

Research & Occasional Paper Series: CSHE.8.07

CSHE | Center for Studies in Higher Education

University of California, Berkeley
<http://cshe.berkeley.edu>

**A Student Experience in the Research University (SERU)
Project Research Paper****

**INSTITUTIONAL VERSUS ACADEMIC DISCIPLINE MEASURES OF
STUDENT EXPERIENCE: A Matter of Relative Validity***

May 2007

Steve Chatman
University of California, Berkeley

Copyright 2007 Steve Chatman, all rights reserved.

ABSTRACT

The University of California's census survey of undergraduates, UCUES, presents an opportunity to measure both disciplinary and institutional differences in students' academic experience. Results from nearly 60,000 responses (38% response rate) from the 2006 administration found greater variance among majors within an institution than between equivalent majors across institutions. Cluster analysis techniques were employed to establish disciplinary patterns, with traditional distinctions between hard and soft sciences generally supported. Reporting practices called into question range from institutional comparisons that ignore academic program mix and discipline to campus performance comparisons that do not recognize pedagogical differences by academic major. More specifically, these results suggest that calls for comparable institutional performance measures, as proposed by the Spellings Commission, must take into consideration disciplinary differences in instruction.

Introduction

There is tremendous appeal in the idea that a series of aggregate institutional measures of performance, expressed in comparative context, will lead to educational improvement or at least will uncover less productive use of public and student revenue. It is an attractive notion but it is very likely misleading and counterproductive in application. This

** The SERU Project is a collaborative study based at the Center for Studies in Higher Education at UC Berkeley and focused on developing new types of data and innovative policy relevant scholarly analyses on the academic and civic experience of students at major research universities. One of the main products of the SERU Project has been the development and administration of the University of California Undergraduate Experience Survey (UCUES). For further information on the project, see <http://cshe.berkeley.edu/research/seru/>

* A version of this paper was presented at the SERU Project Symposium, Assessing the Undergraduate Experience in the Postmodern University, April 25, 2007, Berkeley, CA.

paper will concentrate on one example of publicly reported institutional performance, student survey outcomes, but will raise questions that extend to related outcomes measures. The performance measure in question is institution-level academic experience factor scores as measured by limited-response questionnaire items asked of current students. The specific example is a survey of all undergraduates attending any campus of a state university system, the *University of California Undergraduate Experience Survey*[†] (UCUES). The problems noted appear to be inherent in similar enterprises, *The College Student Report* of the National Survey of Student Engagement (NSSE) for example. It will be asserted that sufficient evidence exists to justify rejection of these measures as valid performance indicators and reconsideration of the institutional comparison effort prescribed by the Spellings Commission until such time as more valid measures are developed or data collection methodology is changed.

The problem derives from a commonly accepted but largely unsubstantiated premise – that the undergraduate experience at most campuses shares sufficient common characteristics to be fairly and accurately measured by single aggregate scores. More explicitly, the problem results from a belief that there should be sufficient components in common that single scores could be valid measures and could be used to assess relative performance. An interesting dialogue in *Inside Higher Ed* between Banta on one hand and Klein, Shavelson, and Benjamin on the other is illustrative. This paper will support Banta's position. In the exchange, Banta wrote first and issued a warning about the Spellings Commission's call for "the use of standardized tests of general intellectual skills to compare the effectiveness of colleges and universities (2007, p. 1)." Banta referred to her and her colleagues' considerable record in assessment and noted many well-established problems with sample-based institutional scores on standardized instruments. Banta proposed as a more viable alternative electronic portfolios and measures based in academic disciplines. It is Banta's recognition of variance by academic disciplines that is supported by this paper.

Responding to Banta, Klein, Shavelson, and Benjamin (2007), who identify themselves as being affiliated with the Collegiate Learning Assessment (CLA) program, wrote that the CLA measures abilities that "cut across academic disciplines and...assesses these competencies with realistic open-ended measures that present students with tasks that all college graduates should be able to perform (p. 2)." They go on to assert the public interest in performance data to determine whether "the students at a given school are generally making more or less progress in developing these abilities than are other students (p. 2)" and conclude by stating that the CLA is the best currently available source of that information. Their argument in support of comparative sample-based summary scores generally, the CLA specifically, and against measuring those skills as taught and learned in academic disciplines appears to be two-fold: first, that these are "broad competencies that are mentioned in college and university mission statements (p. 2)" and second, that legislators, college administrators, many faculty, college-bound students and their parents, the general public, and employers want evidence of competencies regardless of academic major.

Whether or not the conventional wisdom/public interest argument made by Klein et al. or the experience of Banta and colleagues is asserted, a more basic issue may be the lack

[†] The *University of California Undergraduate Experience Survey* (UCUES) is the principal data collection effort of UC Berkeley's Center for Studies in Higher Education project, Student Experience in the Research University (SERU).

of common course experience by undergraduates. There is very little general education in common at large public research universities. One illustration of the variance in student experiences at a large public research university is provided by Chatman (2004) who examined general education policy and student behavior at one institution. He found more than one thousand courses and millions of combinations of courses that might satisfy general educational requirements, and only four courses were taken by a majority of students. Perhaps that should not be surprising for a campus with a cafeteria system and more than one hundred undergraduate academic majors. Given such a large number of majors and courses that can be counted toward satisfying requirements, the notion of a widely shared, common experience would seem to be an invalid premise on its face. And yet, it is a recurring theme from both inside and outside the academy.

The external call for comparable performance measures most recently includes Education Department Secretary Spellings' *Commission on the Future of Higher Education*. On page 25 of the *Test of Leadership: Charting the Future of U.S. Higher Education* (2006), under recommended changes to accrediting standards, is the following (emphasis added):

Accreditation agencies should make performance outcomes, including completion rates and student learning, the core of their assessment as a priority over inputs or processes. A framework that aligns and expands existing accreditation standards should be established to (i) *allow comparisons among institutions regarding learning outcomes and other performance measures, ... In addition, this framework should require that the accreditation process be more open and accessible by making the findings of final reviews easily accessible to the public* and increasing public and private sector representation in the governance of accreditation organizations and on review teams. Accreditation, once primarily a private relationship between an agency and an institution, now has such important public policy implications that accreditors must continue and speed up their efforts toward transparency as this affects public ends.

