# Students' Evaluations of University Teaching: Ratings of Teachers, Universities and Departments

## HEA Conference 19 May 2011

## Herbert W. Marsh
## University of Oxford

# Overview

- **Use of students' evaluations of TEACHERS (SETs) in individual university classes**

- **UK National Student Survey (NSS) &Australian Course Evaluation Questionnaire (CEQ) responses by undergraduates to evaluate UNIVERSITIES & DEPARTMENTS**

- **Australian Postgrad Research Experience Questionnaire (PREQ) to evaluate research training of research students in Australian UNIVERSITIES and DEPARTMENTS**

# Students' Evaluations of University Teaching (SETs)

# Purposes of SETs

- **diagnostic feedback to teachers about the effectiveness of their teaching that will be useful for the improvement of teaching;**

- **a measure of teaching effectiveness to be used in personnel decisions;**

- **information for students to use in the selection of courses and teachers; and**

- **an outcome or a process description for research on teaching.**

## The first purpose is nearly universal, but the next three are not.

Marsh, (2007). Students' evaluations of university teaching: A multidimensional perspective. In R. P. Perry & J C. Smart (Ed.), The Scholarship of Teaching and Learning in Higher Education: An Evidence-Based Perspective (pp.319-384). New York: Springer.

# Summary Conclusions

**My research has led me to conclude that SETs are:**

- **Multidimensional;**

- **Reliable and stable;**

- **Primarily a function of the instructor who teaches a course rather than the course that is taught;**

- **Valid in relation to a variety of indicators of effective teaching;**

- **Relatively unaffected by a variety of variables hypothesized as potential biases; and**

- **Seen to be useful by students for use in course selection, by administrators for use in personnel decisions, by teachers as feedback about teaching**

Marsh, (2007). Students' evaluations of university teaching: A multidimensional perspective. In R. P. Perry & J C. Smart (Ed.), The Scholarship of Teaching and Learning in Higher Educ: An Evidence-Based Perspective (pp.319-384). NY: Springer.

Marsh, H.W. & Roche, L.A (1997). Making students' evaluations of teaching effectiveness effective. Am Psychol, 52, 1187-1197.

# Student Rating Dimensions

# Dimensionality: The SEEQ Factors

**Learning/Value:** You found course intellectually challenging/stimulating;

**Instructor Enthusiasm:** Instructor dynamic/energetic in conducting course;

**Organisation:** Course materials were well prepared/carefully explained;

**Individual Rapport:** Instructor was friendly towards individual students;

**Group Interaction:** Students encouraged to participate in class discussions;

**Breadth of Coverage:** Presented background/origin of ideas/concepts;

**Examinations/Grading:** Feedback valuable from exams/graded materials;

**Assignments/Readings:** Readings, homework, etc. contributed to appreciation and understanding of subject;

**Workload/Difficulty:** Relative course difficulty (very easy...medium…very hard).

Marsh, , Muthén, B., Asparouhov, T., Lüdtke, O., Robitzsch, A., Morin, A. J. S., & Trautwein, U. (2009). Exploratory Structural equation modeling, integrating CFA and EFA: Application to students' evaluations of university teaching. Structural Equation Modeling,16, 439-476.

Marsh, , & Hocevar, D. (1991). Multidimensional students' evaluations of teaching effectiveness: Factor structure stability across academic discipline, instructor level, and course level. *Teaching and Teacher Education, 7,* 9-18.

# RELIABILITY AND STABILITY

# SEEQ Reliability

The reliability of SETs is most appropriately determined from studies of interrater agreement (i.e., generalizability of ratings over students in the same class).

The reliability of the class-average response depends upon the number of students rating the class; it is about

- .95 for the average response from 50 students,
- .90 from 25 students,
- .74 from 10 students, and
- .60 from five students.

Given a sufficient number of students, SET reliability compares favourably with the best objective tests. However, even for small classes good reliability can be obtained by averaging results from several classes.

Marsh, (1987). Students' evaluations of university teaching: Res findings, methodological issues, and directions for future research. *International J of Educ Res, 11,* 253-388 (whole issue).

