

**Making Assignments Count and Other Strategies for  
Increasing Critical Thinking in Political Theory**

Journal:	<i>Journal of Political Science Education</i>
Manuscript ID	JPSE-2017-0036.R2
Manuscript Type:	Scholarship of Teaching and Learning
Keywords:	Critical thinking, political theory, Undergraduate education, domain specificity, Assessment

SCHOLARONE™  
Manuscripts

Review Only

# Making Assignments Count and Other Strategies for Increasing Critical Thinking in Political Theory

## Abstract

Political theory lags behind other subfields in political science in rigorously testing what helps foster critical thinking (CT). Yet some of the greatest temptations to engage in motivated reasoning can be found in normative political contexts. This study uses multiple regression analysis to explore 9 semesters of data from an introductory course in political theory. Three results stand out. 1. As the stakes of an assignment decrease, so do CT scores. 2. As the number of course assignments increase, CT scores fall. 3. When preparatory exercises matter, theory-specific CT exercises outperform generic ones. A theory of student rationality is put forth. Implications for course design, program assessment, and future research are discussed.

## Introduction

Colleges and universities in the United States often see themselves as places where students develop their critical thinking (CT) skills<sup>1</sup>. In political science, multiple articles have outlined and evaluated strategies for promoting CT in the classroom. These include: opportunities for critical self-reflection through writing (Cavdar and Doe 2012), engaging with social science methodology (Marks 2008, Olsen and Statham 2005), group discussion (Blings and Maxey 2016, Williams and Lahman 2011), group debate (Omelicheva 2007, Oros 2006), and the avoidance of political controversy (Fitzgerald and Baird 2011).

---

<sup>1</sup> According to a report from the Association of American Colleges and Universities, the ability for critical analysis and logical thinking is mentioned in fully 25% (75 of 301) of mission statements of the Princeton Review's "331 Best Colleges and Universities in the United States" that express an educational vision. This places critical thinking in fifth place on the list of most common educational themes in university mission statements.(Gaff and Meacham 2014)

1  
2  
3  
4  
5  
6  
7  
8  
9 Although introductory courses, American government courses, and courses in  
10 comparative politics are well-represented in this literature, strategies for developing CT  
11 in political theory have not been empirically assessed. Given evidence that CT may have  
12 important subject specific components (McPeck 1981, Smith 2002) and moreover that  
13 moral reasoning skills probably differ from other kinds of reasoning skills (Mason 2007,  
14 Phillips and McMillian 2010) a gap in the literature exists with respect to what promotes  
15 CT in the subfield of political theory. By leveraging data collected from 451 essays over  
16 the course of 9 semesters in a required junior-level course in political theory, this paper  
17 begins to fill this gap.  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

### 33 **Why Critical Thinking?**

34  
35 The ability to think critically, that is to say the ability of individual citizens to construct  
36 and evaluate conclusions from available evidence and assumptions (Williams and Worth  
37 2001), is a *sine qua non* of a functioning democracy and economy. The idea that  
38 democracy fares roughly as well as the ability of its citizens to critically evaluate the  
39 rhetoric of its leaders is one that can be traced as far back as Plato. The idea that  
40 analytical skill (a component of CT) is important to the job performance of employees is  
41 also well-cemented in today's economy, including among American employers (NACE  
42 Job Outlook, 2016)  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7 The role of colleges and universities in preparing individuals for their roles in these  
8 contexts is hardly disputed either. Faculty have long rated CT among their most cherished  
9 pedagogical objectives (Bok 2006; Paul 1991; Resnick and Peterson 1991) and  
10 universities around the country have focused their attention on developing CT for some  
11 time now.  
12  
13  
14  
15  
16

17  
18  
19 Research projects, such as project CAT (Stein et. al 2010), the California Critical  
20 Thinking Disposition Inventory (Facione, Facione, and Sanchez 1994), or the  
21 Washington State University Critical Thinking Project (Kelly-Riley. et. al. 2008) have  
22 developed assessment tools to try and measure CT. Exams like the California Critical  
23 Thinking Skills Test (CCST) and the Collegiate Learning Assessment (CLA) are  
24 routinely administered by universities to check for progress in CT in their graduates  
25 (though not without controversy: see Klein, Shavelson, and Bolus 2007).  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35

36  
37 Political theorists also want to develop CT skills in their students. In Matthew Moore's  
38 2011 survey "How (and What) Political Theorists Teach: Results of a National Survey,"  
39 fully 96.3% of political theorists surveyed thought that inculcating CT skills in their  
40 courses was either "important" or "very important." This is unsurprising since a wide  
41 range of viewpoints in political theory agree on the importance of CT.  
42  
43  
44  
45  
46  
47  
48  
49

50 Critical theorists, Feminists, and Marxists want to peel back the illusions of social norms,  
51 power, gender, race, and class. To do this, practitioners must be capable of critically  
52 reflecting on the assumptions, categories, and processes that characterize the modern  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 world. Liberals and perfectionists cherish CT because of its connection to autonomy: the  
7 capacity of individuals to set ends for themselves. On these views, the twin-projects of  
8 democracy and individual liberty remain radically incomplete if citizen choices are the  
9 product of unquestioning acceptance of existing cultural norms or power structures.  
10  
11  
12  
13

## 14 15 16 17 18 **Operationalizing Critical Thinking** 19

20 Despite near-universal agreement on its importance, there are competing definitions of  
21 CT. Pascarella and Terenzini (2005) analyze different conceptualizations and find the  
22 following broad regions of agreement. “Most attempts to define and measure critical  
23 thinking operationally focus on an individual’s capability to do some or all of the  
24 following: identify central issues and assumptions in an argument, recognize important  
25 relationships, make correct references from the data, deduce conclusions from  
26 information or data provided, interpret whether conclusions are warranted based on given  
27 data, evaluate evidence of authority, make self-corrections, and solve problems.” (156)  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39

40 CT in this study was measured using a modified version of the Washington State  
41 University critical thinking rubric. Over a period of 9 semesters, 451 1,000-1,500 word  
42 essays were assessed in a required political theory course aimed at junior political science  
43 majors at a midsize southern public university.  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7 Six categories, broadly reflective of the consensus articulated above, are independently  
8 scored. The Focus score evaluates how well a student conceptualized the question being  
9 asked. The Interpretation score looks at how fairly students interpreted their sources in  
10 those essays. The Thesis score assesses the logical consistency of the argument presented.  
11 The Objections score reflects how well a student understood and responded to the  
12 arguments raised by alternative perspectives. The Evidence score measures whether  
13 appropriate sources, facts, and normative principles were marshalled in support of the  
14 thesis. Finally, the Conclusion score measures the ability to assess the significance and  
15 implications of the argument made (for public policy, political theory, and so on)<sup>2</sup>.  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27

