

Application of the TETRAD II Program to the Study of Student Retention in U.S. Colleges

Marek J. Druzdzal
University of Pittsburgh
Department of Information Science
Pittsburgh, PA 15260
marek@lis.pitt.edu

Clark Glymour
Carnegie Mellon University
Department of Philosophy
Pittsburgh, PA 15213
cg09+@andrew.cmu.edu

Abstract

We applied TETRAD II, a causal discovery program developed in Carnegie Mellon University's Department of Philosophy, to a database containing information on 204 U.S. colleges, collected by the *US News and World Report* magazine for the purpose of college ranking. Our analysis focuses on possible causes of low freshmen retention in U.S. colleges. TETRAD II finds a set of causal structures that are compatible with the data.

One apparently robust finding is that student retention is directly related to the average test scores and high school class standing of the incoming freshmen. When test scores and class standing are controlled for, factors such as student faculty ratio, faculty salary, and university's educational expenses per student are all independent of both retention and graduation rates, and, therefore, do not seem to directly influence student retention. Furthermore, simple linear regression applied to test scores, class standing, and retention data showed that the test scores and class standing explain 52.6% of the variance in freshmen retention rate and 62.5% of the variance in graduation rate (test scores alone explain 50.5% and 62.0% respectively). This result becomes even stronger when computed for the group of top ranking colleges — regression applied to a group of 41 top ranking colleges showed explanation of 68.3% of the variance in freshmen retention rate and 77.0% in graduation rate (66.6% and 75% respectively for test scores alone).

As the test scores and class standing are both indicators of the overall quality of the incoming students, we predict that one of the most effective ways of improving student retention in an individual college is increasing the college's selectivity. High selectivity will lead to higher quality of the incoming students and, effectively, to higher retention rate.

1 Introduction

Even though some American colleges achieve a student retention rate of over 90%, the mean retention rate tends to be close to 55% and in some colleges fewer than 20% of the incoming students graduate (see Figure 1 for the distribution of graduation rates across a set of 200 U.S. national universities). Low student retention usually means a waste in effort, money, and human potential. Retention rate is often thought to indicate student satisfaction with their university program and, hence, indirectly, the quality of the

university. Indeed, a significant correlation can be observed between university ranking and retention rate — universities close to the top of ranking lists tend to have high retention rates. Is a university's low student retention rate an indication of shortcomings in the quality of education, facilities available to students, tuition costs, university's location, or perhaps wrong admission policies? More importantly, what action can the university take to improve the student retention rate? Can such actions as higher spending on student facilities, increasing the student/faculty ratio, increasing quality standards for teaching faculty, or modifications to admission policies make a difference?

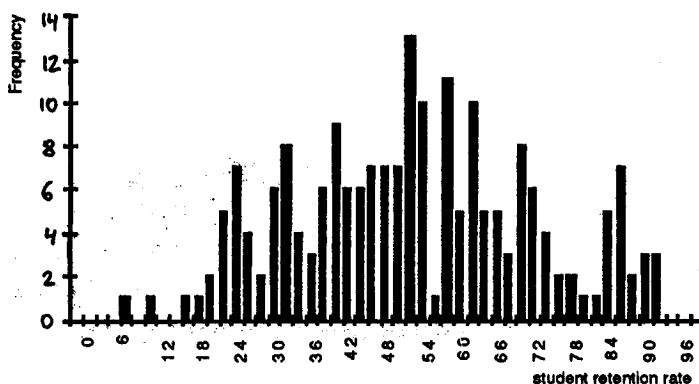


Figure 1: Histogram of the graduation rates for 200 U.S. national universities (Source: *U.S. News and World Report*).