# What is the Relative Importance of the Teacher vs. Course Effects

## How highly correlated are SETs in:

- two different courses taught by the same instructor
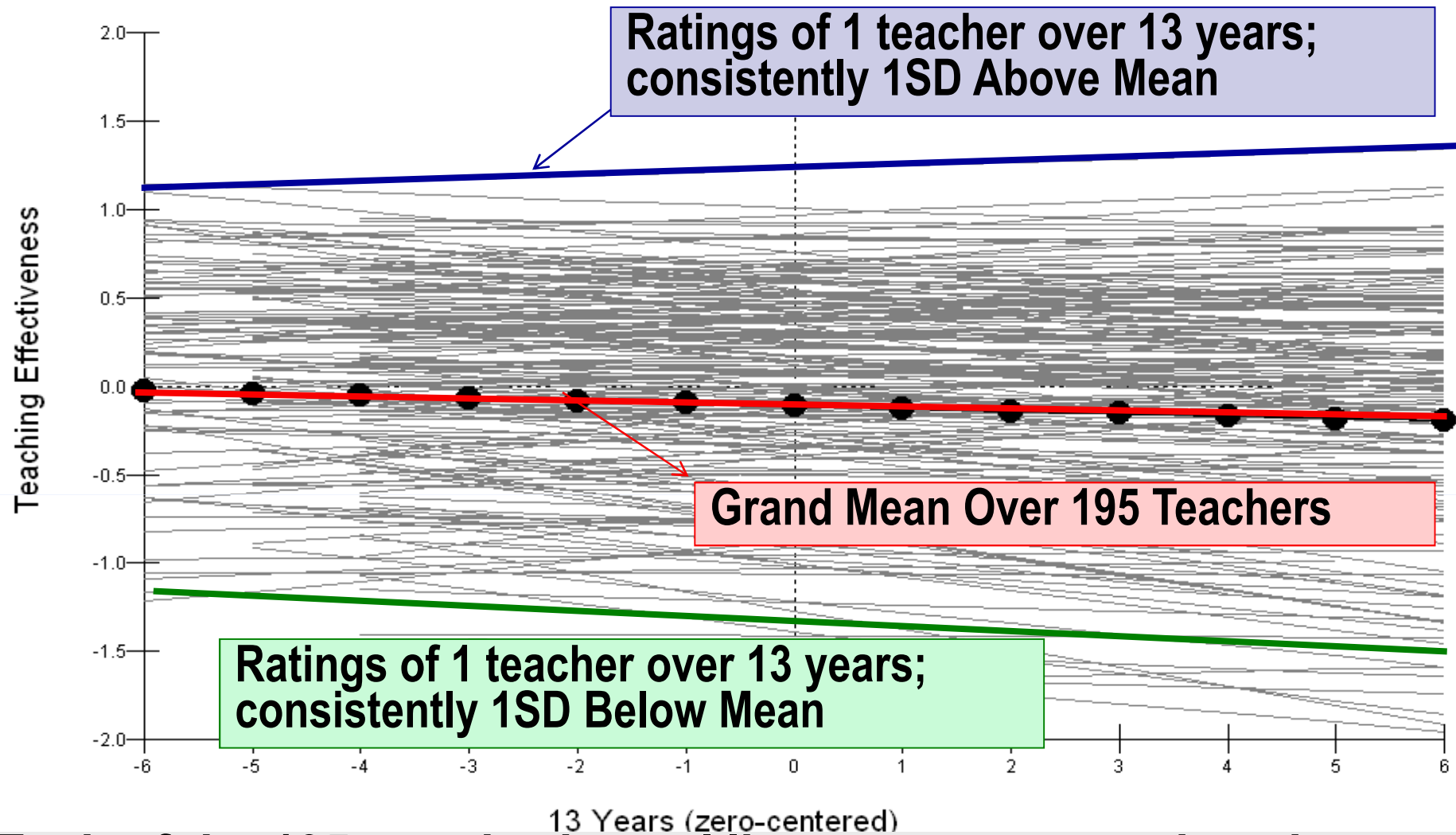- same course taught by different teachers on two different occasions?

## For Overall Instructor Ratings of:

- same instructor teaching same course on two occasions ($r$ = .72) [teacher & course effect],
- same instructor teaching two different courses ($r$ = .61) [teacher effect],
- same course taught by two different instructors ($r$ = -.05) [course effect].

**SETs primarily reflect the teacher who is doing the teaching, not the course that is being taught.**

Marsh, (2007). Students' evaluations of university teaching: A multidimensional perspective. In R. P. Perry & J C. Smart (Ed.), The Scholarship of Teaching and Learning in Higher Education: An Evidence-Based Perspective (pp.319-384). New York: Springer.

Ratings of 195 Teachers Over Time

Ratings of 1 teacher over 13 years; consistently 1SD Above Mean

Grand Mean Over 195 Teachers

Ratings of 1 teacher over 13 years; consistently 1SD Below Mean

Teaching Effectiveness

13 Years (zero-centered)

Each of the 195 grey horizontal lines represents ratings by one teacher over 13 years.

For most teachers there is no systematic increase or decrease in ratings over the 13 years.

# VALIDITY

# In Support of the Validity of SETs

SETs are positively related to many criteria of teaching effectiveness, including:

- the ratings of former students;

- student achievement in multisection validity studies;

- teacher self-evaluations of their own teaching effectiveness; and

- observations of trained observers on specific processes (e.g., teacher clarity).

Marsh, (2007). Students' evaluations of university teaching: A multidimensional perspective. In R. P. Perry & J C. Smart (Eds.), *The Scholarship of Teaching and Learning in Higher Education: An Evidence-Based Perspective* (pp.319-384). New York: Springer.

Marsh, (1982). Validity of students' evaluations of college teaching: A multitrait-multimethod analysis. Journal of Educational Psychology, 74, 264-279.

# Multisection Validity Paradigm: Validating SETs in Relation to Student Learning

- Many sections of the same course;
- Same materials in each section (e.g., course outline, textbooks, objectives, final exam);
- Random assignment (and pre-test measures);
- SETs collected prior to final exam/course grade;
- Common final exam;

Research Question: Are SETs valid in relation to objective measures of student learning (when plausible counter explanations are not viable)?

Marsh, (1982). Validity of students' evaluations of college teaching: A multitrait-multimethod analysis. Journal of Educational Psychology, 74, 264-279.

# Meta-Analysis

Cohen conducted a classic meta-analysis of multisection validity studies. Student achievement was consistently correlated with SETs:

For a subset of 41 "well-designed" studies, correlations between achievement and SETs were more substantial:

Structure (.55), Interaction (.52), Skill (.50), Overall Course (.49), Overall Instructor (.45), Learning (.39), Rapport (.32), Evaluation (.30), Feedback (.28), Interest (.15), and Difficulty (-.04).

Cohen (1987, p. 12) concluded that

"I am confident that global ratings of the instructor and course, and certain rating dimensions such as skill, rapport, structure, interaction, evaluation, and student's self-rating of their learning can be used effectively as an integral component of a teaching evaluation system."

