hort" shows that 41 students declared a horticulture major upon entering the university (see the line total on the far right). Of this number, 26 retained the horticulture major and 29 remained in the Agriculture College. Shifting attention to the "Hort" column, we read from the bottom entry that 71 students were finally horticulture majors, 44 of whom began in the Agriculture College. If an administrator reviewing this particular sheet wanted to know more about where the three students leaving horticulture went, the more detailed departmental matrix could be consulted.

Uses of the Model

The model was first used to answer specific questions about certain departments. For example, a few years ago, the Packaging Department suddenly experienced a many-fold increase in majors via upper-division major changes. An analysis of the major flow printout disclosed that the surge in enrollments was due to overflow from the Business and Engineering Colleges, both of which had placed enrollment restrictions on their programs. Apparently the Packaging Program offered the mix of technology and business applications that made it attractive to both groups of students who had been refused entry in the limited programs. The knowledge that packaging enrollments were closely tied to those in engineering and business was useful in guiding later policy decisions. Subsequently, it was decided that since the increase in majors in the Packaging Program was closely related to the overflow from engineering and business, it was likely to be of short duration and that temporary faculty rather than tenure-stream faculty should be hired to accommodate the overload.

Another study involving the flow model led to the discovery that students who enter the university with no major in mind make fewer changes in major than do those students entering with a declared major. This surprising result threw into question the efficacy of encouraging lower-division students to declare a major, and it may eventually result in policy changes.

At the university level, there customarily is only one persistence ratio (%) used. However, when we look at persistence at the department level via the persistence version of the student flow model, there are several ratios (percentages) that provide slightly different in-

Table 3
Three Types of Persistence Ratios

1. persistence within the packaging major	=	$\frac{\text{original majors graduating or still enrolled in PKG}}{\text{total original packaging majors}}$	=	$\frac{8+7}{27}$	=	56%
2. persistence within the university by original packaging majors	=	$\frac{\text{original majors graduating or still enrolled in university}}{\text{total original packaging majors}}$	=	$\frac{10+10}{27}$	=	74%
3. persistence in packaging by all final majors	=	$\frac{\text{all final PKG majors graduating or enrolled in PKG}}{\text{all students with PKG as their last major}}$	=	$\frac{160+251}{476}$	=	86%

sights. Three such ratios, as calculated for the packaging (PKG) major, are shown in Table 3.

The first percentage is concerned with the persistence of the original majors as contrasted to the third percentage which shows the persistence of all students ending up in the department. In this example, the second ratio is not particularly illuminating; however, once the tracking period is extended to six years so that "persistence" becomes "persistence to a bachelor's degree," then the second ratio becomes more significant for it enables one to gain a feel for the quality of students attracted to the university by the department. If, for example, a department has values of 10% and 15% for percentages (1) and (2) respectively, a planner concerned about the low persistence rate must first consider the possibility that the department attracts students with marginal academic capabilities who are largely unable to earn a degree in any program. If, however, these same persistence ratios are 10% and 80%, then one must seek an explanation elsewhere since the original majors were obviously very able students; they simply did not find their original choice of major to be what they wanted or expected.

Lately, the persistence/attrition model has been used to develop persistence trend data for engineering and business major cohorts as well as to discover where the engineering and business majors tended to go when not admitted into those respective programs at the junior level, where enrollment restrictions are set.

Most recently, the model has been used in an attempt to predict shifts in student interests. For example, we have hypothesized that shifts may occur within the student body currently enrolled before similar shifts can be observed within the incoming freshmen cohort. Already there are signs of a growing interest in education as a major, as evidenced by a recent, significant increase in students changing their majors to education. No such trend has yet developed with new students. If this conjecture is borne out by events, we may, with the help of the student flow model, be able to gain a year or two lead time on large-scale shifts in student demand.

The Data Base

The data used by the more detailed persistence version of the student flow model can vary from a very few elements to quite a refined data base. This variance is due entirely to how specific a student cohort must be selected. If one is satisfied with choosing cohorts no more complex than "all students entering the institution as freshmen during a specified year," the only data

needed is a code indicating entry status (transfer or first-time freshman) and the entry date. However, far more data is needed if one wishes to select cohorts such as "all black women who entered the university in 1978 as transfer students in engineering but later switched majors into a science discipline," an example that is by no means too contrived. Regardless of the nature of the cohort selected, the data required to construct the output matrix is the same, and very minimal, consisting of first major, last major, and the students' status as of the end of the tracking period (graduated, yet enrolled, not enrolled).