28 The choice was made to report results for each dimension of CT separately as well as an  
29 aggregate measure. There is no a priori reason to believe that the interventions or course  
30 designs that (for example) help students better select evidence for their argument also  
31 help them with their logical reasoning or their ability to interpret the meaning of the  
32 question. My experience using this grading rubric is also that students may write very  
33 cogent papers (high Thesis scores) that are completely off topic (low Focus scores). They  
34 may develop their perspective very well (with respect to logical consistency and  
35 evidence) and yet ignore obvious objections. This variation in subcategory scores seems  
36 worth preserving.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

---

53 <sup>2</sup> A more detailed discussion of the WSU grading rubric and adaptations made for political theory  
54 is available in Appendix A.  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Several checks on the reliability of the CT measures were conducted as a post-test using a random sample of 24 essays scored by 3 different graders. A summary of the reliability checks is available in Appendix B.

## Stimulating Critical Thinking

The data set permits us to test the relative effectiveness of several different types of exercises for stimulating CT including critical review essays (CREs), critical thinking exercises (CTEs), and logic game exercises (LGEs).

### Logic Games

An important component of CT is the ability to draw accurate inferences from available information. Plausible induction and valid deduction are central to the application of CT to any given context. In that sense, logic is a kind of “generic” CT skill, since the rules of formal logic can be applied successfully in any context. A popular textbook on logic claims, “the study of logic is one of the best ways to increase students’ skill in critical thinking” and “The study of logic increases one’s ability to understand, analyze, evaluate, and construct arguments.” (Howard-Snyder, Howard-Snyder, and Wasserman 2012: xv)

Because of the question of domain specificity raised earlier, whether CT skills (or any kind of reasoning skills) are easily transferred from one field of inquiry to another is still

1  
2  
3  
4  
5  
6 a topic of debate. On the transferability side of the debate see Ennis (1989) and Bailin  
7 (2002). On the “domain specificity” side see Halpern (2001) and Van Gelder (2005).  
8  
9

10  
11  
12 During three semesters, students practiced LGEs of the kind one finds in puzzle books or  
13 the Law School Admissions Test (LSAT) every two weeks for one class period. At the  
14 beginning of the semester, students were subjected to a short introduction to formal logic.  
15 Before the first LGE, students were walked through some examples of logic puzzles and  
16 common types of operations seen in those puzzles (particularly *modus ponens* and *modus*  
17 *tollens*) After most students had solved the games, the class as a whole was asked to  
18 assess of the potential soundness of the conclusions reached.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

30 Though LGEs incorporated themes from political theory, they did not provide  
31 opportunities for moral reasoning but instead consisted of a series of deductive logical  
32 operations that led to a unique set of possible conclusions. As such, they are very unlike  
33 the more nuanced, probabilistic, self-aware, and contextual kind of reasoning  
34 characteristic of high quality essays in political theory. Still, if CT is a generic skill that  
35 can be both practiced and applied to political theory, then LGEs should help with the  
36 rigor of the essays if nothing else<sup>3</sup>.  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

### 48 **Critical Review Essays and Critical Thinking Exercises**

49 Two other interventions stand out as domain-specific counterparts to the attempt to  
50 cultivate generic CT skills: critical review essays (CREs) and critical thinking exercises  
51  
52  
53

---

54 <sup>3</sup> The thesis score stands out as the most likely candidate for improvement.  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6 (CTEs). For the CREs, each student received an anonymous essay written by a classmate  
7  
8 that semester and was tasked to write a review essay outlining its strengths and  
9  
10 weaknesses.  
11

12  
13  
14  
15 Though not directly an opportunity for self-reflection as suggested by Çavdar and Doe  
16  
17 (2012), CREs ask students to revisit a topic they have written on and see it through the  
18  
19 eyes of another student. This indirectly provides them with an opportunity for self-  
20  
21 examination. One might expect the Interpretation, Objections, and Evidence subscores  
22  
23 on subsequent essays to benefit from the exercise.  
24

25  
26  
27  
28 Like CREs, CTEs are also domain specific. CTEs used in this study are a series of  
29  
30 scenarios depicting moral and political decision-making that the students were asked to  
31  
32 logically analyze. Students were then asked (as homework) to identify the moral  
33  
34 principles in play and think about the implications of their opinions. Other types of  
35  
36 exercises included the analysis of hidden premises in pared down arguments found in  
37  
38 internet memes. The students were asked to identify flaws in the logic of the argument  
39  
40 proposed, imagine what would have to be true for the argument to be sound, reconstruct  
41  
42 the argument in a more plausible way, and assess the implications. One would expect  
43  
44 CTEs to have an impact on all the dimensions of critical thinking.  
45  
46  
47  
48

### 49 **Course Designs**

50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Aside from comparing the relative impact of rival strategies for promoting CT in the course, the paper also has the opportunity to examine the impact of different course designs.

One plausible route to improving CT is to try to increase the overall quantity of CT assignments. If proficiency is a function of practice, then more practice is likely better. In political science, this “scaffolding” approach is endorsed (in various forms) by Fitzgerald and Baird (2011), Ewell and Rodgers (2014) and Mulcare and Schwedel (2015). On the other hand, if students only have so many hours to devote to a course, it is possible that a surfeit of lower stakes assignments causes less careful and well thought out reasoning for each assignment – and especially for the CT essays.

Another feature of the course design instructors have control over is the weight of each assignment in the final grade. The literature on test taking and student effort suggests that the stakes of the test impact scores. If a small number of points are in play, many students reduce their effort and test and writing assignments scores fall. (Barry et. al. 2010, Wise and Demars 2005, and Elbow 1997). Sewell (2004) also found that dropping a quiz or homework grade negatively affected performance on a final exam. But CT scores are often thought of as measures of ability, not effort, so a question arises as to whether CT scores are sensitive to the stakes in play. Both the percentage value of an essay in the final grade and whether or not it was announced that the worst essay would be dropped are tested.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Other course design features tested include: the presence or absence of exams and whether or not course slides were made available to students.