# Teacher Self-Evaluations

In two studies, teachers evaluated their own teaching using SEEQ and were evaluated by their students:

- separate **factor analyses** of teacher and student responses identified the SEEQ factors;

- **student-teacher agreement** on all dimensions was significant (median $r$s of .49 & .45), supporting convergent validity;

- **Multitrait-multimethod analyses** indicated student/teacher agreement was specific to each SEEQ factor, supporting discriminant validity;

- **mean differences** between student & teacher responses were small (student ratings not systematically higher/lower).

> Good student/teacher agreement supports the validity of student ratings. The specificity of student/teacher agreement to each factor supports multidimensionality

Marsh, (1982). Validity of students' evaluations of college teaching: A multitrait-multimethod analysis. Journal of Educational Psychology, 74, 264-279.

# Improving Teaching Effectiveness

# Improving Teaching Effectiveness

## Many SET Feedback studies in which:

- Teachers randomly assigned to experimental (feedback) and control (no feedback) groups;
- SETs collected; Experimental Teachers get SETs feedback;
- Groups compared subsequent SETS (and other variables).

## In a meta-analysis of these studies:

- Feedback teachers .33 SD higher than control teachers
- Feedback+consultation produced much larger effects.

Marsh, (2007). Do university teachers become more effective with experience? A multilevel growth model of students' evaluations of teaching over 13 years. *Journal of Educational Psychology,* 99, 775-790.

# Two Early SEEQ Feedback Studies Using Multiple Sections of the Same Course

Study 1: results from an abbreviated survey were simply returned to teachers; impact of the feedback was positive, but very modest.

Study 2: we met with teachers in feedback group to discuss the SETs and strategies for improvement. Students in feedback group

- rated teaching effectiveness more favourably at end of the term;
- performed better on the final examination; and
- experienced more favourable affective outcomes (i.e., feelings of course mastery, plans to pursue and to apply the subject).

Study 2 was important because it demonstrated that augmented feedback improves student learning and subject affect as well as subsequent SETs.

Marsh, (2007). Do university teachers become more effective with experience? A multilevel growth model of students' evaluations of teaching over 13 years. *Journal of Educational Psychology, 99, 775-790.*

Marsh, & Overall, J U. (1980). Validity of students' evaluations of teaching effectiveness: Cognitive and affective criteria. *Journal of Educational Psychology, 72, 468-475.*

# Prototype SEEQ Feedback/Consultation Intervention

Teachers randomly assigned to groups.

At T1 (middle of semester 1), T2 (end of semester 1) and T3 (end of semester 2) all teachers:
- evaluated themselves & rated importance of each SEEQ factor;
- were evaluated by students on SEEQ;

At T2 Feedback Teachers selected *target* SEEQ factors that:
- were important to the teacher (teacher self-evaluations);
- had low ratings (needed improvement);
- were "appropriate areas to target improvement efforts."

**Teaching idea packets** given to teachers for each targeted SEEQ factor. Each packet contained up to 40 strategies (based on interviews with outstanding teachers). Teacher (with consultant) selected a few strategies to implement for each target SEEQ factor.

Control teachers received no feedback until end of the study.

# Results For Feedback Teachers

- Because teachers only targeted one or a few scales, interpretations of overall ratings most straight forward; effects significant for all 4 overall ratings (effect sizes .4 to .5).

- The feedback group had higher ratings for all 12 SEEQ scores; 8 were statistically significant.

- Now, lets see how the intervention worked with the target scales (that teachers chose for the intervention) compared to non-target scales.

Marsh, , & Roche, L.  (1993). The use of students' evaluations and an individually structured intervention to enhance university teaching effectiveness. American Educational Research Journal, 30, 217-251.

Results

**Intervention Post-test ratings were all much higher; Target scales were now similar to non-Target scales. The intervention improved target scales much more than non-target scales**

**Control Group Post-test ratings of Target scales were still much lower than non-target scales**

Rating

7.2

7

6.8

6.6

Endterm Feedback

**At Pretest the ratings of Target scales were much lower than non-target scales for all groups – part of**

3

— ▲ — Endterm grp: Target Areas

— ● — Control Group: Target Areas

— △ — Endterm grp: Nontargeted

— ○ — Control Group: Nontargeted

**Consistent with the rationale for the study, ratings of targeted scales improved substantially relative to nontargeted areas for experimental groups, but not for the control group.**

# Discussion

The most important results of the investigation were to provide varying degrees of support for a priori predictions that:

- SEEQ feedback and the feedback/consultation provided an effective means of improving university teaching;

- Effects stronger for the initially less effective teachers;

- In support of multidimensional SEEQ perspective, improvement largest for targeted SEEQ scales;

- Important for teachers to specifically target particular scales.

- Teaching packets: even if teachers motivated to improve their teaching, they apparently do not know how to do so. Need concrete strategies to facilitate teaching improvement efforts.

However, few universities implement teaching improvement programmes as part of the collection of SETs even though clear evidence that they work. I would like to pursue a large-scale test of this process to include all UK universities.

# Overall Summary Conclusions

In conclusion, let me return to my original conclusion that SETs based on the teacher as the unit of analysis are:

- **multidimensional;**

- **reliable and stable;**

- **primarily a function of the teacher who teaches a course rather than the course that is taught;**

- **valid in relation to a variety of indicators of effective teaching;**

- **relatively unaffected by a variety of variables hypothesized as potential biases; and**

- **seen to be useful by students for use in course selection, by administrators for use in personnel decisions, by teachers as feedback about teaching**

# Using Student Ratings To Benchmark Universities

# The National Student Survey: A Multilevel Analysis of Discipline Effects

**Cheng & Marsh (2010). UK National Student Survey: Are differences between universities and courses reliable and meaningful. *Oxford Rev of Educ*, 36, 6, 693-712.**

# NSS Overview

| 2005 | 2006 |
|---|---|
| N = 171, 290 (285, 445) | N = 157, 341 (278, 796) |
| Female = 97,356<br>Male = 73, 964 | Female = 93,704<br>Male = 63, 667 |
| Mean age = 22.0 yrs | Mean age = 21.2 yrs |
| 140 Universities | 144 Universities |

- At the level of the individual student, responses had good psychometric properties, but

- How much variance explained by university & discipline-within-university groups ?