These are admirable standards that higher education would likely embrace if it were confident that it could effectively measure and then communicate the complexity of higher education. Modern public research universities are academically diverse and, by publicly supported agreement, serve extremely diverse populations. The accountability strategies that have been at least partially successful in improving elementary and secondary education cannot be easily generalized to postsecondary study because postsecondary education is more complex by at least an order of magnitude. Elementary schools offer few course choices, secondary schools several more within a few program tracks, and postsecondary institutions a 100 or more academic majors with thousands of courses. Is there cause for concern that the Spellings' Commission would subject higher education to reporting that could only grossly oversimplify performance?

On page 23, the Spellings report cites NSSE as an example of student learning assessment, stating the following (emphasis added):

Administered by the Indiana University Center for Postsecondary Research, the National Survey of Student Engagement (NSSE) and its community college counterpart, the Community College Survey of Student Engagement (CCSSE), survey hundreds of institutions annually about student participation and

engagement in programs designed to improve their learning and development. The measures of student engagement – the time and effort students put into educational activities in and out of the classroom, from meeting with professors to reading books that weren't assigned in class – *serve as a proxy for the value and quality of their educational experience. NSSE and CCSSE provide colleges and universities with readily usable data to improve that experience and create benchmarks against which similar institutions can compare themselves.*

NSSE is one of three examples offered, but attention is focused on the NSSE example here because it shares similarities with the source of data for this study, UCUES.

Whether striving to accurately assess a performance construct or to assess relative institutional performance by comparison, too little consideration is given by the Spellings Commission, and others who would hold higher education accountable, to the question of whether institution-level statistics are valid measures for the proposed purposes. At least in this area of assessment, recent evidence provided by NSSE researchers Nelson Laird et al. (2005, 2006) and UCUES researchers Brint (2006) and Chatman (2006, 2007) indicates that institution-level measures of student academic experience may be too crude to reflect real differences in performance, especially for large institutions offering a wide range of majors and courses, because they do not account for disciplinary differences in students' academic experience.

Relevant Research

From NSSE

Two publications reporting reliable disciplinary differences from NSSE and Faculty Survey of Student Engagement (FSSE) administrations are Nelson Laird, Shoup and Kuh's 2005 AIR paper on deep learning, *Deep Learning and College Outcomes: Do Fields of Study Differ?* and Nelson Laird, Schwarz, Kuh, and Shoup's 2006 AIR paper, *Disciplinary Differences in Faculty Members' Emphasis on Deep Approaches to Learning*. Deep learning, from an information processing perspective, refers to student generated efforts to increase the number and organization of associations formed between new information and information already in memory. The first paper, using student responses and a deep learning scale derived from 13 NSSE questionnaire items, found the following disciplinary differences for senior respondents:

- Students in social sciences, arts and humanities, professional programs (e.g., architecture, urban planning, nursing), and education scored higher on deep learning. Business, physical science and engineering scored lower on the deep learning scale. Biological sciences majors were midrange.
- Subscale high-order learning favored professional and engineering students.
- Both other subscales, integrative learning and reflective learning, were highest for social science and arts and humanities students and were lowest for physical science and engineering students.

These findings were generally supported when the same analytical strategy was applied to faculty responses on the FSSE:

- Education, arts and humanities, and social science faculty described using pedagogical practices that emphasized deep learning more often, and engineering and physical science faculty used the practices less often.
- Higher-order learning techniques were used less frequently in biological sciences and were uniformly more frequent in the other fields.
- Use of pedagogical practices to encourage integrative learning was highest in education, arts and humanities, and social sciences and was lowest in the physical sciences.
- Reflective learning was more frequently used in education, arts and humanities, and social science and was less frequently used in engineering and physical science.

The most common pattern, where arts and humanities and social sciences scored higher and science and engineering scored lower, was consistent from Nelson Laird et al.'s NSSE (2005) and FSSE (2006) studies. Based solely on these findings, it would be reasonable to assert that social sciences and arts and humanities graduates would have experienced a better education than science and engineering graduates. Of course, it would be a more persuasive argument if social science and arts and humanities students were in greatest demand at graduation and were able to command the highest salaries.

Nelson Laird et al. cite several publications reporting advantages of deep learning processing (2005) and conclude based on the analysis of observed variance in scores that there is room for improvement in every field of study and that there are good examples of how to improve within each disciplinary area. While there were serious limitations with both studies (disproportionate participation by discipline in the first and faculty self-selection of a single course to describe in the second), this paper will not belabor the argument whether deep learning is a valuable and valued construct. This study is concerned with the use of this institutional measure, or very likely any other institutional measure of student academic experience, as an indicator of comparative institutional performance. Unless it is assumed that all academic majors should be taught using the same strategies, then the data provided by Nelson Laird et al. (2005, 2006) show that an institutional outcome measure of engagement would reflect program mix.

From UCUES

Two University of California researchers, Brint and Chatman, have examined difference in student academic experience by major using UCUES results. The first study used data from the 2006 UCUES administration where more than 150,000 students across a university system were invited to participate in the survey; 38% responded overall and more than 32% responded at each campus. Brint's study examined responses by upper-division students completing the academic core component that is common to the various UCUES forms. (UCUES is comprised of a common academic core and one of four or five randomly assigned modules, depending on campus choice.) Using factor analysis to operationally define dimensions of student academic engagement (n~28,000), Brint (2006) found two types of student engagement, one that he asserted to be more typical of humanities and social sciences and the other more typical of the sciences. These hypotheses were confirmed. "Students in the arts, humanities, and social sciences score higher than students in other majors on the 'humanities culture' scale,

and they score much lower than other students on the 'sciences culture' scale ..." (p. 13) In addition to expected differences by major, he found the following results:

- SAT verbal was a significant predictor of humanities culture score and SAT math score was a predictor of sciences culture score.
- Campus was a minor explanatory factor for sciences culture and was not associated with humanities culture.
- GPA was positively associated with the humanities culture score but with lower study time.
- Sciences culture score was not related to GPA but was associated with study time.