This paper describes a preliminary effort to see what, if anything, aggregate data for many U.S. universities can tell us about the problem. Policy decisions demand that we know the causal structure of the system that we want to manipulate, and we therefore believe that determining the interactions among different relevant variables, including the direction of these interactions, is the necessary first step in addressing the problem. As university-scale experiments may be too expensive, ethically suspect, or otherwise impractical, our research needs to rely mainly on observations. The analysis has to be practically limited to extracting patterns from large collections of measurements of relevant variables. Our analysis involves data concerning 204 U.S. colleges, collected annually by *U.S. News and World Report* for the purpose of their college ranking (the data available to us is for 1992). In our analysis, we apply TETRAD II [3], a program embedding recently developed methods for causal discovery from observations. These methods, described in [4], are closely related to those employed in the induction of probabilistic models from data (e.g., [1]). While we are far from giving clear cut answers to the questions posed above, we believe that our analysis provides some interesting insight into the problem. The available data suggests that the main factor in student retention among the studied variables is the average test scores (or other measures of academic ability) of incoming students. The test scores of matriculating students are a function of the quality of the applicants and the university's selectivity. High selectivity leads to high average test scores of the incoming students and effectively to higher freshmen retention and graduation rates. Factors such

as student faculty ratio, faculty salary, and university's educational expenses per student do not seem to be directly causally related to freshmen retention. This hypothesis should be checked using data internal to any particular university, especially since the national data are aggregated to include both academic and non-academic dropouts. If the national pattern is confirmed locally, we would suggest that, wherever possible, steps aimed at making the university more selective be taken. Improving the comparative image of the school, and therefore increasing the number of applicants, increasing the selectivity of the admission process, increasing the chance that good applicants will accept admission offer rather than choosing another university, should improve student retention in the long run.

The remainder of the paper is structured as follows. We describe the analyzed data set (Section 2) and our assumptions about this data (Section 3). Then we summarize our view of the system that will provide us with prior information about the problem, useful in causal discovery procedures (Section 4). The results of our analysis are presented in Section 5. Section 5.1 presents the results of TETRAD II's search for possible causal structures that generated the data and Section 5.2 reports the results of applying simple regression to selected interactions identified by TETRAD II. We finish with a discussion of these results and policy suggestions (Section 6).

2 The Data

The data used in our study consists of a set of statistics concerning 204 U.S. national universities and national liberal arts colleges¹ collected by the *U.S. News and World Report* magazine for the purpose of college ranking. To prepare the data for its annual ranking of colleges,² *U.S. News* each year goes through a laborious process of data collection from several hundred of U.S. colleges. The data is collected from various university offices, such as admissions or business office, by means of surveys prepared by outside companies. It is subsequently verified by the schools representatives. The process of collecting the data and combining them into the final college ranking is described in [2].

We started with four spreadsheet files for 204 national universities provided by *U.S. News and World Report: Instructional Resources Ranking, Selectivity Ranking, Retention Ranking, and Financial Resources Ranking*. Each of the four spreadsheets contained the 204 universities ranked from the best to the worst in the respective category. To bring together various measurements and to relate the two variables of interest, freshmen retention rate and graduation rate to such indicators as colleges' selectivity, financial and instructional resources, we combined the four spreadsheets into one large spreadsheet containing over 100 variables measured for each of the 204 universities. Many of these variables were analytical derivatives of other variables (e.g., retention rate was simply the ratio of graduating seniors to incoming freshmen, both numbers included separately in the spreadsheet).

¹Defined as major research universities and leading grantors of doctoral degrees.

²The data available to us are for the year 1992.