- How reliable are NSS responses? How is this related to sample size at different levels?

# Design: Variance Components

## 3-Level Model

- L1 = students,
- L2 = Departments/discipline groups,
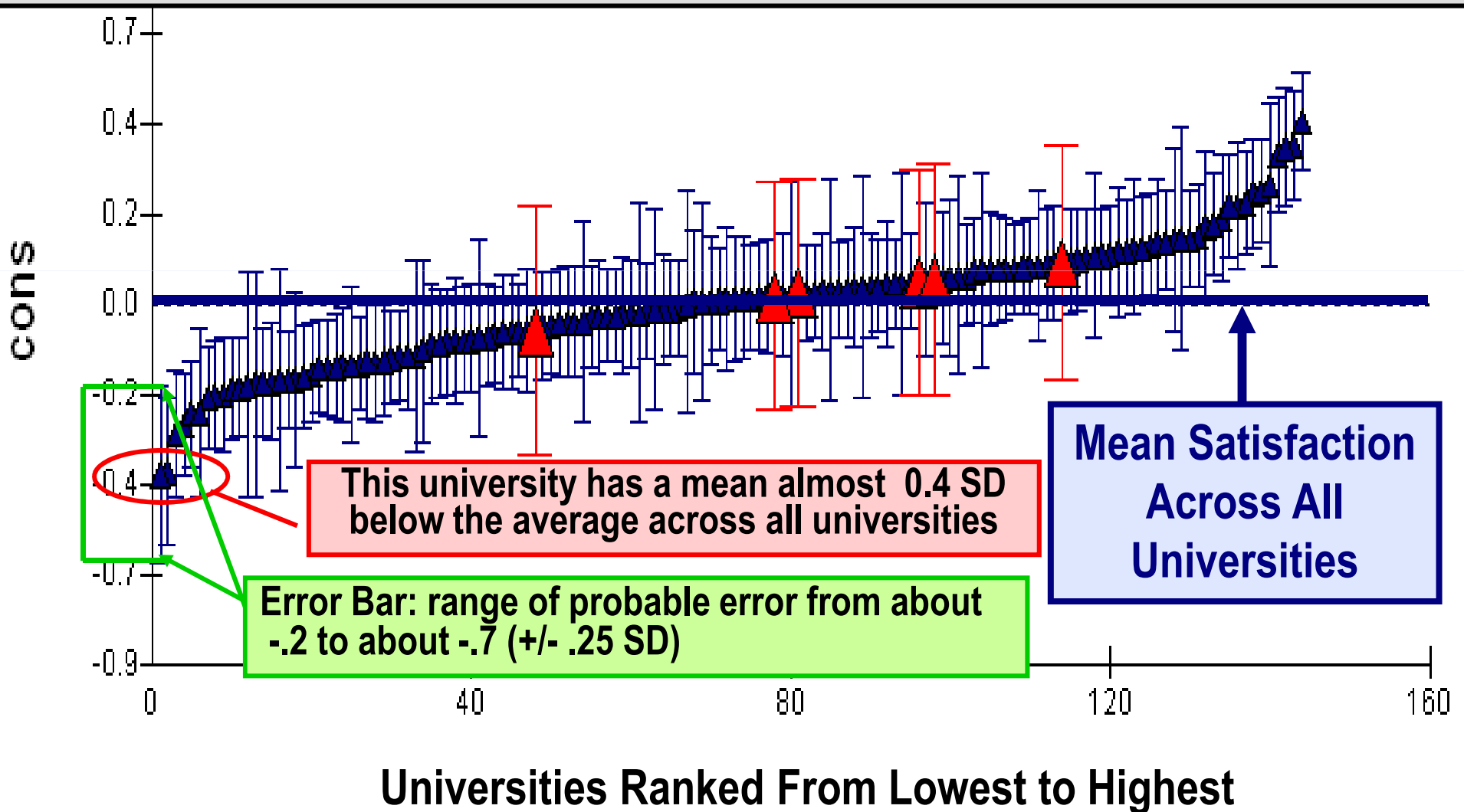- L3 = university

## Variance component Models

- How much variance is explained by each level?
- How much is this changed by controlling fixed effects (student characteristics & discipline)?