## Key Control Variables

A group of control variables was selected from several meta-analyses of CT studies in higher education (Abrami. et. al 2015, Pithers and Soden 2000, Tiruneh, Verbugh, and Elen 2014) as well as from plausible theoretical assumptions about student behavior. These variables seek to account for three principal sources of variation in CT scores: the individual characteristics of the student, the specific characteristics or context of the assignment, and features of the course-design itself – the latter two being the most pedagogically interesting since they can be manipulated by the instructor.

Control variables collected include: GPA, ACT scores, the percentage of online readings accessed, the percentage of online slides accessed, attendance grades, gender, age, the number of prior essays, and the number of pages of reading since the last essay. Cross-correlations for all these variables can be found in Table 1 and Table 2. Descriptive statistics, collection methods, and a brief discussion of validity and reliability issues pertaining to each variable are discussed in Appendix C.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Finally, when levels of CT are being considered, it is important to adjust scores for topic difficulty. Each of the 23 topics given was therefore assigned a dummy variable so as to produce topic insensitive estimates.

## Methods

If CT is subject-specific, then the best predictors in the literature from other disciplines (or even other subfields of political science) might not all apply for political theory. A backwards elimination method and a (forward) stepwise regression method were therefore used to create regression specifications that combine the best predictors available from the pool of independent variables. Because of the similarity of the results, only the backwards elimination method is presented<sup>4</sup>.

The only theoretical stricture imposed on the stepwise procedure was the inclusion of the bloc of topic dummy variables.

Since the goal of the study is not to obtain the most parsimonious model in order to leverage a small n relative to a large number of potential regressors, the stepwise selection rules adopted were relatively inclusive. Each model was constructed using a p value cut off of .10<sup>5</sup>.

---

<sup>4</sup> The forward stepwise specifications are available upon request. The differences are minor.

<sup>5</sup> Other decision criteria were considered, including minimum AICc and BIC, but were rejected on the grounds that parsimony and fit were less important than determining whether there was some statistically significant that relationship could be detected for the purposes of future investigation.

## Results

### **Preparatory Assignments, Volume of Assignments, and Manipulating the Stakes**

The view that the stakes matter for CT scores fares well in this investigation. The aggregate measure (Table 5) and all the subcategory scores bar one are negatively associated with the stakes of the assignment (Table 3; Table 4). Announcing that the worst essay grade from the semester will be dropped is associated with a further decline in CT scores on the Interpretation category.

The number of CT assignments in a semester is also negatively associated with the aggregate measure of CT (Table 5) as are the measures for selecting evidence and assessing implications (Table 3). In the latter two cases, effect sizes are important, with the number of assignments accounting for a greater quantity of variation than even GPA or ACT scores.

Though more assignments overall are apparently not a path to better scores, subject specific CT assignments fare better than generic ones. CREs help some subcategories (Objections, Evidence, Conclusions) as do CTEs (Conclusions, Objections). The Prior Essays variable was also positive and significant for the Focus score. LGEs, on the other hand, fare poorly: always carrying a negative sign and sometimes appearing to be significantly associated with lower CT scores (Evidence, Conclusion).

### **Essay Topics and Prompts**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

If we focus on the subcategories, we can see that topics account for between 8.3 and 15.7% of the variation in student CT scores, by far the largest single source of variation<sup>6</sup> (Table 3; Table 4). The magnitude of the effects for individual topics was not usually trivial, routinely pushing scores 20 or 30 points or more in either direction, though the sample sizes for each topic were sometimes small.

### Other results

As expected, a student's GPA at the start of the semester, ACT Scores, and the percentage of readings accessed all have a strong association with CT scores overall and for each dimension (Table 3; Table 4; Table 5). The sign for the dummy variable indicating female gender is also negatively associated with CT scores overall and the Objections score. On the other hand, class attendance, the number of pages assigned between assignments, hours of college credit, and majoring in political science are never statistically significant.

That being a political science major and hours of prior college credit are not significantly correlated with CT scores in political theory probably boosts the case for the existence of a domain-specific moral reasoning skill distinct from other forms of reasoning practiced in political science classes or elsewhere in the college curriculum.

---

<sup>6</sup> These figures are an average for all topics. Some individual topics may account for more, others less – hence the need for more fine grained analysis than I can give the topic here.

## Discussion

### A Theory of Rational Irrationality

One of the most salient results of the study is that students appear to adjust their level of CT to the stakes of each assignment. That students reduce their level of effort as the stakes in play decline is not a result that will surprise most instructors. What is perhaps not as obvious is that CT scores also drop as the stakes are lowered. But of course thinking clearly requires effort and discipline that is often in scarce supply.

The idea that individuals vary the quality of their reasoning based on the costs and benefits of thinking well is an idea pioneered by economist Bryan Caplan. Caplan calls this phenomenon “rational irrationality,” treating “cognitive inadequacy as a choice” (2007, 123) at least within some bounded region. Caplan draws a relationship between two familiar categories from epistemology: “epistemic rationality” and “instrumental rationality.”

Huemer (2016) defines epistemic rationality as “forming beliefs in truth-conducive ways – accepting beliefs that are well-supported by evidence, avoiding logical fallacies, avoiding contradictions, revising one’s beliefs in light of new evidence against them and so on,” noting in passing that “This is the kind of rationality that books on logic and critical thinking aim to instill.” (461) Instrumental rationality, on the other hand, is defined as “adopting suitable means to one’s ends” (Stanford Encyclopedia of

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Philosophy) or “judging an action by its anticipated consequences” (Hoppe 2003: 120),  
Because exercising epistemic rationality is costly in terms of time and effort, Caplan  
thinks it will not always be instrumentally rational to do so.

This trade-off between instrumental and epistemic rationality may be in evidence when  
students approach CT assignments. When the price of not thinking clearly is low,  
students respond by indulging in less rigorous thinking - for example by applying little  
scrutiny to facts, launching into chains of reasoning that flatter their self-image, or failing  
to examine alternative perspectives to their own. When the price of intellectual  
indulgence is higher, they respond by tightening up their standards of evidence and  
crafting own arguments more carefully.

If this theory of “rational irrationality” is correct, then important implications for course  
design and program assessment follow.

### **Low-Stakes Preparation, Assignment Overload, and Course Design**

As discussed in the literature review, a popular educational strategy consists of assigning  
low-stakes critical thinking practice assignments to build up student CT capabilities. If  
students are instrumentally rational about deploying their CT skills, this strategy could  
backfire in at least three ways.