The sample size, redundancy of the variable set, and missing values for various quantities, made it important to reduce the number of variables studied.³ We selected the following nine variables for our analysis: average percentage of freshmen retention (*apret*), average percentage of graduation (*apgra*), rejection rate (*rejr*), average test scores of the incoming students (*tstsc*), class standing of the incoming freshmen (*top10*), which is percentage of the incoming freshmen who were in top 10% of their high school graduating class, percentage of admitted students who accept university's offer (*pacc*), total educational and general expenses per student (*spend*), which is the sum spent on instruction, student services, academic support, including libraries and computing services, student teacher ratio (*strat*), and average faculty salary (*salar*). Describing each of over 100 remaining variables and discussing why we have not considered them for our analysis would make this paper unacceptably long. We limit ourselves to a few remarks. The values of a large number of the variables were included indirectly in the nine chosen variables. Average test scores of incoming students (*tstsc*), for example, is a normalized compilation of values of 14 variables, including a breakdown of average results for various parts of SAT and ACT tests. Average percentage of freshmen retention (*apret*) and average percentage of graduation (*apgra*) express the essence of all 14 variables in the *Retention Ranking* file. Rejection rate (*rejr*) and percentage of admitted students who accept university's offer (*pacc*) express, along with the average test scores (*tstsc*) and class standing (*top10*), selectivity of the school. We chose the total educational and general expenses per student (*spend*), student teacher ratio (*strat*) and average faculty salary (*salar*) as indicators of the quality of school's teaching and financial resources.

From the complete set of 204 universities, we removed 23 universities that had missing values for any of the nine variables of interest. This resulted in a set of 181 data points.

3 The Assumptions

Although TETRAD II's algorithms are independent on the actual distribution of the variables, they rely on the outcomes of a series of statistical tests. The necessary tests are especially powerful if we can assume normally distributed, linearly related variables. We studied how reasonable this assumption is for the available data set by plotting histograms of each of the nine variables and scatter plots of each pair of the nine variables. By visual inspection of the histograms and scatter plots, we removed six data points from the set of 181 data points that appeared to be outliers. The resulting data set, consisting of 175 data points, reasonably satisfies the normality and the linearity assumptions. All histograms were close to symmetric unimodal distributions (see Figure 2 for an example), with the exception of two positively skewed variables, *spend* and *strat*. The interactions

³A reviewer asked why any variables were omitted at all, and why covariances were not computed by simply skipping missing data points. The power of statistical tests and the reliability of search depend on the ratio of the number of sample points to the number of variables: The higher the ratio, the better. Including variables with missing values and calculating covariances by skipping a particular unit for a particular variable, as the reviewer suggested, would undermine the theoretical reliability of statistical tests. Testing partial correlations involves multiple correlations from the correlation matrix and, since

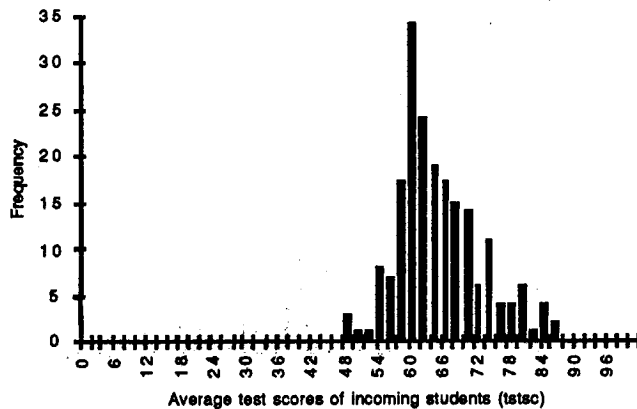


Figure 2: Histogram of the test scores *tstsc* for the 175 data points.

between different pairs of variables could be viewed as approximately linear (see Figure 3 for an example).

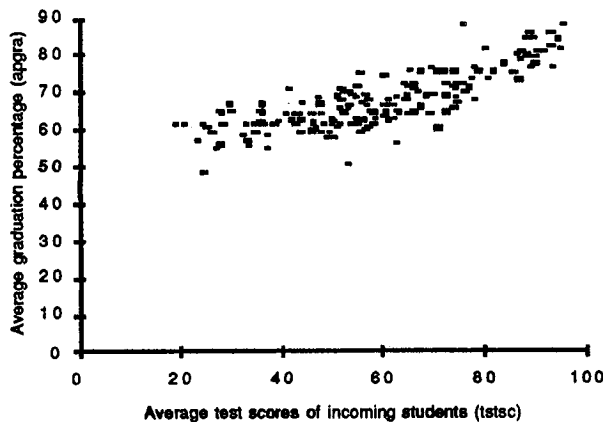


Figure 3: Interaction between *tstsc* and *apgra* for the 175 data points.