## Means & Probable Error (Error Bars)

- For each group (university or department) there is an mean level of satisfaction and a range of probable error (error bar around the mean)
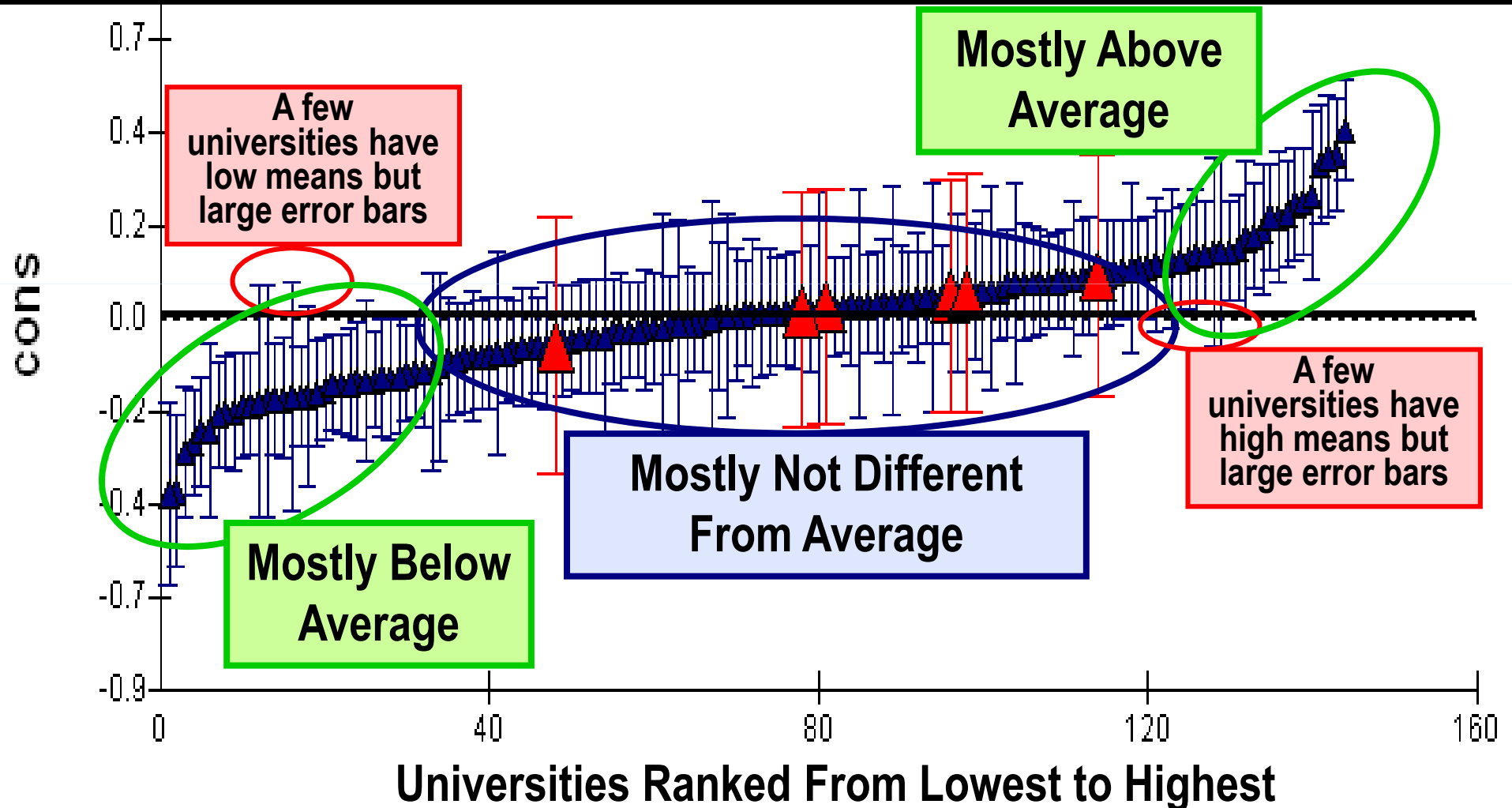
# Differences Between Universities: Caterpillar Plots

Each triangle is the mean satisfaction ratings for one university. The vertical line that goes above and below the mean is an error bar (range of probable error); longer bars represent more error.

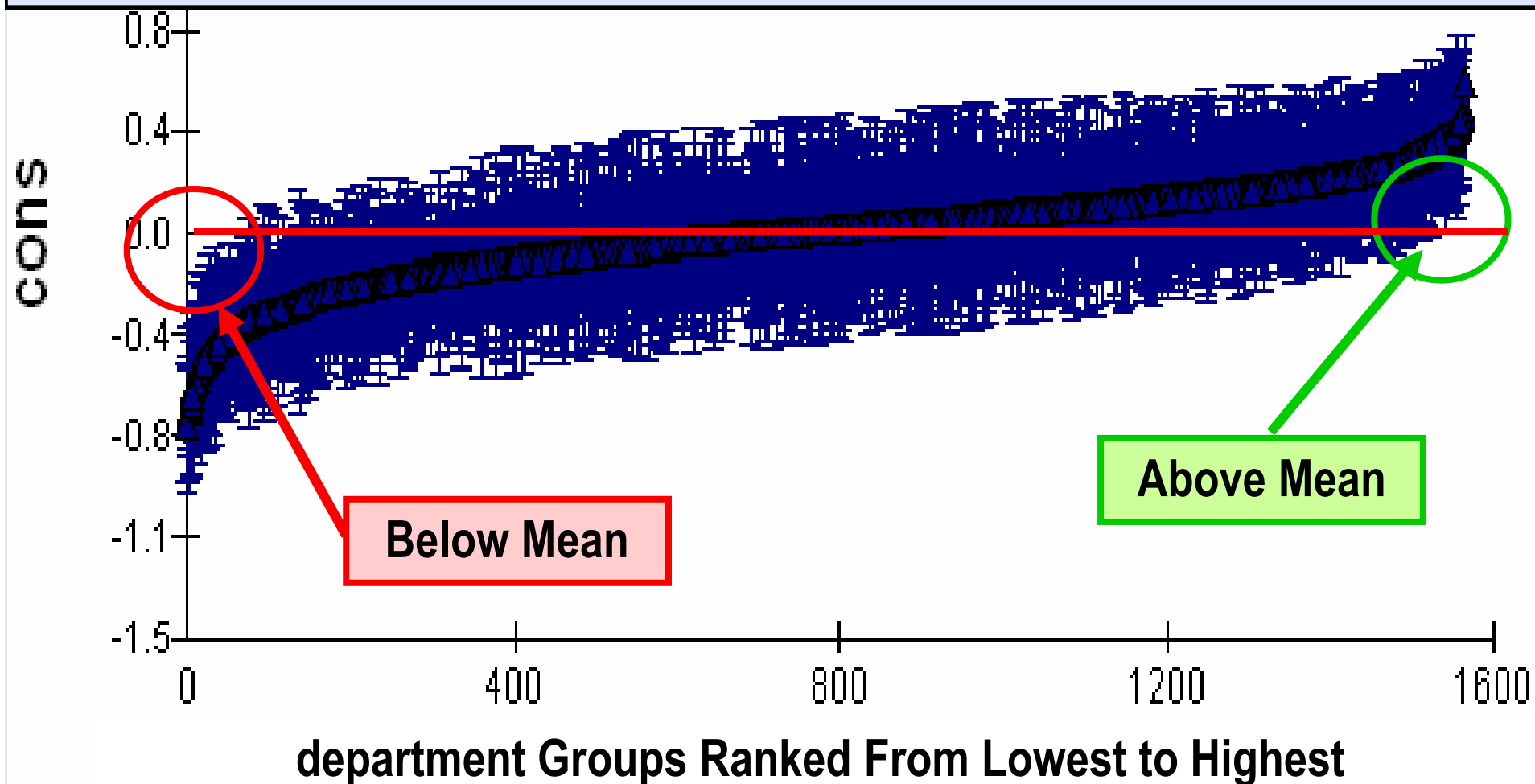This university has a mean almost 0.4 SD below the average across all universities