One possibility is that small stakes practice means low quality practice. As with other  
skill-based activity, poor quality practice can lead to poor quality performance. Since the



1  
2  
3  
4  
5  
6 total grade in the class can only ever add up to 100%, this would help make sense of the  
7  
8 decline in CT scores as the number of assignments goes up. To practice well, students  
9  
10 would need to be properly incentivized, and this would put a cap on how many  
11  
12 scaffolding exercises an instructor could assign without seeing the quality of CT begin to  
13  
14 drop (through declining stakes).  
15  
16

17  
18  
19 A second possibility is less dramatic, and makes better sense of the results pointing to  
20  
21 CTEs and CREs having some positive effect on CT scores. On this view, low stakes  
22  
23 practice could be helpful, but only if it leaves enough weight to the assignments where  
24  
25 CT is measured. In some course designs, the quantity of CT exercises could simply be  
26  
27 crowding out the essays, leading to lower measured CT scores. CT skills could be  
28  
29 improving at the same time as the disposition to showcase those skills is being depressed  
30  
31 by the lower stakes attached to the essays.  
32  
33

34  
35  
36 Finally, thinking well may also simply take time. Too many assignments and students  
37  
38 might not have enough time between assignments to devote to thinking about their  
39  
40 essays. In that case, the proper spacing of assignments and especially giving students  
41  
42 enough time to think about the essay topic may help raise CT scores.  
43  
44  
45

46  
47  
48 Clearly more research is needed to disentangle these different hypotheses, but whatever  
49  
50 turns out to be the case, more low-stakes practice is not always pedagogically optimal.  
51  
52

### 53 54 **Manipulating the Stakes, Essay Topics, and Program Assessment**

55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7 Variations in stakes and topics should also give us pause when thinking about program  
8 assessment. If Department A conducts an exit exam with a CT component using a  
9 standardized test worth some percentage of the final grade in a senior seminar and  
10  
11 Department B requires the very same exit exam to graduate, but fails to attach a  
12  
13 consequence to low performance, statistics comparing the programs will fail to accurately  
14  
15 measure the differences in their graduates.  
16  
17  
18  
19

20  
21 The same will be true if students are given a CT exam upon entering college with one set  
22  
23 of stakes attached, but then given a CT exit exam in their senior year with different stakes  
24  
25 attached. This second scenario is further complicated by the fact that, controlling for  
26  
27 individual and course characteristics, different topics produce very different average  
28  
29 scores, as well different degrees of subscore variation. One can readily imagine a student  
30  
31 having made great strides in their ability to articulate alternative perspectives, yet have  
32  
33 that improvement go unmeasured if the post-test is much more difficult on that particular  
34  
35 dimension of critical thinking.  
36  
37  
38  
39

40  
41 Does it matter what the stakes are so long as they are consistent? According to the theory  
42  
43 of rational irrationality, the answer is yes. If the stakes are low, then only students who  
44  
45 have strong internal motivation will exercise the intellectual discipline necessary to  
46  
47 achieve high scores. If the stakes are high, both the internally motivated *and* those  
48  
49 capable of being externally motivated towards higher levels of CT will achieve higher  
50  
51 scores. If this second population is non-trivial in size, instructors and departments may be  
52  
53 underestimating their students' capabilities (and even their own performance) by some  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 margin when they attach no real stakes to CT assessments. The coefficients obtained in  
7  
8 this study suggest that the effect of changing the stakes is non-trivial. Doubling the stakes  
9  
10 of an essay from 5 to 10% of the final grade in the class could yield up to 20 extra  
11  
12 percentage points in measured CT.  
13

14  
15  
16  
17 Several researchers have suggested a need to focus higher education on developing  
18  
19 students who have the disposition (as well as the ability) to engage in critical thought  
20  
21 (Giancarlo and Facione 2001; Pascarella and Terenzini 2005). Looking at the findings  
22  
23 from this study, there could indeed be large gains in measured CT from improving the  
24  
25 disposition to exercise CT at low stakes.  
26  
27

## 28 29 30 **Directions for Further Research** 31 32

33  
34  
35  
36 Aside from experimental confirmation of the impact of stakes on measured CT scores,  
37  
38 several other areas for further research stand out.  
39

40  
41  
42 Differences between high- and low- achieving students are of special interest. This study  
43  
44 measures the average effect of different factors across all students, but there may be  
45  
46 subgroups of students whose demonstrated CT is particularly helped or harmed by  
47  
48 different contexts and instructional strategies. Other studies have found this sort of effect  
49  
50 among psychology students. (Williams et. al 2003, 2004)  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7 The negative coefficients for female students also merits further research. In other  
8 contexts and disciplines, gender has been found to have both positive, negative, and  
9 insignificant associations with CT. (Facione 1990, Walsh and Hardy 1999, Rudd Baker  
10 and Hoover 2000, Friedel et. al 2008) It would be helpful to understand the mechanism  
11 through which gender influences CT scores in political theory.  
12  
13  
14  
15  
16  
17  
18

19 One hypothesis (supported by informal class surveys of interest and the fact that average  
20 scores by topic vary sharply according to gender) is that the curriculum in the class is  
21 simply (on average) less interesting to female students. This would support Nussbaum's  
22 concern that political theory has done a poor job of catering to the interests of women  
23 (1999). Another explanation might be that the types of moral reasoning that are normal  
24 and accepted in political theory are excessively "masculine." (Held 1987) Either of these  
25 criticisms alone, if substantiated, would be motive enough for rethinking the traditional  
26 curriculum in this and other theory courses.  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38

39 Finally, Fitzgerald and Baird (2011) suggest that topics for CT assignments on normative  
40 issues should fall outside domains that are likely to be familiar and politically charged for  
41 them. By assigning topics that are highly ideologically charged (where students may be  
42 tempted to "consume" ideas or beliefs that they prefer to believe), one would expect CT  
43 scores to fall. Examining what makes some topic averages worse than others more  
44 generally seems like a promising avenue for future research and may provide some  
45 tangible payoffs to researchers and programs looking to develop topic-insensitive  
46 measures of CT for assessment purposes.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Works Cited

- Bailin, Sharon. 2002. "Critical Thinking and Science Education." *Science & Education* 11(4): 361–75.
- Barry, Carol L. et al. 2010. "Do Examinees Have Similar Test-Taking Effort? A High-Stakes Question for Low-Stakes Testing." *International Journal of Testing* 10(4): 342–63.
- Bataineh, Ruba Fahmi, and Lamma Hmoud Zghoul. 2006. "Jordanian TEFL Graduate Students' Use of Critical Thinking Skills (as Measured by the Cornell Critical Thinking Test, Level Z)." *International Journal of Bilingual Education and Bilingualism* 9(1): 33–50.
- Blings, Steffen, and Sarah Maxey. 2016. "Teaching Students to Engage with Evidence: An Evaluation of Structured Writing and Classroom Discussion Strategies." *Journal of Political Science Education* 13(1): 15–32.
- Bok, Derek. 2008. *Our Underachieving Colleges: A Candid Look at How Much Students Learn and Why They Should Be Learning More*. Princeton: Princeton University Press.
- Brown, J. D., T. Hilgers, and J. Marsella. 1991. "Essay Prompts and Topics: Minimizing the Effect of Mean Differences." *Written Communication* 8(4): 533–56.
- Caplan, Bryan Douglas. 2008. *The Myth of the Rational Voter: Why Democracies Choose Bad Policies*. Princeton, NJ: Princeton Univ. Press.