An important assumption made by TETRAD II is that the causal structure that generated the data points is acyclic. This assumption is not necessarily true in our data set. For example, most of the variables considered influence the image of the university. The image, in turn, can be argued to influence all of the nine variables. We still think that the acyclicity assumption is reasonable in our data set, as all feedback processes that we can think of in this context are extremely slow acting (at least on the order of decades as opposed to the interaction of our interest between the measured factors and retention

these would not be based on a fixed sample size, the sample size used in the tests would be indeterminate.

rate), so that in the snapshot provided by the 1992 data points they can be assumed negligible.

An assumption frequently made in causal modeling is causal sufficiency, which is an assumption that the analyzed variables form a self-contained structure — there are no latent common causes. An equivalent of this assumption is the assumption that all error terms are independent. TETRAD II allows for search with both the causal sufficiency assumption and without it. As it is unlikely that the selected variables form a self-contained structure, we have run TETRAD II without making the causal sufficiency assumption. Several control runs with causal sufficiency assumption did not reveal anything that would put our main conclusions in question.

4 Prior Knowledge

Interactions between some of the considered variables are reasonably well known. For example, we know the formula for computing the rejection rate, acceptance rate, retention, and graduation rates. We know what determines the tuition amount, the number of accepted students, the average faculty salary, etc. In several discussions between us and our colleagues, we developed a reasonable consensus on the causal graph that involves the analyzed variables (see Figure 4). We believe that a variable that we named *image*

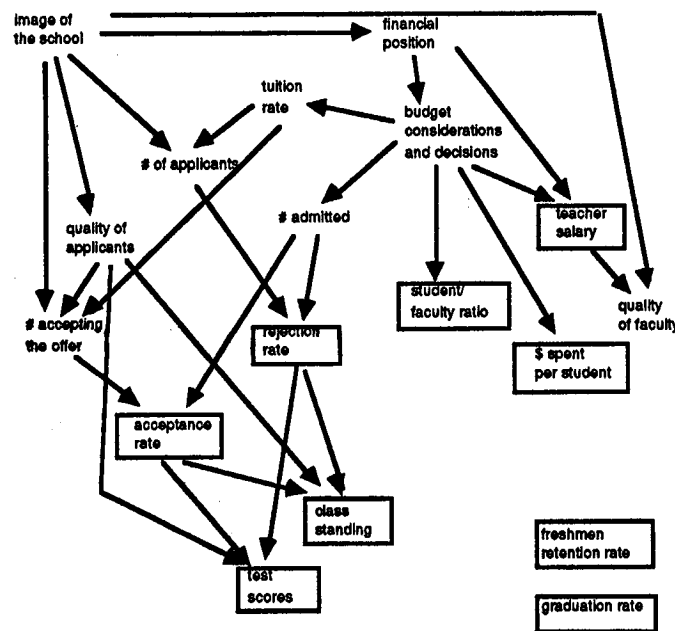
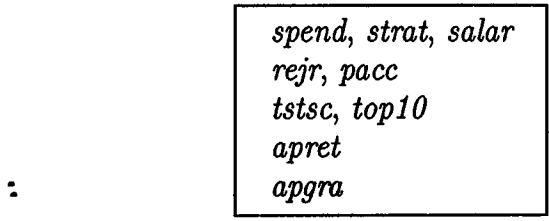


Figure 4: Initial assumptions about causal interactions in the system. Variables inside rectangles are analyzed in our study.

of the school directly influences the number of applicants, their quality, and the number of them that will accept the university's admission offer. It also influences the quality of

the faculty and the financial situation of the college (consider, for example, endowment funds and higher likelihood of external funding). Financial status of the school influences indirectly the tuition rate, the number of admitted students, student teacher ratio, average faculty salary, and quality of the faculty. Rejection rate and admission acceptance rate are determined by the number of applicants, the number of admitted students, and the number of them accepting the university's offer. The average test scores of incoming freshmen are determined by the overall quality of the applicants, the rejection rate, and the admission acceptance rate. Finally, our only assumption about how freshmen retention and graduation rates fit into this structure is that they do not cause any other variables considered.