Error Bar: range of probable error from about -.2 to about -.7 (+/- .25 SD)

Mean Satisfaction Across All Universities

Universities Ranked From Lowest to Highest

For 19 Discipline Categories **1565** Discipline-Within University Groups.

Differences between groups are larger than differences between universities (more variance explained).

However, the error bars are VERY LARGE so only a few extreme department groups differ significantly from mean.

cons

Below Mean

Above Mean

department Groups Ranked From Lowest to Highest

# Australian Course Evaluation Questionnaire (CEQ)

Marsh, Ginns, Morin, Nagengast, Martin, (in press). The Course Evaluation Questionnaire (CEQ): Use of student ratings to benchmark Australian universities. Journal of Educational Psychology
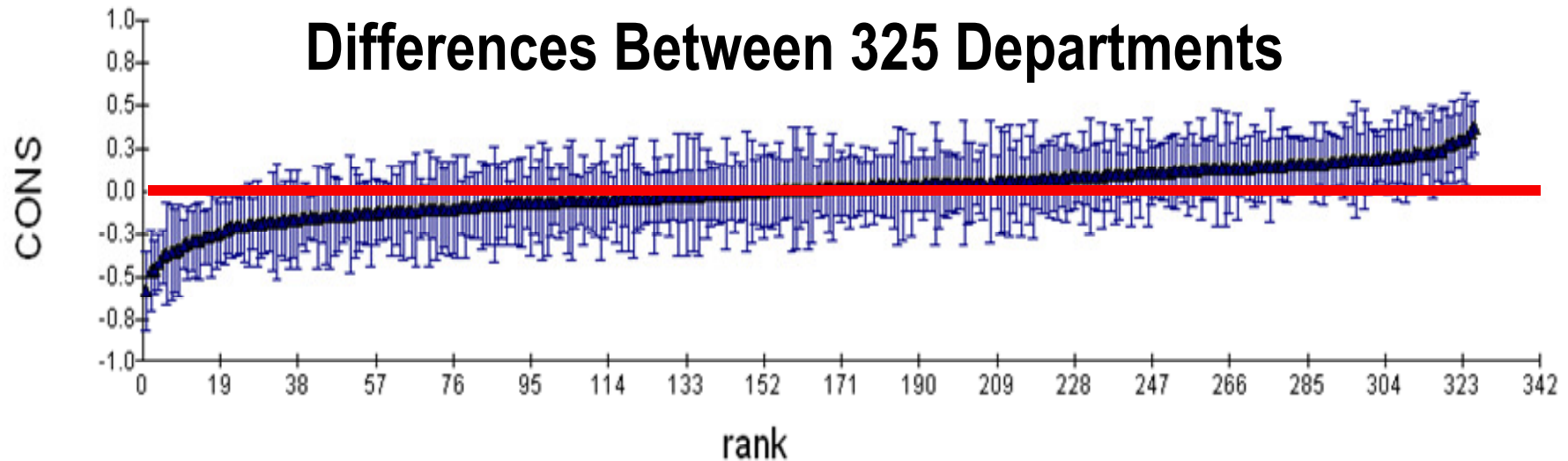
# Benchmarking Australian Universities

- Australian government & universities cooperate to collect standardized data of many kinds that are used to compare universities – a benchmarking exercise.

- In Australia, the Course Experience Questionnaire (CEQ) is used to compare undergraduate teaching in different universities – to benchmark teaching effectiveness.

- The Australian CEQ is the oldest of the university/department experience type instruments and one basis for the NSS. In his review, Richardson indicated that (in 2005) it was the only one to have been widely researched.