The usefulness of the model increases with the ability to select more constrained cohorts. With this in mind, the researcher should endeavor to build a model that can use as much of the available student data as possible, even if the usage of the model, as initially contemplated, might require only a small fraction of the stored data. In time, the model's latent potential will surely be appreciated and used.

Despite the foregoing encouragement towards constructing the most flexible possible model, one must realize that any model of this type, no matter how rudimentary, is better than no model at all. Even the most simplistic student flow model will prove to be a valuable tool. The most basic version is well within the data support capabilities of nearly all institutions, as it requires only the following data for each student: entry status, entry date, last major, and first major. By adding a fifth data element—the student's status—the departmental persistence version of the model can be created. Of these data elements, the one most likely to be absent from the data base is the student's first major. The replacement of the student's former major with the most recent one anytime a change occurs seems to be a common practice; the correction of this situation would involve augmenting the data base to include every student's first declared major.

To avoid entering the university student data base every time we wished to run the student flow model, we created a scaled-down version of the official data base. Since the Statistical Analysis System (SAS) was used to construct the student flow model (more on this point later), selected data elements were copied from the university data base and placed in a SAS data file. This file has easily supported all of our uses of the student flow model up to this point and clearly has the capability of supporting the next generation of model refinements. Even though the SAS data base is extremely large, containing 1.2 million records for each of the 130,000 students enrolled at the university from 1973 to the

present, it is a simple sequential file with none of the sophisticated linkages that one expects to see under today's standards. Every time the model is run, the program must read the entire data file (arranged by student number) from first to last entry, extracting data pertaining to each student whose characteristics fit those of the specified cohort. This mode of operation, for such a large file, is clearly inefficient from a theoretical point of view, but it has not been a problem with respect to cost or operational time. Certainly, a smaller institution would find this straightforward approach completely adequate.

Our data base currently has support capabilities that outstrip present usage, and yet it contains but a moderate number of data elements and few that would not be found in nearly all student data bases. A data collection consisting of the following elements would be capable of supporting a very refined version of the departmental persistence model:

1. First and last majors
2. Status code pertaining to each academic year from entry to last
3. Record (graduated, yet enrolled, not enrolled)
4. Gender code
5. Ethnic code
6. Entry code (transfer, first-time freshman, special student)
7. Resident status (in-state, out-state, foreign)
8. Standardized test scores (SAT, ACT)
9. High school grade point average (GPA)
10. University orientation/placement test scores
11. Last recorded university GPA.

(The addition of *every* major declared by each student would enable one to augment the student flow model with an option to tabulate *all* major changes.)

Developing the Model

The programming needed to build the student flow model can be greatly reduced by using one of the standard data manipulation software programs included within SPSS or SAS. We used SAS to construct both the student data base and the model. The ease with which the program can be modified has led to many refinements being made to the model.

The model is quite simple in concept, consisting of two well-defined steps. The user sets a series of codes that specify the characteristics of the student cohort desired. The model then searches the data base from start to end, looking for students with the cohort characteristics. Once such a student has been found, the program extracts three data elements—status as of the end of the tracking interval, first major, last major—and stores them in a temporary file. When this phase has been completed, the model records the cohort data in a frequency count table that reports the number of students changing from major M_1 to major M_2 , from M_1 to M_3 , . . . , from M_{99} to M_{100} . Both SPSS and SAS have subroutines that compile, format, and fully label the columns and rows of such frequency tables. The format displayed in Table 3 is a copy of one such table as produced by SAS. (SAS has two subroutines that can be used: PROC TABULATE and PROC SUMMARY.)

Since many questions about student flow require several runs of the model—perhaps using different cohorts or different tracking time intervals—it is clear that the model had to be easy to use, which is to say that it had

to be easy for the user to select a different cohort or tracking period. This was accomplished by using a standard device, the setting of a series of parameter switches at the beginning of the program. An illustration of a specification sheet for the switches that must be set for each run of the model appears as Figure 2. With but a little practice, a user quickly progresses to the point where ten or more runs of the model can be made in five minutes.

Currently there is a SAS package designed for use with a microcomputer. This enables a researcher to build the student flow model on a small desk-top computer as long as the student data base is not too large. There are also alternative ways to use a microcomputer in conjunction with a mainframe computer which enable a user to surmount the storage problem presented by a large student data base.