The only purpose for showing Figure 4 is to make explicit the time order among the studied variables. In particular, the average spending per student (*spend*), student teacher ratio (*strat*), and faculty salary (*salar*) are determined based on budget considerations and are not influenced by any of the five remaining variables. Rejection rate (*rejr*) and percentage of students who are offered admission (*pacc*) precede the average test scores (*tstsc*) and class standing (*top10*) of incoming freshmen. The average freshmen retention rate (*apret*) precedes average graduation rate (*apgra*) because graduation rate depends on freshmen dropouts but also on dropouts in later years. We used only the temporal ordering of variables captured below as information to restrict the model search for TETRAD II.



5 The Results

While applying, for example, simple regressions to the data would allow us to make predictions about the value of a variable of interest given the values of other variables, this would not be sufficient for our purpose. What we want is to predict the effects of external manipulations of the system by means of new policies aimed at improving the retention rate. For this, we need information about the underlying causal structure of the system. We describe the results of the search for a class of causal structures that could possibly have generated the analyzed data set by means of a causal discovery program, TETRAD II, in Section 5.1.

In Section 5.2 we describe the results of measuring the strength of the most important causal connections suggested by the data: from the average test scores and class standing to retention rate and from test scores and class standing to the graduation rate. We apply simple linear regression to obtain a quantitative estimate of the interaction between these variables. We emphasize that we used regression only to estimate the coefficients in a linear model obtained by the TETRAD II search. If regression were used instead to search for the variables influencing retention and graduation, it would include variables

that TETRAD II says have no direct influence on the outcome, and that are conditionally independent of the outcome variables.

5.1 TETRAD II

When TETRAD II is run on normally distributed data with the linearity assumption, it converts the raw data into a correlation matrix. The values of the elements of this matrix is all that matters in discovery. The correlation matrix for all 175 data points is reproduced in Figure 5.

	apret	apgra	rejr	tstsc	pacc	spend	strat	salar	top10
apret	1.00000								
apgra	0.78122	1.00000							
rejr	0.53434	0.54303	1.00000						
tstsc	0.70576	0.79334	0.67515	1.00000					
pacc	-0.28385	-0.26149	-0.00739	-0.11191	1.00000				
spend	0.52424	0.56882	0.61999	0.73886	-0.11454	1.00000			
strat	0.40727	0.47905	0.39634	0.55430	-0.17285	0.72463	1.00000		
salar	0.66202	0.65033	0.65577	0.75969	-0.29412	0.71291	0.44534	1.00000	
top10	0.68521	0.66603	0.68243	0.82430	-0.15524	0.67249	0.43016	0.68265	1.00000

Figure 5: Matrix of correlations among the analyzed variables (175 data points).

When making decisions about independence of a pair of variables conditional on a subset of the remaining variables, TETRAD II uses statistical tests (in the normal-linear case, standard z -test for conditional independence). The search begins with a complete undirected graph. Edges in this graph are removed by testing for appropriate conditional independence relations. If two variables a and b become independent when conditioned on a subset \mathcal{S} of the remaining variables, there is no direct causal connection between them — all interactions between a and b take place through intermediate variables included in \mathcal{S} . This is a simple consequence of two assumptions known as *Markov condition* and the *faithfulness condition* [4]. Orientation of the remaining edges is based on a theorem proven in [4]. For example, suppose that two variables a and b are not directly connected (i.e., there exists a subset of the remaining variables \mathcal{S} that makes a and b conditionally independent) and there is an edge between a and c and an edge between b and c . If a and b are independent conditional on \mathcal{S} and dependent conditional on $\mathcal{S} \cup c$, then a and b are both direct causal predecessors of c . In other words, the edges can be oriented from a to c and from b to c . Both, the process of removing edges and the process of orienting edges, can be aided by prior information about the underlying graph. TETRAD II allows for specifying presence or absence of direct connections between pairs of variables and also temporal precedence among the variables. Knowledge of temporal precedence allows for limiting the number of tests for conditional independence and, under certain circumstances, aids in orienting the edges of the graph. If, for example, variables a and b are directly connected, there is no latent common cause of a and b , and a precedes b in time, then the edge can be oriented from a to b . The details of TETRAD II's search algorithm are given in [4].