- 1  
2  
3  
4  
5  
6 Carbonaro, William. 2005. "Tracking, Students' Effort, and Academic Achievement."  
7  
8 *Sociology of Education* 78(1): 27–49.  
9
- 10 Cavdar, Gamze, and Sue Doe. 2012. "Learning through Writing: Teaching Critical  
11  
12 Thinking Skills in Writing Assignments." *PS: Political Science & Politics*  
13  
14 45(02): 298–306.  
15  
16
- 17 Condon, William, and Diane Kelly-Riley. 2004. "Assessing and Teaching What We  
18  
19 Value: The Relationship between College-Level Writing and Critical Thinking  
20  
21 Abilities." *Assessing Writing* 9(1): 56–75.  
22  
23
- 24 Elbow, Peter. 1997. "High Stakes and Low Stakes in Assigning and Responding to  
25  
26 Writing." *New Directions for Teaching and Learning* 1997(69): 5–13.  
27
- 28 Ennis, Robert H. 1989. "Critical Thinking and Subject Specificity: Clarification and  
29  
30 Needed Research." *Educational Researcher* 18(3): 4.  
31
- 32 Ewell, William Henry, and Robert R. Rodgers. 2014. "Enhancing Student Preparedness  
33  
34 for Class through Course Preparation Assignments: Preliminary Evidence from  
35  
36 the Classroom." *Journal of Political Science Education* 10(2): 204–21.  
37  
38
- 39 Facione, N C, P A Facione, and C A Sanchez. 1994. "Critical Thinking Dispositions as a  
40  
41 Measure of Competent Clinical Judgment: the Development of the Critical  
42  
43 Thinking Disposition Inventory. ." *Journal of Nursing Education* 33(8): 345–50.  
44
- 45 Facione, Peter A. 1990. *The California Critical Thinking Skills Test--College Level*.  
46  
47 *Technical Report No. 3. Gender, Ethnicity, Major, Ct Self-Esteem, and the Cctst*.  
48  
49 Education Resources Information Center.  
50
- 51 Facione, Peter A., and Noreen C. Facione. 1994. "Critical Thinking Ability: A  
52  
53 Measurement Tool." *Assessment Update* 6(6): 12–13.  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7 Fitzgerald, Jennifer, and Vanessa A. Baird. 2011. "Taking a Step Back: Teaching Critical  
8 Thinking by Distinguishing Appropriate Types of Evidence." *PS: Political*  
9 *Science & Politics* 44(03): 619–24.
- 10  
11  
12  
13 Friedel, Curt et al. 2008. "Overtly Teaching Critical Thinking and Inquiry-Based  
14 Learning: A Comparison of Two Undergraduate Biotechnology Classes."  
15 *Journal of Agricultural Education* 49(1): 72–84.
- 16  
17  
18  
19 Gaff, Jeffrey, and Jack Meacham. 2014. Learning Goals in Mission Statements:  
20 Implications for Educational Leadership *Learning Goals in Mission Statements:*  
21 *Implications for Educational Leadership*. Association of American Colleges &  
22 Universities. rep. [https://www.aacu.org/publications-](https://www.aacu.org/publications-research/periodicals/learning-goals-mission-statements-implications-educational)  
23 [research/periodicals/learning-goals-mission-statements-implications-educational](https://www.aacu.org/publications-research/periodicals/learning-goals-mission-statements-implications-educational)  
24 (July 21, 2017).
- 25  
26  
27  
28  
29  
30  
31  
32 Gelder, Tim Van. 2005. "Teaching Critical Thinking: Some Lessons From Cognitive  
33 Science." *College Teaching* 53(1): 41–48.
- 34  
35  
36  
37 Giancarlo, Carol Ann, and Peter A. Facione. 2001. "A Look Across Four Years at the  
38 Disposition Toward Critical Thinking Among Undergraduate Students." *The*  
39 *Journal of General Education* 50(1): 29–55.
- 40  
41  
42  
43 Halpern, Diane F. 2001. "Assessing the Effectiveness of Critical Thinking Instruction."  
44 *The Journal of General Education* 50(4): 270–86.
- 45  
46  
47  
48 Harvey, A. C. 1976. "Estimating Regression Models with Multiplicative  
49 Heteroscedasticity." *Econometrica* 44(3): 461–66.
- 50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6 Held, Virginia. 1987. "Non-Contractual Society: a Feminist View." In *Science, Morality,*  
7  
8 *and Feminist Theory*, Canadian Journal of Philosophy, Supplementary Volume,  
9  
10 eds. Marsha P. Hanen and Kai Nielsen. Calgary: University of Calgary Press.  
11  
12  
13 Hoppe, Robert. 2003. "Dealing with Multiple Rationalities: A Rejoinder to Ventris,  
14  
15 Tenbenschel and Snellen." *Administrative Theory & Praxis* 25(1): 119–27.  
16  
17 [http://www.jstor.org/stable/10.2307/25610593?ref=search-](http://www.jstor.org/stable/10.2307/25610593?ref=search-gateway:e554e4e2a58720e747c1a2eeb18b2d07)  
18  
19 [gateway:e554e4e2a58720e747c1a2eeb18b2d07](http://www.jstor.org/stable/10.2307/25610593?ref=search-gateway:e554e4e2a58720e747c1a2eeb18b2d07) (March 21, 2017).  
20  
21  
22 Howard, Larry W., Thomas Li-Ping Tang, and M. Jill Austin. 2014. "Teaching Critical  
23  
24 Thinking Skills: Ability, Motivation, Intervention, and the Pygmalion Effect."  
25  
26 *Journal of Business Ethics* 128(1): 133–47.  
27  
28  
29 Howard-Snyder, Frances, Daniel Howard-Snyder, and Ryan Wasserman. 2012. *The*  
30  
31 *Power of Logic*. New York: McGraw-Hill.  
32  
33 Huemer, Michael. 2016. "Why People Are Irrational About Politics." In *Philosophy,*  
34  
35 *Politics, and Economics: an Anthology*, eds. Jonathan Anomaly, Geoffrey  
36  
37 Brennan, Michael C. Munger, and Geoffrey Sayre-McCord. New York: Oxford  
38  
39 University Press. essay, 456–67.  
40  
41 "Job Outlook 2016: The Attributes Employers Want to See on New College Graduates'  
42  
43 Resumes." [http://www.naceweb.org/s11182015/employers-look-for-in-new-](http://www.naceweb.org/s11182015/employers-look-for-in-new-hires.aspx)  
44  
45 [hires.aspx](http://www.naceweb.org/s11182015/employers-look-for-in-new-hires.aspx) (March 19, 2017).  
46  
47  
48 Kelly-Riley, Diane, Gary Brown, Bill Condon, and Richard Law. 2008. Washington  
49  
50 Center University Critical Thinking Project: Resource Guide *Washington Center*  
51  
52 *University Critical Thinking Project: Resource Guide*. Washington State  
53  
54 University Critical Thinking Project. rep.  
55  
56  
57  
58  
59  
60