Brint explained the GPA, study time and scale score associations as reflecting disciplinary differences in grading practices. Brint, like Nelson Laird, proposed overcoming the observed differences but unlike Nelson Laird, Brint recognized that there were limitations in each culture. Brint also identified about 10% of students in both fields as very engaged, hard working, and active learners who were exemplars.

Chatman (2005) attempted to replicate Laird et al. (2005) using UCUES census-based results for a single campus instead of NSSE sample-based results across many campuses. Over a five-factor varimax solution, Chatman found patterns similar to those reported by Laird, essentially higher scores for engagement in letters and social sciences, lower academic engagement scores for engineering and physical sciences, and biological sciences in a middle range. Chatman also described an example of earlier UCUES results (2004) where students in engineering at one campus scored lower on a long list of academic items than did the other students at the same campus but scored essentially the same as engineering students at other campuses—an applied example of the fact that variance is greater across disciplines than across campuses. In this engineering instance, intra-institutional comparison would have led to a dramatically different summative judgment of performance than would inter-institutional comparison made using the same academic discipline at other campuses.

Impact of Disciplinary Patterns on Performance Scores and Interventions

Collectively, these NSSE and UCUES results suggest that there are real disciplinary differences in academic engagement specifically and academic experience generally. Given valid and reliable disciplinary patterns, institutional summary scores would appear to be poor measures for campuses with diverse majors. How might program mix impact the validity of deep processing as an institutional measure? Here are a few questions with answers that can be inferred from the extant research to illustrate the point:

- Question 1: Why would liberal arts institutions be expected to score higher than state schools?
 - Liberal arts schools have relatively more social science and humanities majors, and social science and humanities students have higher scores. Conversely, liberal arts schools often do not have lower-scoring engineering and business majors.
- Question 2: Explain how institutional scores can mask program deficiencies or areas of strength when comparing two institutions? (Give one example of a masked area of strength and one example of a masked deficiency.)

- First, if the deficit occurs at campus A in a field with higher scores on average, the campus mean could be the same as B if there were more students at A in that field.
- Second, if A has an area of strength in a field that is expected to score lower, A and B could still score the same overall if B had fewer students in the same field or more students in higher scoring fields.
- Question 3: Explain why comparing the average score for one major to the campus average is misleading?
 - Without knowledge of an expected score for the major, it is not possible to separate disciplinary effects from performance.

Institutional efforts to intervene to improve scores at a campus with lower scores would necessarily be diffuse if the campus were ignorant of relative performance by major. Such interventions would probably be unsuccessful because most faculty would rightly assume that they were not part of the problem. Sample-based statistics will not identify these patterns unless students are sampled at the level of the major and will likely provide erroneous information leading to misdirected intervention. It is akin to confounding within-group effects with between-group effects and thereby conveying little of importance (Zwick, Brown & Sklar, 2004). Given the importance of academic program, sample-based statistics are of questionable value in a high stakes environment.

Because there are known academic engagement differences by major and little evidence of common experience among students at large institutions, this paper asserts that institution-level measures of academic engagement are of limited use and mask more valid measures at the level of academic discipline. In fact, institution-level measures might well be a better reflection of program mix than campus performance. The obvious alternative to sample-based study or to a census study conducted at a single campus is census-based collection across multiple campuses. Until recently the resource expenditure to survey more than 100,000 students distributed across a state would have been prohibitive, but Internet delivery and email contact make multi-campus census surveys a viable alternative. In addition, a logging-in process can be used to identify responses for the purpose of linking questionnaire data with other student records. The resulting merged record is an exceptional resource for academic inquiry and administrative needs.

Methodology

The recently completed 2006 UCUES survey, which included all undergraduate students attending the University of California system (~153,000), attained a 38% response rate overall (~58,000 responses). Each student received a common core set of items and one of five randomly assigned modules: academic experience, civic engagement, student development, student services, or a campus-specific module (optional). Because the campuses share many similarities, including programs offered and selective admissions, these data should provide a unique opportunity to determine the extent to which academic experience varied by academic program and, if variance is observed, the extent to which programs can be combined based on similarity of student responses into fewer clusters. The process required two clusterings: a reduction of survey items into factor scores and a clustering of academic majors based on those factor scores. The analysis used the work of Luan, Zhao, and Hayek (2005) as a model, and focused on academic core items as the most salient assessment dimension. Institutional differences were controlled by restricting study to the undergraduate student bodies of

eight similar institutions of one university system. Analysis was further restricted to upper-division students with declared majors. These actions increase the likelihood of useful results but may limit generalization to large public research universities.

Results

Factor Scores

The UCUES factor analysis of the upper-division academic core was a statistically driven “consensus of judgment” process. The bulk of the analysis was performed by a seven-person team of faculty, institutional research and UCUES project representatives during a day-long working session where alternatives were considered in real time by running the programs and examining results collectively. The solution was done in two stages. The first stage identified principal components and used orthogonal solutions. The second stage was performed within each principal component set and used oblique solutions as it was understood that items within a principal component would be correlated. Again, consensus judgment regarding the best solution was used. The resulting solutions very closely followed empirical results but final placement was supplemented by judgment-based movement of a handful of items from one subfactor to another. The first session was followed by two shorter meetings during which factor names were attached and minor revisions were made. The final result was a solution with seven principal components. The factor names and their internal consistency (Cronbach’s Coefficient Alpha) were:

Factor 1: Satisfaction with Educational Experience (.92)

Factor 2: Current Skills Self-Assessment (Nonquantitative) (.91)

Factor 3: Gains in Self-Assessment of Skills (Nonquantitative) (.89)

Factor 4: Development of Scholarship (.89)

Factor 5: Understanding Other Perspectives (.85)

Factor 6: Research Experiences (.69)

Factor 7: Quantitative Professions (.64)