COHORT SPECIFICATION INSTRUCTIONS	
IN EACH BLANK SPACE PUT NO MORE THAN ONE ENTRY. EVERY VARIABLE DOWN TO 12A AND 13A MUST HAVE AN ENTRY.	
1. WHAT REPORT DO YOU WANT?	(REPORT = 0 FOR ATTRITION STUDY, 1 FOR MAJOR FLOW, 2 FOR BOTH MODELS)
2. ENTRY STATUS	(ENTRY = 1 IF 1ST TIME FRESHMAN, 2 IF TRANSFER UNDERGRAD, 3 IF GRADUATE STUDENT, 4 IF ALL ENTERING STUDENTS)
3. GENDER	(GENDER = 1 IF F, 2 IF M, 3 IF BOTH)
4A. COHORT ENTRY INTERVAL	(INTERVAL = 1 IF FALL TERM, 2 IF AY (F,W,S), 3 IF AN (F,W,S), 4 IF MULTIPLE YEARS)
B. END YEAR OF COHORT	(IF INTERVAL = 4 THEN ENDY = LAST TWO DIGITS OF LAST YEAR OF THE COHORT INTERVAL, IF INTERVAL NE 4 THEN ENDY = 99)
C. BEGINNING YEAR OF COHORT	(BEGINY = LAST TWO DIGITS OF FIRST YEAR OF THE COHORT INTERVAL)
D. STATUS AT END OF TRACKING PERIOD	(STATUS = 0 IF NOT ENROLL, 1 IF ENROLLED BUT NOT GRADUATED, 2 IF GRADUATED, 3 IF ALL)
6. TRACKING PERIOD: STARTING YEAR	(STARTY = LAST TWO DIGITS OF STARTING YEAR)
	LENGTH OF STUDY (LENGTH = # YEARS)
7. RESIDENT STATUS	(RESID = 0 IF MICH., 1 IF U.S. (NOT MICH.), 2 IF FOREIGN, 3 IF ALL U.S., 4 IF ALL)
8. HIGH SCHOOL GPA:	ALL X GE (HSOPL = X.XX DECIMAL PT. IMPLIED, 0 IF IGNORED) ALL X LE (HSOPL = X.XX DECIMAL PT. IMPLIED, 0 IF IGNORED)
9. LAST MSU GPA:	ALL X GE (MSUGPL = X.XX FORMAT, 0 IF IGNORED) ALL X LE (MSUGPL = X.XX FORMAT, 0 IF IGNORED)
10. HIGH SCHOOL ATTENDED	(HS = SEVEN DIGIT H.S. CODE, 9999999 IF IGNORED)
11. ETHNIC STATUS	(ETHNIC = 0 IF ALL (I.E. IGNORE), 1 IF WHITE, 2 IF BLACK, . . . , 6 IF ASIAN/PACIFIC, 7 IF OTHER)
12A. RESTRICT COHORT TO STUDENTS ENTERING A CERTAIN UNIT	(EUNIT = 1 IF DEPT, 2 IF A COLLEGE, 3 IF IGNORE)
B. IF RESTRICTION 12A APPLIES, SPECIFY THE UNIT	(ECURR = 'XX' AND EMAJ = XX IF EUNIT = 1 ECURR = 'XX' AND EMAJ = 0 IF EUNIT = 2 ECURR = ' ' AND EMAJ = 0 IF EUNIT = 3)

Figure 2. An illustration of cohort specification parameters.

Later Refinements

Once a researcher has built a first student flow model and used it for a period of time, he or she will know whether or not the programming of additional features is a desirable project. If the initial model is strictly the first-to-last version, then the next step would quite likely be the expansion into the departmental persistence/attrition model. An investment consisting of the addition of one more data element to the base—a student's status for each year from entry to last record—and a small amount of new programming would make the transition into the more powerful version. At this point, if the student data base has been a rather minimal one, to the extent that the usefulness of the model in per-

forming certain studies has been limited, the researcher may wish to focus his or her efforts on upgrading the data base. By adding only a few more student biographical data elements, the model's capacity for selecting more specific student cohorts can be increased manyfold, since the number of combinations of student characteristics increases geometrically.

Once a satisfactory data base has been developed, one might consider developing a model option that tabulates *all* major changes. This decision will undoubtedly be conditioned by whether the information regarding all former majors of each student is available and, if it is, how time consuming it would be to have it added to the model's data base.