- 1  
2  
3  
4  
5  
6 Klein, S., R. Benjamin, R. Shavelson, and R. Bolus. 2007. "The Collegiate Learning  
7 Assessment: Facts and Fantasies." *Evaluation Review* 31(5): 415–39.  
8  
9  
10 Marks, Michael P. 2008. "Fostering Scholarly Discussion and Critical Thinking in the  
11 Political Science Classroom." *Journal of Political Science Education* 4(2): 205–  
12 24.  
13  
14  
15  
16  
17 Mason, Mark. 2007. "Critical Thinking and Learning." *Educational Philosophy and*  
18 *Theory* 39(4): 339–49.  
19  
20  
21 McPeck, John E. 1981. *Critical Thinking and Education*. New York: St. Martin's Press.  
22  
23  
24 Moore, Matthew J. "How (and What) Political Theorists Teach: Results of a National  
25 Survey." *Journal of Political Science Education* 7, no. 1 (2011): 95-128.  
26  
27 doi:10.1080/15512169.2011.539921.  
28  
29  
30 Mulcare, Daniel M., and Allan Schwedel. 2016. "Transforming Bloom's Taxonomy into  
31 Classroom Practice: A Practical Yet Comprehensive Approach to Promote  
32 Critical Reading and Student Participation." *Journal of Political Science*  
33 *Education* 13(2): 121–37.  
34  
35  
36  
37  
38  
39 Musekamp, Frank, and Jacob Pearce. 2016. "Student Motivation in Low-Stakes  
40 Assessment Contexts: an Exploratory Analysis in Engineering Mechanics." *Assessment & Evaluation in Higher Education* 41(5): 750–69.  
41  
42  
43  
44  
45  
46 Nussbaum, Martha Craven. 1999. *Sex & Social Justice*. Oxford: Oxford University Press.  
47  
48  
49 Olsen, Jonathan, and Anne Statham. 2005. "Critical Thinking in Political Science:  
50 Evidence from the Introductory Comparative Politics Course." *Journal of*  
51 *Political Science Education* 1(3): 323–44.  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7 Omelicheva, Mariya Y. 2007. "Resolved: Academic Debate Should Be a Part of Political  
8  
9 Science Curricula." *Journal of Political Science Education* 3(2): 161–75.
- 10  
11 Oros, Andrew L. 2007. "Let's Debate: Active Learning Encourages Student Participation  
12  
13 and Critical Thinking." *Journal of Political Science Education* 3(3): 293–311.
- 14  
15 Pascarella, Ernest T., and Patrick T. Terenzini. 2005. *How College Affects Students*. San  
16  
17 Francisco, CA: Jossey-Bass.
- 18  
19 Paul, R. 1991. "Critical Thinking: What Every Person Needs To Survive in a Changing  
20  
21 World." *NASSP Bulletin* 75(533): 120–22.
- 22  
23  
24 Phillips, John L., and Laura Mcmillian. 2010. "Liberal Civic Education, Religious  
25  
26 Commitment, and the Spillover Thesis: What Psychology Can Teach Us."  
27  
28 *Politics and Religion* 4(01): 27–48.
- 29  
30 Pithers, R.t., and Rebecca Soden. 2000. "Critical Thinking in Education: a Review."  
31  
32 *Educational Research* 42(3): 237–49.
- 33  
34  
35 Resnick, Daniel, and Natalie Peterson. 1991. "Evaluating Progress toward Goal Five: A  
36  
37 Report to the National Center for Education Statistics."  
38  
39 <https://eric.ed.gov/?id=ED340764>.
- 40  
41 Rudd, Rick, Matt Baker, and Tracy Hoover. 2008. "Undergraduate Agriculture Student  
42  
43 Learning Styles and Critical Thinking Abilities: Is There a Relationship?"  
44  
45 *Journal of Agricultural Education* 41(3): 2–12.
- 46  
47  
48 Sewell, Ellen. 2004. "Grade Dropping: An Empirical Analysis." *The Journal of*  
49  
50 *Economic Education* 35(1): 24–34.
- 51  
52  
53 Smith, Gerald. 2002. "Are There Domain-Specific Thinking Skills?" *Journal of*  
54  
55 *Philosophy of Education* 36(2): 207–27.
- 56  
57  
58  
59  
60