The factor solution process and results are described in detail elsewhere (Chatman, 2007) but a brief description of principal factors will be provided here. Satisfaction with Educational Experience was composed of 30 survey items ranging from global satisfaction with GPA, social experience, academic experience, etc., but mostly consisted of items regarding the major (e.g., advising, access, instruction). Current Skills Self-Assessment (Nonquantitative) was 13 self-ratings of general, research, and personal skills. The third factor was the difference between skills at entry and as currently rated for the skills comprising the second factor. Development of Scholarship consisted of a series reflecting Bloom’s taxonomy and includes critical reasoning and assessment, curricular foundations for reasoning and elevated academic effort. The fifth factor concerned development of an appreciation and understanding of the perspectives of others, based on interactions with students of different race, religion, gender, nationality, economic circumstance or sexual preference. Research Experiences was a cluster of six items included to reflect the unique opportunities available to students at a research university. The seventh factor, Quantitative Professions, included quantitative skills, collaborative learning experiences, and three items about choice of major (remuneration, prestige, and fulfillment). One additional scalelet (Pike, 2006) was used, Academic Time (time in class or lab and academic preparation). Factor scores were computed as the standardized mean of standardized item scores. In other words, item responses were first standardized and the mean of those responses was computed for

each student. These first two steps produced the raw factor scores. The raw factor scores were then standardized to produce a reported score with a mean of 5 and a standard deviation of 2.

Academic Major Clusters

Student major was assigned to one of 19 disciplinary clusters using local conventions. The clusters were similar to the level of aggregation achieved using a two-digit CIP code (e.g., communications, engineering, social sciences, biological sciences, letters, agriculture). Factor mean scores by discipline were computed for areas with 100 or more responding students. Those mean area scores were subjected to cluster analysis using an agglomerative hierarchical clustering based on centroid distance. There appeared to be a natural and reasonable cutoff at about 0.7 that produced seven clusters that are shown in Figure 1 and with a more complete description of the mapping of majors to clusters in Table 1.

The resulting academic topology creates an interesting mix, with many clusters confirming conventional wisdom and others raising interesting questions. One of the surprises was that area, ethnic, cultural and gender studies (Area) was quickly distinguished from other majors. (When the scores are shown graphically in the following section, area, ethnic, cultural and gender studies presents a remarkably strong profile from an engagement perspective.) The next content areas to separate from the pack were engineering, business administration, and mathematics and computer science. Physical science and biological sciences joined social sciences, humanities, and an agriculture and architecture cluster pair, as the majority cluster. If an institution were to create academic divisions to reflect this topology, the schools and colleges would probably be agriculture; architecture; humanities and social sciences; biological and physical sciences; area and ethnic studies; mathematics and computer science; business administration; and engineering. This seven-cluster solution was used to illustrate variation in scores by factor score.

Table 1: Factor scores for principal components by disciplinary clusters

Disciplinary Area	Principal Component Factors								#*	%
	F1	F2	F3	F4	F5	F6	F7	Ftb		
Agriculture	5.5	4.9	4.9	5.1	4.7	5.1	5.2	5.2	601	2.5%
Architecture	4.8	5.0	5.3	5.1	5.3	5.0	4.9	5.3	210	0.9%
Agriculture & Architecture	5.1	5.1	5.1	5.0	5.0	5.4	4.7	4.9		17%
Social Sciences	5.3	5.4	5.2	5.1	5.2	4.8	4.4	4.7	5,214	21.6%
Communications	5.2	5.5	5.2	5.0	5.1	4.8	4.3	4.6	542	2.2%
Education	5.0	5.5	5.6	5.1	5.3	4.9	4.9	4.9	78	0.3%
Public Administration	5.5	5.2	5.4	5.0	5.3	5.0	4.1	4.5	111	0.5%
Law	5.3	5.4	5.2	5.5	5.2	4.7	4.8	4.6	175	0.7%
Interdisciplinary Studies	5.2	5.3	5.4	5.1	5.2	5.2	4.6	4.9	950	3.9%
Foreign Languages	5.6	5.2	4.8	5.0	5.2	4.8	3.9	4.8	622	2.6%
Letters	5.4	5.5	4.8	5.2	5.0	4.8	3.8	4.7	1,631	6.8%
Psychology	5.0	5.1	5.1	4.9	5.0	5.5	4.5	4.7	2,175	9.0%
Fine Arts	5.1	5.5	4.9	5.0	5.1	5.0	4.2	5.2	1,415	5.9%
Humanities & Social Science	5.3	5.4	5.1	5.1	5.2	4.8	4.2	4.8		40%
Biological Sciences	4.9	4.7	4.9	5.0	4.9	5.6	5.4	5.2	2,660	11.0%
Physical Sciences	5.1	4.6	4.7	5.1	4.7	5.6	5.9	5.3	1,068	4.4%
Biological & Physical Sciences	4.9	4.7	4.9	5.0	4.9	5.6	5.6	5.2		15%
Area and Ethnic Studies	5.7	5.7	5.9	5.6	5.9	5.4	4.0	5.1	555	2%
Mathematics	4.9	4.3	4.4	4.8	4.6	4.5	5.9	5.1	539	2.2%
Computer Science	4.7	4.7	4.5	4.4	4.2	4.6	6.2	5.3	634	2.6%
Mathematics & Computer Science	4.8	4.5	4.5	4.6	4.4	4.5	6.1	5.2		5%
Business Administration	4.8	5.0	5.1	4.6	5.0	4.4	6.6	4.6	1,096	5%
Engineering	4.7	4.6	4.8	5.0	4.7	5.3	6.6	5.3	3,878	16%
Minimum	5.7	5.7	5.9	5.6	5.9	5.6	6.6	5.3		
Maximum	4.7	4.3	4.4	4.4	4.2	4.4	3.8	4.5		
Range	1.1	1.4	1.5	1.2	1.7	1.2	2.8	0.8		