A further refinement would be to increase the flexibility of the tracking period. In most uses of the model, a cohort is specified—say "all students entering as freshmen in fall 1978," and then tracked from first to last

major from 1978 to some future date such as 1984. In such a model run, the start of the 1978-1984 tracking interval coincides with the entry year of the cohort. One can imagine, however, a situation where it might be useful to track such a cohort during a single year only—say its sophomore year 1979-80. Such flexibility can be gained by programming in additional parameters (start and length of tracking period) that the user specifies before each model run.

Frequent usage has shown that, once the model has selected a cohort to track, it is often useful to have the model also print out some ancillary statistics concerning the cohort. Two that have been used often are a frequency table for the final grade point averages of the students and a frequency table for the number of students making no changes in major, one change, two changes, etc. Both of these tables are easy to obtain using standard SAS routines.

**See page 8 for
AIR Professional File
Ordering Information**

The *AIR Professional File* is intended as a presentation of papers which synthesize and interpret issues, operations, and research of interest in the field of institutional research. Authors are responsible for material presented. The *File* is published up to four times per year by the Association for Institutional Research.

Editor-in-Chief: John A. Lucas
Director, Planning & Research
William Rainey Harper College
Algonquin & Roselle Roads
Palatine, IL 60067

Managing Editor: Jean C. Chulak
Administrative Director
The Association for
Institutional Research
314 Stone Building
Florida State University
Tallahassee, FL 32306

©1987 The Association for Institutional Research

THE AIR PROFESSIONAL FILE

The AIR Professional File 1-16: Institutional Research Issues & Applications, 1978-1983—bound volume including the following issues—\$2 per copy, prepaid:

1. Organizing for Institutional Research (J.W. Ridge)
2. Dealing with Information Systems: The Institutional Researcher's Problems and Prospects (L.E. Saunders)
3. Formula Budgeting and the Financing of Public Higher Education: Panacea or Nemesis for the 1980s? (F.M. Gross)
4. Methodology and Limitations of Ohio Enrollment Projections (G.A. Kraetsch)
5. Conducting Data Exchange Programs (A.M. Bloom & J.A. Montgomery)
6. Choosing a Computer Language for Institutional Research (D. Strenglein)
7. Cost Studies in Higher Education (S.R. Hample)
8. Institutional Research and External Agency Reporting Responsibility (G. Davis)
9. Coping with Curricular Change in Academe (G.S. Melchiori)
10. Computing and Office Automation—Changing Variables (E.M. Staman)
11. Resource Allocation in U.K. Universities (B.J.R. Taylor)
12. Career Development in Institutional Research (M.D. Johnson)
13. The Institutional Research Director: Professional Development and Career Path (W.P. Fenstermacher)
14. A Methodological Approach to Selective Cutbacks (C.A. Belanger & L. Tremblay)
15. Effective Use of Models in the Decision Process: Theory Grounded in Three Case Studies (M. Mayo & R.E. Kallio)
16. Triage and the Art of Institutional Research (D.M. Norris)

The AIR Professional File—single issues (8pp.)—\$2 each, prepaid:

17. The Use of Computational Diagrams and Nomograms in Higher Education (R.K. Brandenburg & W.A. Simpson)
18. Decision Support Systems for Academic Administration (L.J. Moore & A.G. Greenwood)
19. The Cost Basis for Resource Allocation for Sandwich Courses (B.J.R. Taylor)
20. Assessing Faculty Salary Equity (C.A. Allard)
21. Effective Writing: Go Tell It on the Mountain (C.W. Ruggiero, C.F. Elton, C.J. Mullins & J.G. Smoot)
22. Preparing for Self-Study (F.C. Johnson & M.E. Christal)
23. Concepts of Cost and Cost Analysis for Higher Education (P.T. Brinkman & R.H. Allen)
24. The Calculation and Presentation of Management Information from Comparative Budget Analysis (B.J.R. Taylor)
25. The Anatomy of an Academic Program Review (R.L. Harpel)
26. The Role of Program Review in Strategic Planning (R.J. Barak)
27. The Adult Learner: Four Aspects (L. Jurand/G.D. Kuh & L.W. Cracraft/J.F. Campbell, D. Hentschel & L.M. Spiro/M.V. Mehallsis)

To order any of the above items, send your name, address, list of items, and full payment to The Association for Institutional Research, 314 Stone Building, Florida State University, Tallahassee, FL 32306. (No purchase orders or unpaid requests, please.)