- 1  
2  
3  
4  
5  
6  
7 Stein, Barry et al. 2010. "Faculty Driven Assessment of Critical Thinking: National  
8  
9 Dissemination of the CAT Instrument." *Technological Developments in*  
10  
11 *Networking, Education and Automation*: 55–58.
- 12  
13 Terenzini, Patrick T., Leonard Springer, Ernest T. Pascarella, and Amaury Nora. 1995.  
14  
15 "Influences Affecting the Development of Students' Critical Thinking Skills."  
16  
17 *Research in Higher Education* 36(1): 23–39.
- 18  
19 Tiruneh, Dawit T., An Verburgh, and Jan Elen. 2014. "Effectiveness of Critical Thinking  
20  
21 Instruction in Higher Education: A Systematic Review of Intervention Studies."  
22  
23 *Higher Education Studies* 4(1).
- 24  
25 Walsh, Catherine M, and Robert C Hardy . 1999. "Dispositional Differences in Critical  
26  
27 Thinking Related to Gender and Academic Major." *Journal of Nursing*  
28  
29 *Education* 38(4): 49–155. <http://eric.ed.gov/?id=EJ583065> (April 13, 2017).
- 30  
31 Williams, Leonard, and Mary Lahman. 2011. "Online Discussion, Student Engagement,  
32  
33 and Critical Thinking." *Journal of Political Science Education* 7(2): 143–62.
- 34  
35 Williams, Robert L., and Stephen L. Worth. 2001. "The Relationship of Critical Thinking  
36  
37 to Success in College." *Inquiry: Critical Thinking Across the Disciplines* 21(1):  
38  
39 5–16.
- 40  
41  
42  
43 Williams, Robert L., Renee Oliver, and Susan L. Stockdale. 2004. "Psychological Versus  
44  
45 Generic Critical Thinking as Predictors and Outcome Measures in a Large  
46  
47 Undergraduate Human Development Course." *The Journal of General Education*  
48  
49 53(1): 37–58.  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6 Wise, Steven L., and Christine E. Demars. 2005. "Low Examinee Effort in Low-Stakes  
7  
8 Assessment: Problems and Potential Solutions." *Educational Assessment* 10(1):  
9  
10 1-17.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Peer Review Only

## Appendix A: The Grading Rubric

The papers in the class were graded along 6 dimensions of critical thinking, inspired by the grading rubric developed by the Washington State University Critical Thinking Project:

[Insert Table A1]

Each of these 6 categories is scored on a spectrum from 0 to 6 where 0 means absent, 1 means minimal, 2 means “emerging,” 3 means “developing,” 4 means “competent,” 5 means “effective,” and 6 means “mastering.”

Although the grading rubric used in this class departs from the original WSU categories, the WSU critical thinking project allows and even encourages variation and adaptation of the assessment instrument to different subjects. (Kelly-Riley. et. al 2008)

The WSU rubric separates Issue Identification from Identifying Contexts and Assumptions, whereas I combine them. The WSU Rubric combines Interpretation and Evidence, whereas I separate them. The justification for doing this initially was to make room for textual interpretation – which is a domain much cherished by political theorists. Combining Issue Identification and Identifying Context and Assumption was done because in normative contexts it is often hard to separate “determining what the issue is” from “determining what the different positions on the issue are.”

1  
2  
3  
4  
5  
6 This rubric also contains no Communication category. This category was scored but is not  
7 reported here. The reasons for this are twofold. On theoretical grounds it is not clear that  
8 command of the English language is conceptually related to CT. Moreover, a factor analysis  
9 conducted in earlier versions of the project determined that while the sources of variation in  
10 the six CT categories overlapped quite a bit – a single distinct source of variation could  
11 account for 94% of the variation in Communication scores. The distinction between CT and  
12 writing skill is something the authors of the WSU rubric also seem to concede (Condon and  
13 Kelly-Riley, 2004)  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24

25 For ease of interpretation, scores on the 7-point scale (0-6) are converted to a 100 point scale  
26 in the study.  
27  
28  
29  
30  
31

## 32 **Appendix B: Test of Intercoder Reliability**

33  
34  
35  
36  
37

### 38 **Methodology**

39 A post-test study of intercoder reliability was conducted between the original grader and two  
40 senior undergraduates applying to graduate school in political science. Although ideally a  
41 fully random and sample of papers would be used, time and budget constraints led to  
42 choosing a random sample of 4 topics – each of which had been used over multiple semesters.  
43 The topics used were Topics 103, 107, 110, and 112. It was felt that using the full range of  
44 topics would make the burden of grading unduly onerous on the undergraduates. Instead,  
45 three random semesters were selected and two anonymous random papers from each semester  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

were selected using an online random number generator. The test therefore comprised a total of 24 papers per grader (6 per topic).

In this instance, the principal worry is about overall agreement between different coders (especially that high rating and low ratings were consistent). Two measures are created. Mean Absolute Differences (MAD) and Pearson's  $r$ . The first establishes how far away, on average, each rating is from another. The second is a measure of correlation for each of the dimensions of critical thinking (as well as overall). Here obviously the higher the correlation, the better.

[Table B1]

### **Discussion**

The overall critical thinking score is satisfactory (considering that two of the graders were not subject matter experts) with a MAD score of 4.38 out of a possible 43 points (10.42% deviation) and an overall Pearson's  $r$  of .61. The thesis category appears to be the most problematic. For that category, Pearson's  $r$  correlations are weak and Mean MAD for the graders are over 20%. The evidence category has weak correlations but much better MAD scores, indicating close scores but uneven co-variation between the graders – which is somewhat puzzling. All other categories exhibit statistically significant covariation and acceptable MAD scores.

## **Appendix C – The Independent Variables**

1  
2  
3  
4  
5  
6  
7 [Table C1]

8  
9 **GPA:** Studies have found that students with higher GPAs do better on critical thinking  
10 measures. (Bataineh and Zghoul 2006; Howard, Li-Ping Tang, and Austin 2014) The  
11 inclusion of GPA is also important for another reason: weaker students drop out of the  
12 class as the semester progresses, leading to fewer observations for those students than for  
13 stronger students. Controlling for GPA will limit the effect of any bias arising out of  
14 missing data.  
15  
16  
17  
18  
19  
20  
21  
22

23  
24 **ACTs:** Prior studies have found a relationship between ACT scores and critical thinking  
25 over and above the impact of GPA. (Howard, Li-Ping Tang, and Austin 2014) A  
26 downside of using ACT scores as a control variable in this particular study is that many  
27 students at the university are not required to take any kind of standardized test for  
28 admissions. Using ACT scores causes 146 observations to drop from the study (35% of  
29 the total). The populations are similar in most respects. However, two statistically  
30 significant differences are worth noting: the mean age of excluded population is much  
31 higher (20.83 vs. 27.75) and the gender distribution skews more female (60.59% male vs.  
32 51.41% male).  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

45  
46 In the end, the decision was taken to include ACT scores in the study on the grounds that  
47 even if the results skewed the population of students towards more traditional students,  
48 this might make the study more comparable to other studies of critical thinking at the  
49 college level. A limitation of the study is that the critical thinking habits of older non-  
50 traditional students must be set aside for further study and that generalization of any  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3  
4  
5  
6 insights to this population would be hazardous at best. Three ACT scores were collected  
7  
8 (English, Math, and Reading). These were averaged into a composite score.  
9