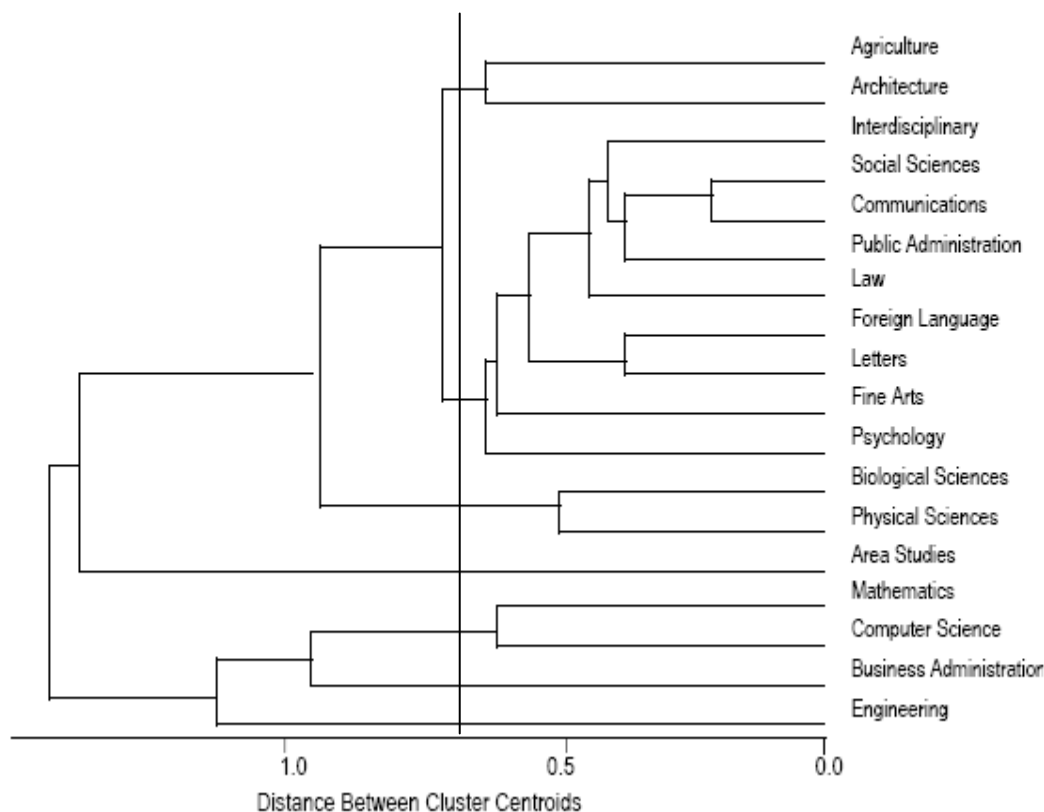
Factor Structure

F1	Factor 1: Satisfaction with Educational Experience
F2	Factor 2: Current Skills Self-Assessment (Nonquantitative)
F3	Factor 3: Gains in Self-Assessment of Skills (Nonquantitative)
F4	Factor 4: Development of Scholarship
F5	Factor 5: Understanding Other Perspectives
F6	Factor 6: Research Experiences
F7	Factor 7: Quantitative Professions
FTb	Factor Time: Subfactor Tb -- Academic Time

* Minimum number of students used in computing a factor score for this discipline.

Figure 1: Empirically Derived Structure of the University

(Centroid Hierarchical Cluster Analysis: Agglomerative at Distance ~0.7)

***Factor Scores of Academic Major Clusters***

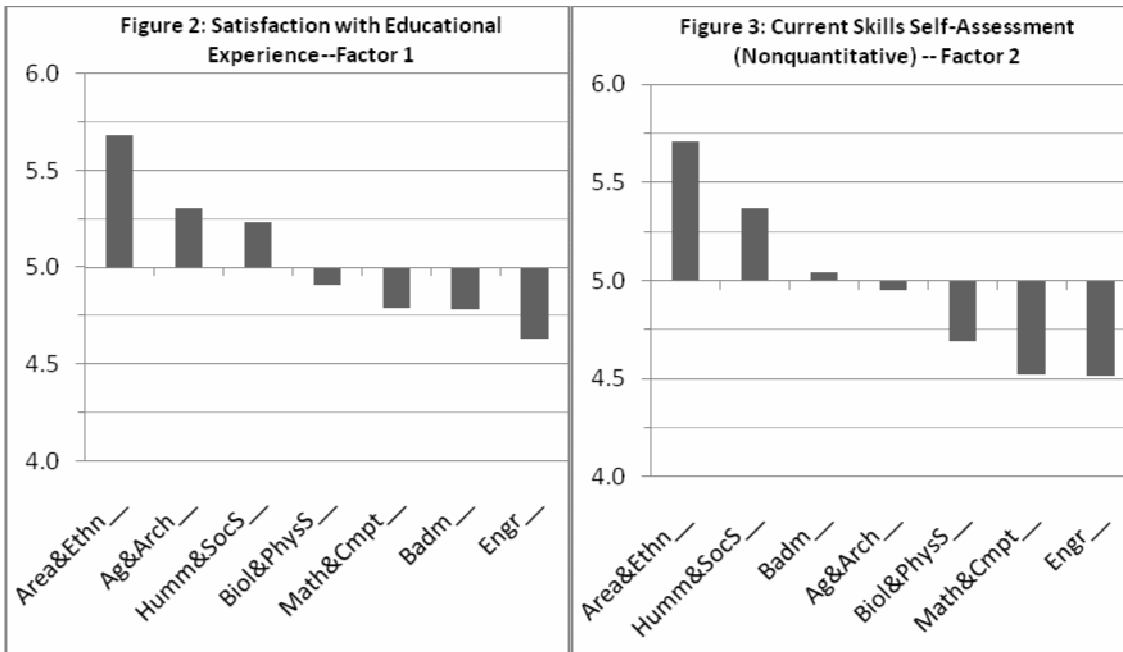
Scores on the first factor, satisfaction with educational experience, were highest in area and ethnic studies, agriculture and architecture, and humanities and social sciences. Satisfaction was lower in mathematics and computer science, business administration, and engineering (Figure 2). With a few position changes, the second factor, current skills self-assessment (nonquantitative), was similarly arranged (Figure 3). Area and ethnic studies and humanities and social sciences were at the upper end and mathematics and computer science and engineering were at the lower end. The profile for the third factor, gains in self-assessment of skills (nonquantitative), was very much like that of the second factor but with more variance at the extremes (Figure 4). Area and ethnic studies was more clearly distanced at the upper end and mathematics and computer science was more clearly distanced at the lower end. The fourth factor, development of scholarship, found four areas to be close to the overall mean: humanities and social sciences, biological and physical sciences, agriculture and architecture, and engineering (Figure 5). Again, distinguished at the upper end was area and ethnic studies. The lower end was held by mathematics and computer science and business administration. The fifth factor, understanding other perspectives, was thankfully highest in area and ethnic studies and unfortunately, but perhaps as expected, lowest in engineering and mathematics and computer science (Figure 6). Research experiences, the sixth factor,