Marsh, , Ginns, P., Morin, A. J. S., Nagengast, B., Martin, A. J. (in press). The Course Evaluation Questionnaire (CEQ): Use of student ratings to benchmark Australian universities. Journal of Educational Psychology

# Results Based on Australian CEQ Responses (2001: 44,000 students, 45 universities)



**Differences Between 325 Departments**

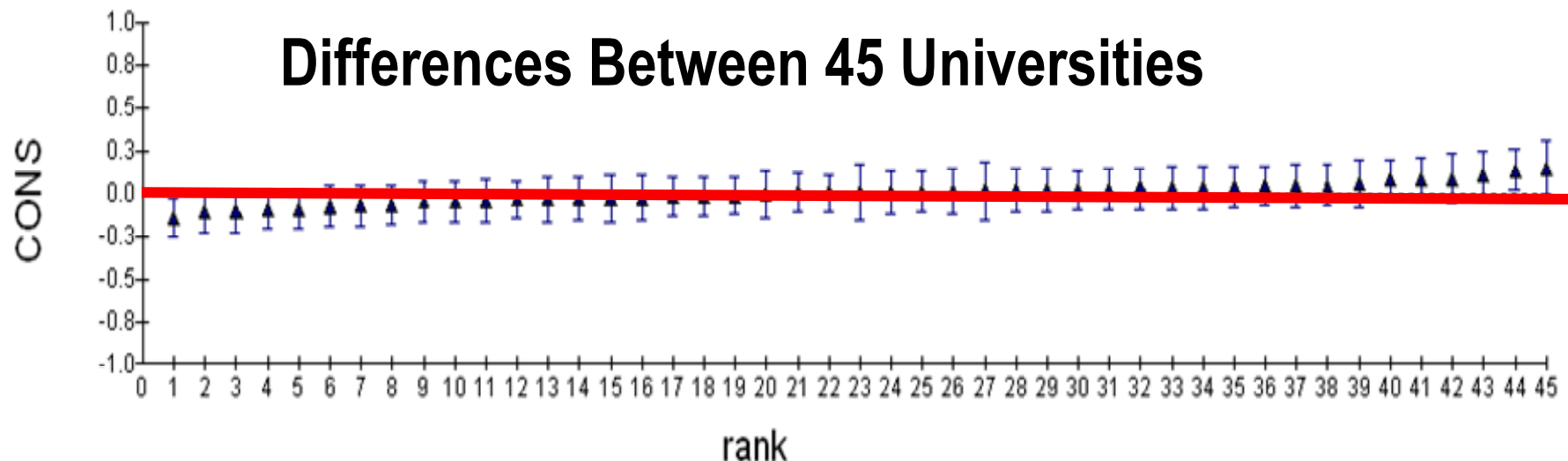**Differences Between 45 Universities**

Marsh, , Ginns, P., Morin, A. J. S., Nagengast, B., Martin, A. J. (in press). The Course Evaluation Questionnaire (CEQ): Use of student ratings to benchmark Australian universities. Journal of Educational Psychology

# Australian PhD Students' Evaluations of Supervision: Benchmarking Universities With PREQ

# Differences Among 35 Universities: Study 2
## (0.4% of Var Explained)



**total score**

**Mean Rating Across all 32 Universities**

**Mean Rating (x) and Range of Probable Error for Each University**

resid normal score

university

Marsh, , Rowe, K., Martin, A. (2002). PhD students' evaluations of research supervision: Issues, complexities and challenges in a nationwide Australian experiment in benchmarking universities. Journal of Higher Education, 73 (3), 313-348.

# Validity/Usefulness of PREQ Responses

PREQ responses are completely unreliable, so they must also be invalid for purposes of differentiating between universities. PREQ responses were unrelated to :

- Research Productivity (publications & grants);
- Number Australian PhD Student Scholarships;
- Attrition Rates.

We concluded that PREQ responses were unlikely to be useful for any of the purposes for which they were designed (including benchmarking and improving PhD programmes)

Marsh, , Rowe, K., Martin, A. (2002). PhD students' evaluations of research supervision: Issues, complexities and challenges in a nationwide Australian experiment in benchmarking universities. Journal of Higher Education, 73 (3), 313-348.

# Summary
# and
# Discussion

For Purposes of Comparison consider data from earlier SEEQ Longitudinal study of 195 Different Teachers (Consistency across an average of 30 Classes per teacher over 13 Years)

Above Average

Below Average

Ratings of individual teachers highly differentiated relative to NSS ratings (large differences between teachers relative to probable error; higher % significantly above & below the mean)

195 Individual Teachers Ranked from Lowest to Highest on Overall Rating

# NSS vs. Student Evaluations of Teaching

1.  Students' evaluations of teaching are good at reflecting teaching effectiveness of individual teachers. Importantly, when coupled with appropriate enhancement interventions, they lead to improved teaching. May also be useful for personnel decisions and student choice of teachers, but not very good a benchmarking universities and department groups.

    SETs are apparently not good at benchmarking whole universities of department – or even classes independent of the teacher who teaches them

2. NSS, CEQ & PREQ  responses provide limited evidence about small differences in student satisfaction with educational experience at the university or department  level.

    NSS, CEQ & PREQ studies have focused almost exclusively on reliability-type issues. Unlike SET research, there is almost no rigorous validity research asking whether the limited differentiation are meaningfully related to other indicators of effect educational effectiveness or useful in leading to the educational improvement.

# Summary
## NSS vs. Student Evaluations of Teaching

In summary, there is no basis for using NSS-type approaches instead of SET-type approaches (or vice versa). They have different purposes

However, there is limited rigorous evidence that NSS-type responses are reliable, valid, or useful for any purposes.

The substantial SET literature in support of their reliability, validity, and usefulness SHOULD NOT be used to justify NSS-type ratings.

NSS-type ratings should only be used with extreme caution for benchmarking purposes: comparisons of ratings across different universities, different departments within the same university or the same department across different universities. Comparisons should be qualified by estimates of probable error as in caterpillar plots.

# References

Cheng, J. H. S. & Marsh, (2010). UK National Student Survey: Are differences between universities and courses reliable and meaningful. Oxford Rev of Educ, , 36, 6, 693-712.