Depending on the significance level used in independence tests, TETRAD II's individual statistical decisions regarding independence may be different and a different class of causal structures may result. It is, therefore, a good practice to run the program at several significance levels. We ran TETRAD II with the following significance levels: $p = 0.2$, 0.15, 0.1, 0.05, 0.01, and 0.001. The core of the structure, i.e., how freshmen retention rate and graduation rate are related to the remaining variables, was insensitive to changes in significance. This suggests that the structure proposed by TETRAD II is robust. The graphs proposed by TETRAD II for significance levels $p = 0.05$ and $p = 0.001$ are presented in Figure 6. The edges of the graph have the following meaning: A single arrow (\longrightarrow) denotes a direct causal influence. A double headed arrow (\longleftrightarrow) between two variables denotes presence of a latent common cause of these two variables. A single arrow with a circle at one end ($\circ\longrightarrow$) expresses TETRAD II's inability to deduce whether there is a direct influence between the two variables (\longrightarrow) or a latent common cause between them (\longleftrightarrow). An edge with circles at both ends ($\circ\text{---}\circ$) expresses TETRAD II's inability to deduce whether there is a direct influence between the two variables and, if so, what is its direction (\longrightarrow or \longleftarrow) or a latent common cause between them (\longleftrightarrow).

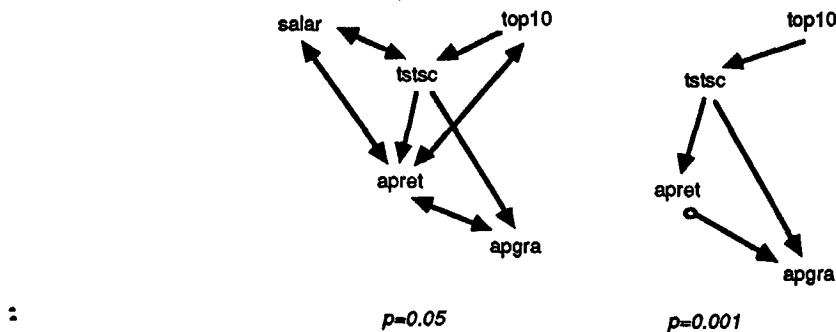


Figure 6: Two relevant parts of causal graphs proposed by TETRAD II for the complete data set of 175 universities (significance levels $p=0.05$ and $p=0.001$).

In both of the graphs in Figure 6 as well as in most of the graphs suggested by TETRAD II any connection between *apret* and *apgra* and variables like *spend*, *strat*, or *salar* is through *tstsc* or *top10*. The "latent common cause" connection between *salar* and *apret*, shown in Figure 6 for $p = 0.05$, disappears at $p < 0.04$. Most graphs contained a direct causal connection between the average test scores and freshmen retention. Also, the graphs contain a direct (or through a common cause) connection between freshmen retention and graduation rate.

TETRAD II's algorithms are much more reliable in determining existence of direct causal links than in determining their orientation. Therefore, prior knowledge supplied to TETRAD II may be critical for the orientation of edges of the graph. We used the temporal sequence described in Section 4, but we also checked the robustness of our result to temporal ordering by running TETRAD II with no assumptions about temporal precedence. Although TETRAD II proposed different orderings of variables, all direct

links, and the direct link between test scores and retention and graduation in particular, were the same in both cases.