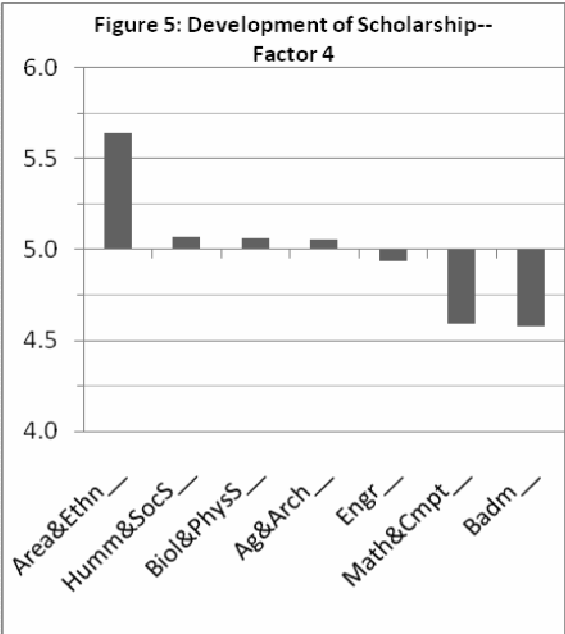
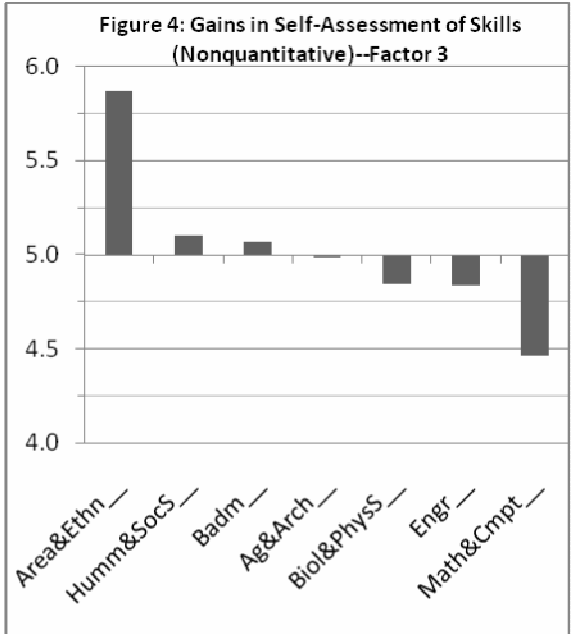
presented the first major reordering with biological and physical sciences, area and ethnic studies, and engineering leading the array. Mathematics and computer science and business administration were at the lower end of the array (Figure 7). Quantitative professions, the seventh factor, confirmed expectations with engineering, business administration and mathematics and computer science leading and humanities and social sciences and area and ethnic studies trailing (Figure 8). The academic time subfactor (treated as a principal factor here) placed science, engineering and mathematics (SEM) fields highest and humanities and social sciences, area and ethnic studies, and business administration lowest (Figure 9).

The relative variance explained by discipline and campus was determined for the eight factors (Table 2). In all cases, disciplinary cluster explained more variance in factor score than did campus, at least about twice as much and much more for the sixth, seventh, and academic time factor. It was also notable that the interaction of discipline and campus was much less important than either main effect.

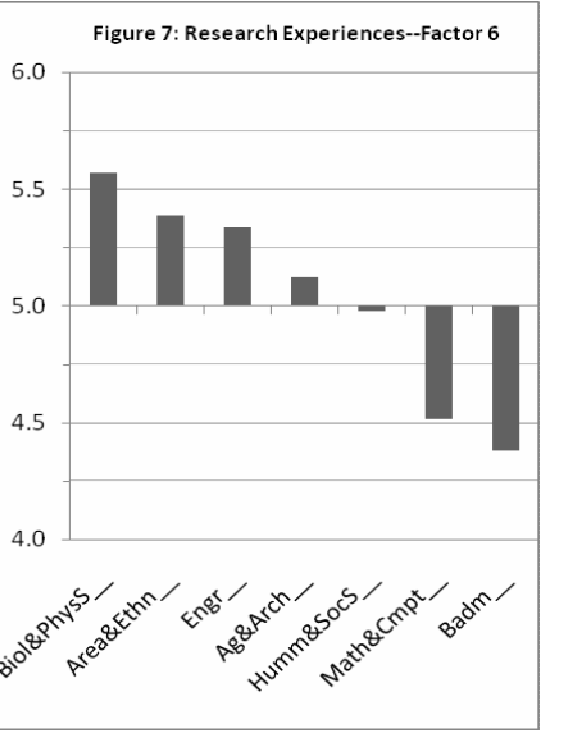
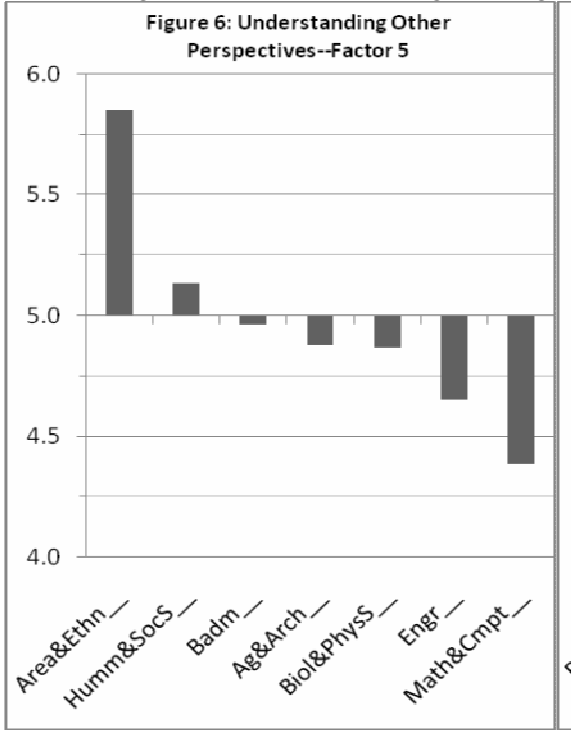
Summary