Marsh, , Ginns, P., Morin, A. J. S., Nagengast, B., Martin, A. J. (in press). The Course Evaluation Questionnaire (CEQ): Use of student ratings to benchmark australian universities.Journal of Educational Psychology

Marsh, (2007). Students' evaluations of university teaching: A multidimensional perspective. In R. P. Perry & J C. Smart (Ed.), The Scholarship of Teaching and Learning in Higher Education: An Evidence-Based Perspective (pp.319-384). New York: Springer.

### Also see

Marsh, , Muthén, Asparouhov, Lüdtke, Robitzsch, Morin & Trautwein. (2009). Exploratory Structural equation modeling, integrating CFA and EFA: Application to students' evaluations of university teaching. Structural Equation Modeling,16, 439-476.

Marsh, , & Hoceva (1991). Multidimensional students' evaluations of teaching effectiveness: Factor structure stability across academic discipline, instructor level, and course level. *Teaching and Teacher Education, 7,* 9-18.

Marsh & Roche (1997). Making students' evaluations of teaching effectiveness effective. *Am Psychol, 52,* 1187-1197.

Marsh (1987). Students' evaluations of university teaching: Res findings, methodological issues, and directions for future research. *International J of Educ Res, 11,* 253-388 (whole issue).

Marsh (1984). Students' evaluations of university teaching: Dimensionality, reliability, validity, potential biases and utility. *J of Educ Psych, 76,* 707-754. (Invited Lead Article).

Marsh (2007). Students' evaluations of university teaching: A multidimensional perspective. In R. P. Perry & J C. Smart (Ed.), The Scholarship of Teaching and Learning in Higher Education: An Evidence-Based Perspective (pp.319-384). New York: Springer.

Marsh, Rowe & Martin (2002). PhD students' evaluations of research supervision: Issues, complexities and challenges in a nationwide Australian experiment in benchmarking universities. *Journal of Higher Education, 73 (3),* 313-348.

Marsh, & Roche. (2000). Effects of grading leniency and low workloads on students' evaluations of teaching: Popular myth, bias, validity or innocent bystanders? *Journal of Educational Psychology, 92,:*202-228.

Marsh & Dunkin (1997). Students' evaluations of university teaching: A multidimensional perspective. In R. P. Perry & J C. Smart (Ed.), *Effective teaching in higher education* (pp. 241-320). New York: Agathon.

Marsh, , & Roche (1993). The use of students' evaluations and an individually structured intervention to enhance university teaching effectiveness. *American Educational Research Journal, 30,* 217-251.

Ginns, Marsh, Behnia, Cheng & Scalas (2009). Using postgraduate students' evaluations of research experience to benchmark departments and faculties: Issues and challenges. *British Journal of Educational Psychology, 79,* 577-598.

Marsh (2007). Do university teachers become more effective with experience? A multilevel growth model of students' evaluations of teaching over 13 years. *Journal of Educational Psychology, 99,* 775-790.

Marsh & Roche (1999). Student Evaluations of Teaching Effectiveness (SETs): Objective Learning, Usefulness, Validity, Bias, and Further Research. *American Psychologist, 54,* 517-518.

Marsh & Bailey (1993). Multidimensional students' evaluations of teaching effectiveness: A profile analysis. *Journal of Higher Education, 64,* 1-18.

Marsh & Hocevar(1991). Students' evaluations of teaching effectiveness: The stability of mean ratings of the same teachers over a 13-year period. *Teaching &Teacher Education, 7,* 303-314.

Marsh, & Hocevar, D. (1984). The factorial invariance of students' evaluations of college teaching. *American Educational Research Journal, 21,* 341-366.

Marsh (1982). SEEQ: A reliable, valid, and useful instrument for collecting students' evaluations of university teaching. *British Journal of Educational Psychology, 52,* 77-95.

Marsh, (1982). Validity of students' evaluations of college teaching: A multitrait-multimethod analysis. *Journal of Educational Psychology, 74,* 264-279.

Marsh, (1980). The influence of student, course and instructor characteristics on evaluations of university teaching. *American Educational Research Journal, 17,* 219-237.

Marsh, & Overal. (1980). Validity of students' evaluations of teaching effectiveness: Cognitive and affective criteria. *Journal of Educational Psychology, 72,* 468-475.

Overall & Marsh (1979). Midterm feedback from students: Its relationship to instructional improvement and students' cognitive and affective outcomes. *Journal of Educational Psychology, 71,* 856-865.

# SEEQ Factor Structure

**How well does the SEEQ factor structure generalize across different disciplines ?**

For a sample of 25,000 courses (~1million students), I conducted factor analyses for:

- the total sample and

- each of 21 subsamples of unique groups of teachers who varied in terms of academic discipline (e.g., psychology, engineering, etc.) and level (e.g., undergraduate and graduate courses).

Factor analyses for total sample and each of the 21 subsamples all identified the same 9 SEEQ factors. The SEEQ factor structure is very robust.

# NSS 22-item Instrument
## (6 specific factors & overall rating item)

- **Teaching**: "staff are good in explaining things";

- **Assessment/Feedback**: "Assessment arrangements and marking have been fair";

- **Support**: "I have received sufficient advice and support with my studies";

- **Organisation/Management**: "The timetable works efficiently as far as my activities are concerned";

- **Resources**: "The library resources and services are good enough for my needs";

- **Personal development**: "The course has helped me to present myself with confidence";

- **Overall satisfaction**: "Overall, I am satisfied with the quality of the course"

**Again, at the level of the individual student, responses had reasonable psychometric properties.**