To check whether the causal structure is the same for the top-ranked universities we prepared two additional data sets for TETRAD II: one with universities that were in the top 50 universities on at least one of the four lists, and one with universities that were in the top 30 on at least one of the four lists. The two data sets contained 74 and 41 data points respectively. The results are similar for each of the three data sets. Any differences can be partially attributed to a significantly smaller number of data points and, hence, higher susceptibility to chance variations. Figure 7 shows two graphs obtained for the set of 41 top ranking universities.

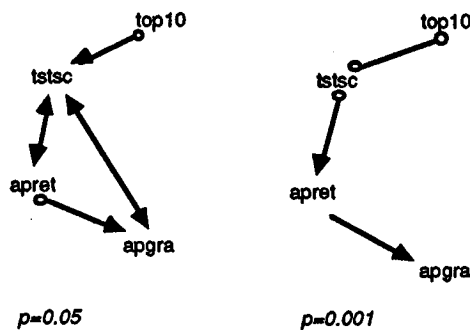


Figure 7: Two relevant parts of causal graphs proposed by TETRAD II for a subset of 41 top ranking universities (significance levels $p=0.05$ and $p=0.001$).

5.2 Linear Regression

We applied linear regression to the relation between the indicators of the quality of incoming freshmen: *tstsc* (average test scores) and *top10* (class rating) and *apret* (freshmen retention rate) and *apgra* (graduation rate) to obtain a quantitative measure of these interactions. In the full data set of 175 data points, linear regression applied to *apret* on *tstsc* results in the following equations:

$$\text{apret} = 33.4 + 0.142 \text{ top10} + 0.634 \text{ tstsc}, \text{ R-sq(adj)} = 52.6\%$$

$$\text{apgra} = -68.4 + 0.0283 \text{ top10} + 1.87 \text{ tstsc}, \text{ R-sq(adj)} = 62.5\%$$

In the restricted set of 74 top universities, the regression equations are:

$$\text{apret} = 49.8 + 0.0702 \text{ top10} + 0.490 \text{ tstsc}, \text{ R-sq(adj)} = 57.5\%$$

$$\text{apgra} = -69.0 - 0.116 \text{ top10} + 2.04 \text{ tstsc}, \text{ R-sq(adj)} = 61.7\%$$

In the restricted set of 41 top universities, the regression equations are:

$$\text{apret} = 53.7 + 0.0494 \text{ top10} + 0.468 \text{ tstsc}, \text{ R-sq(adj)} = 68.3\%$$

$$\text{apgra} = -73.0 - 0.150 \text{ top10} + 2.15 \text{ tstsc}, \text{ R-sq(adj)} = 77.0\%$$

As the coefficient of *tstsc* in all three equations is significantly larger than the coefficient of *top10* (note that it is in the groups of top ranking colleges actually negative), we repeated the procedure for *tstsc* as the only indicator, obtaining:

$$\begin{aligned} \text{apret} &= 13.2 + 1.02 \text{ tstsc}, R\text{-sq}(\text{adj}) = 50.5\% \\ \text{apgra} &= -78.7 + 2.04 \text{ tstsc}, R\text{-sq}(\text{adj}) = 62.0\% \end{aligned}$$

In the restricted set of 74 top universities, the same regression equations are:

$$\begin{aligned} \text{apret} &= 37.7 + 0.713 \text{ tstsc}, R\text{-sq}(\text{adj}) = 57.0\% \\ \text{apgra} &= -61.5 + 1.84 \text{ tstsc}, R\text{-sq}(\text{adj}) = 63.2\% \end{aligned}$$

In the restricted set of 41 top universities, the regression equations are:

$$\begin{aligned} \text{apret} &= 49.2 + 0.574 \text{ tstsc}, R\text{-sq}(\text{adj}) = 66.6\% \\ \text{apgra} &= -59.4 + 1.82 \text{ tstsc}, R\text{-sq}(\text{adj}) = 75.0\% \end{aligned}$$

Although the impact of test scores on the average freshmen retention rate and graduation rate is smaller for top ranking colleges (note a smaller value of the coefficient), these test scores explain more of the variance. In the group of top ranking colleges, the average test scores of incoming freshmen explain as much as 75% of the variance in graduation rates. Average test scores along with class standing explain as much as 77% of the variance in graduation rates.