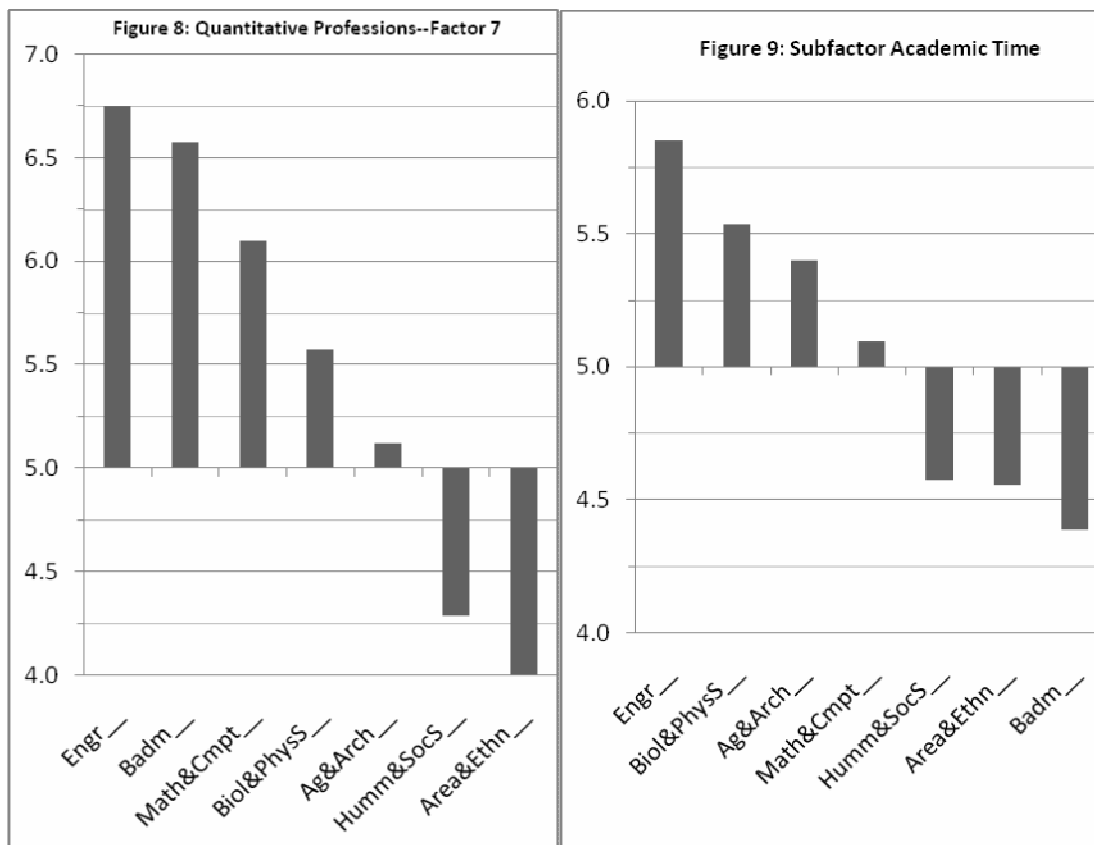
Previous research suggested disciplinary differences in educational engagement specifically and the academic experience generally. This project confirmed that differences do exist across a large public research university system; that the pattern of traditional engagement differences tend to favor social sciences, arts and humanities; and that by including items focused on research and collaborative learning, factors are found that favored students in mathematics, computer science, engineering and business administration fields. The most important result is that academic experience and student engagement varies by program of study in predictable ways. What does this finding mean for instruction?





Ag&Arch Agricultural Sciences and Architecture
 Humn&Soc Social sciences, communication, public administration, law,
 foreign language, letters, fine art
 Biol&Phys Biological sciences, physical science
 Area Area and ethnic studies
 Math&Cmpt Mathematics, computer science
 BAdm Business administration and management
 Engr Engineering





When they reached a similar point in their papers, Laird et al. (2005, 2006) and Brint (2006) began to suggest ways that instruction might be improved in the lower ranking fields (Laird) or that the better aspects of various fields might be used for common improvement (Brint). Both authors suggested that educational experience differences between disciplines should be reduced. That differences should be reduced is not a matter of concern for this paper, although it seems clear that more research is needed to understand why instructional practices differ by academic discipline before recommending that they be changed. After all, many of the programs described here are considered among the best in the country. Instead of suggesting changes, this paper was solely concerned with demonstrating that important differences do exist by academic discipline and that these differences would lead to misleading conclusions when comparing one program to a campus average and when comparing one campus to another. Actions then taken because of erroneous conclusions could hardly succeed. Worse, most institutions of higher education remain ignorant of these real differences because they rely on easily attained statistical samples that could not support analysis at the level of an academic discipline.

There is real danger in embracing the Spellings Commission recommendation to use widely available student engagement assessments to compare performance of one institution with another. Institution-level scores are simply inadequate. Unless the campuses to be compared are composed of the same programs in the same proportions, then the comparison will necessarily be biased by program composition. To illustrate this fact, bachelor degrees awarded by Association of American Universities institutions were clustered into this study's seven areas and assigned the mean values found in this study.