6 Discussion

It seems that none of the variables in the data set are directly causally related to freshmen retention except for test scores and class standing. This result, following directly from the fact that freshmen retention rate and graduation rate are, given average test scores and class standing, conditionally independent of all remaining variables, seems to be robust across varying significance levels, availability of prior knowledge, and data set size. The average test scores seem to have a high predictive power for student retention and graduation rates. For the top 41 ranking colleges, average test scores in combination with class standing explain as much as 68.3% of the variance in freshmen retention rate and 77% of the variance in graduation rate.

Average test scores and class standing of incoming students can be viewed as indicators of the quality of incoming students. It seems that retention rate in an individual college can be improved by increasing the quality of the incoming students. This, in turn, can be improved by increasing the number and the quality of applicants. The better the pool of applicants from which an admission committee can select, the better the accepted students and, hopefully, the better the matriculating students are likely to be. Changing factors such as faculty salary, student/teacher ratio, or spending per student should, according to our result, have no direct effect on freshmen retention and graduation rates.

Theoretically, it is possible to use the regression coefficients between average test scores and retention rate obtained in this study to predict the impact of improvement in the average test scores of incoming students on freshmen retention and graduation. There are, however, potential problems with making predictions of an intervention at one institution, as the coefficients of the regression equations do not need to be identical for each institution.

One limitation in our study is that the available *U.S. News* data do not disaggregate academic from non-academic dropouts. We predict that internal data will show a difference between average test scores of dropouts (academic and non-academic) and graduates.

Another limitation is that our data do not disaggregate between different departments. Some departments may have many academic dropouts, others few. Also, the available data set did not include other variables that may have been relevant, as geographical location (climate, urban/rural, etc.), tuition costs, available academic support, financial situation of the students, prominence of athletics on campus, etc.

Finally, it is possible to apply alternative prior models of interaction of the variables in our data set. One alternative, suggested to us by Steven Klepper, might involve one latent variable influencing all nine variables studied. This model, however, would not account for the strong conditional independences observed in the data, and is in fact rejected by the standard f ratio test (Chi square of 356 with 27 degrees of freedom).

7 Acknowledgments

Considerable data collection effort and generosity in making the collected data available on the part of *U.S. News and World Report* made this study possible. Steven Klepper, Chris Meek, Richard Scheines and Peter Spirtes contributed to our work with valuable suggestions. We thank Felicia Ferko, Kevin Lamb, and Jeffrey Bolton from Carnegie Mellon University's Office of Planning and Budget for enabling us to access the data files and providing insightful background information. Anonymous reviewers prompted us for more details in our presentation. Support for this work has been provided by ONR and NPRDC under grant N00014-93-1-0568 to Carnegie Mellon University.

References

- [1] Gregory F. Cooper and Edward Herskovits. A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, 9:309-347, 1992.
- [2] Robert J. Morse, Senior Editor. U.S. News & World Report's America's Best Colleges Rankings: How it's done. Technical report, U.S. News and World Report, Washington, DC, May 8, 1992.
- [3] Richard Scheines, Peter Spirtes, Clark Glymour and Christopher Meek. *TETRAD II: Tools for Discovery*. Lawrence Erlbaum Associates, Hillsdale, NJ, 1994.
- [4] Peter Spirtes, Clark Glymour and Richard Scheines. *Causation, Prediction, and Search*. Springer Verlag, 1993.