The results were then rank ordered. Using the first factor, Satisfaction with Educational Experience, at Harvard as an example, Harvard would be predicted to score very high because it has one of the highest proportions of humanities and social sciences students and few if any students in business administration, engineering and mathematics and computer sciences. Georgia Tech would be predicted to score low because it has one of the highest concentration of engineering students and a very small proportion of humanities and social sciences students. In other words, the 62 AAU institutions can be rank-ordered based solely on disciplinary composition and the tendency of students in disciplines to respond differently. Here are some of the hypothetical results:

Factor 1: Satisfaction with Educational Experience	
Top Five	Brandeis, Yale, Harvard, Catholic University, NYU
Range	0.44 (an effect size of .22)
Factor 2: Current Skills Self-Assessment (Nonquantitative)	
Top Five	NYU, Brandeis, Yale, Oregon, Emory
Range	0.64 (an effect size of .32)
Factor 3: Gains in Self-Assessment of Skills (Nonquantitative)	
Top Five	Brandeis, NYU, Yale, Emory, Oregon
Range	0.31 (an effect size of .16)
Factor 4: Development of Scholarship	
Top Five	Brandeis, Yale, Princeton, Cal-Davis, Harvard
Range	0.24 (an effect size of .12)
Factor 5: Understanding Other Perspectives	
Top Five	Brandeis, Yale, NYU, Emory, North Carolina
Range	0.41 (an effect size of .21)
Factor 6: Research Experiences	
Top Five	Cal Tech, Cal-Davis, Princeton, Case Western, Duke
Range	0.48 (an effect size of .24)
Factor 7: Quantitative Professions	
Top Five	Georgia Tech, MIT, Cal Tech, Purdue, Case Western
Range	1.77 (an effect size of .89)
Academic Time	
Top Five	Cal Tech, Georgia Tech, MIT, Case Western, Purdue
Range	0.95 (an effect size of .48)

The point of this example is that substantive differences in scale scores can occur as a result of nothing more than disciplinary composition. Even when two campuses are composed of the same programs in the same proportions, the summary score will most likely not reflect relative performance at the level of interest to faculty and student, the academic major or discipline. Simple measures to respond to public accountability desires may be more easily constructed for elementary schools and even for secondary schools because of curricular similarities, but the curriculum and curricular offerings of postsecondary schools appear to be too complex to be effectively reduced to a few numbers. If public accountability demands comparative performance, then the unit of analysis for performance should be the academic discipline.

An obvious limitation of this study results from the academic structure used to initially combine academic majors into a smaller number of units (equivalent to two-digit CIP). The same arguments that this paper made about the dangers of aggregation could extend to combining majors within any group. For example, there might be important

differences between civil and mechanical engineering, or the combination of programs within agriculture may mask the same type of differences seen at the campus level.

Setting those concerns aside for the moment, the relative validity of measures from derived disciplinary clusters and from institutional samples is important to understanding the student experience in higher education and whenever survey outcomes are used as accountability measures by which institutional performance may be compared. Perhaps the most valuable contribution of disciplinary-based measures is in program review, because program review happens at the level of the major where faculty recognize and bear responsibility for the academic experience.

Once it is recognized that institution-level measures are of questionable validity, leading to erroneous conclusions and offering little if any direction for improvement, it is obvious that accountability demands more. Imagine reporting to Proctor and Gamble (P&G) shareholders that consumers of its products were less satisfied than those who used Unilever's products. P&G produces about 100 brands distributed over about 25 categories, not so different from a large public research university. Unilever has about 30 brands, many competing for the same markets. Imagine that your research was based on a sample of P&G consumers and you are not able to report satisfaction by product line or to express relative satisfaction by product line for competing products. How would P&G begin to address the problem? Which division head would acknowledge that his or her brand was partially responsible for the lower score and should therefore be the one to improve? What reception would your report receive? More importantly, what reception should your report receive? Universities faced with the Spellings recommendation need to think about these types of questions.

REFERENCES

- Banta, T. W. (2007). A warning on measuring learning outcomes. *Inside Higher Ed* (Jan. 26).
- Brint, S. (2006). The Two Cultures of Undergraduate Academic Engagement and How to Bridge Them. Distinguished Educational Thinkers Speaker Series. UC Davis School of Education and Graduate Group in Education. Davis, CA.
- Chatman, S.P. (2004). General Education Coursework of Students Graduating in Spring of 2003. SARI Report No. 331. Davis, CA.
- Chatman, S.P. (2007). A Common Factor Solution for UCUES 2006 Upper-Division Core Items. Center for Studies in Higher Education. Berkeley, CA.
- Hartigan, J.A. (1975). *Clustering Algorithms, 99th Edition*. John Wiley & Sons, Inc. New York, NY.
- Klein, S., Shavelson, R., and Benjamin, R. (2007). Setting the record straight. *Inside Higher Ed* (Feb. 8).

- Nelson Laird, T.F.N., Shoup, R., & Kuh, G.D. (2005). Deep Learning and College Outcomes: Do Fields of Study Differ? Paper presented at the Annual Conference of the California Association for Institutional Research, San Diego.
- Nelson Laird, T.F.N., Schwarz, M.J., Kuh, G.D., Shoup, R. (2006). Disciplinary Differences in Faculty Members Emphasis on Deep Approaches to Learning. Paper presented at the Annual Forum of the Association for Institutional Research, May, 2006, Chicago, IL.
- Luan, J., Zhao, C-M, & Hayek, J. (2005). Exploring a New Frontier in Higher Education: A Case Study Analysis of Using Data Mining Techniques to Create NSSE Institutional Typology. Paper presented at the Annual Forum of the Association for Institutional Research, San Diego.
- NSSE (2004). *Student Engagement: Pathways to Collegiate Success*. NSSE annual report. Center for Postsecondary Research, Bloomington, IN.
- Pike, G. R. (2006). The convergent and discriminant validity of NSSE scalet scores. *Journal of College Student Development*, 47, 550-563.
- Ratcliff, J., & Associates (1988). Development and testing of a cluster-analytic model for identifying coursework patterns associated with general learning abilities of students. (Progress Report No. 6). Iowa State University, College of Education, Ames, IA.
- The Secretary of Education's Commission on the Future of U.S. Higher Education (2006). *A Test of Leadership: Charting the Future of U.S. Higher Education*, Education Publications Center, U.S. Department of Education, Jessup, MD.
- Zwick, R., Brown, T. & Sklar, J. (2004). California and the SAT: A reanalysis of University of California admissions data." Center for Studies in Higher Education, Research & Occasional Paper Series, Berkeley, CA.