10  
11  
12 **Reading Percentage/Slide Percentage:** The reading percentage and slide percentage  
13 variables measure the percentage of online readings and slides accessed from the  
14 beginning of class up until a given scored essay<sup>1</sup>. Although we cannot be sure that a  
15 student who accessed the readings/slides actually read them, we can be fairly confident  
16 that a file that was never visited was never read.  
17  
18  
19  
20  
21  
22

23  
24  
25  
26 Only a portion of the course readings are online. The reading measure does not capture  
27 engagement with the textbooks and only very few classic texts. Some students may read  
28 regularly from the textbooks but access online content only seldom. Other students may  
29 never buy the textbooks, but reliably access any content available for free online. All  
30 students who did not have access to the slides during the semester in which they took the  
31 course were coded as zero.  
32  
33  
34  
35  
36  
37  
38

39  
40 **Age:** Age proxies for maturity and experience. Age in the study is measured in years. If a  
41 student turned a different age during the semester, they were coded as having the age they  
42 began the course with.  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52

53  
54 

---

<sup>1</sup> A related measure, the percentage of readings accessed since the last assignment was also  
55 computed, but dropped from the study since it produced very similar overall results  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7 **Attendance:** Attendance grades are a combination of class attendance and attendance at  
8 extra-curricular events. Every class missed equals 5 points less on the attendance grade.  
9

10 The failure of class attendance to strongly correlate with critical thinking measures in the  
11 study could be due to its imprecise measurement.  
12  
13

14  
15  
16  
17 **Gender:** Students who requested to be treated as a belonging to a different gender than the  
18 one suggested by their first name were assigned the gender of their choice in the study (3 of  
19 151), but this had no impact on results.  
20  
21  
22

23  
24  
25 **Hours of College Credit:** An alternate measure of college experience was also constructed in  
26 the form of an ordinal variable separating students into First Year Students, Sophomores,  
27 Juniors, Seniors, and Super-seniors (over 120 college credits) but did not prove any more  
28 helpful in predicting CT.  
29  
30  
31  
32

33  
34  
35 **Exams:** Exams were short answer exams, covering terms and concepts from the readings  
36 and lectures, as well as some questions asking students to apply principles to novel  
37 situations.  
38  
39  
40  
41  
42

43  
44 **Previous Essays/Critiques/CTEs:** These variables track the student, not the number of  
45 assignments to date in that semester. For example, if a student fails to hand in essay #1,  
46 the variable is coded 0 for that student even though all other students might be coded 1  
47 for that essay. CTEs were worth between 2.5 and 5% of the final grade. Critiques were  
48 worth 5-10% of the final grade.  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Appendix D: Discussion of Limitations Due to Cross-Sectional Design

One danger presented by the multiple regression methodology for data spanning across semesters in this way is that individual measures of critical thinking are not independent of one another. Although some effort was made to control for class and assignment level variables, it is possible that each semester exerted an influence not captured by the control variables. If the main sources of semester variation are captured this should not be a problem – but omitted variable bias is always a threat<sup>2</sup>.

A Durbin Watson test was effected for each model to check for serial autocorrelation, with results between 2.06 and 2.28 showing no statistically significant autocorrelation (Table D1). This suggests that the problems associated with non-independence of observation are minimal.

Visually inspecting the pattern of residuals through a Q-Q plot reveals well-behaved residuals for the Focus, Interpretation, Objections, Evidence, and Overall Critical Thinking. Those for the Thesis and Conclusion models present potential problems for the assumption that errors are normally distributed. A Shapiro-Wilk test confirms this. (Table D1) This could suggest bias in the coefficients for those models due to one or more omitted variables. It could also signify non-random sources of bias in grading.

---

<sup>2</sup> A regression using a semester dummy variable and a standard set of individual and assignment level controls finds that semester level effects are significant in the aggregate and account for around 3.5% of the variation in aggregate critical thinking scores.

1  
2  
3  
4  
5  
6 The problematic models were re-run using log-linear variance regression (Harvey 1976) to  
7 simultaneously fit the expected value of the dependent variables and studentize the residuals  
8 by fitting the variance of the responses (thus producing a model with heteroscedasticity-  
9 consistent standard errors). The percentage of slides variable drops from the Conclusion  
10 model. The essay value variable drops from significance in the Thesis model. Otherwise, the  
11 results are consistent.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 1**  
**Correlation Matrix – Part 1 (Pearson's r)**

	GPA	ACT Avg	Age	Gend.	Reading%	Slide %	Attend	Hours	Maj.
GPA	1.00								
ACT Avg	0.39***	1.00							
Age	-0.24***	-0.12**	1.00						
Gender	0.08	0.03	-0.13**	1.00					
Reading %	0.28***	-0.03	0.06	0.18***	1.00				
Slide %	-0.05	-0.17***	0.19***	0.00	0.25***	1.00			
Attend.	0.31***	0.00	0.05	0.02	0.28***	0.15***	1.00		
Hours	0.15**	0.07	0.44	0.00	0.09	0.10**	0.08	1.00	
Major	0.11*	-0.03	0.09	-0.02	0.11*	0.06	0.09	0.20**	1.00
Essay Drop	-0.10*	-0.16***	-0.06	0.11*	0.20***	0.49***	0.04	0.03	0.09
Essay value	-0.03	0.11*	0.15**	-0.09	-0.12**	-0.08	0.17***	-0.10*	0.13**
# of Assignments	-0.10*	-0.01	-0.07	0.06	0.13**	0.48***	-0.06	0.08	0.15***
LGs	0.03	-0.03	0.13**	-0.14**	-0.16***	-0.14**	-0.06	-0.01	-0.16***
Slide dummy	-0.10*	-0.06	0.03	-0.01	0.23***	0.61***	0.17***	-0.06	0.04
# of Pages	0.04	0.02	-0.03	0.05	0.09	0.02	0.07	-0.03	0.04
Exams	-0.05	-0.09	0.09	0.02	0.08	0.40***	0.00	0.06	-0.05
CTEs	-0.02	-0.10	0.08	0.04	0.17***	0.54***	0.00	0.13**	-0.02
Prev. Essays	0.00	0.00	-0.01	-0.02	-0.08	0.08	-0.01	0.07	0.00
CREs.	0.08	0.11	-0.12**	-0.12**	-0.17***	-0.37	0.01	-0.03	0.08

p<.1 \*\*: p<.05 \*\*\*p<.01

**Table 2**  
**Correlation Matrix – Part 2 (Pearson’s r)**

	Essay value	# of Assignments	LGs	Slide dummy	# of Pages	Exam	CTEs	Prev. Essay	CREs
<b>Essay Drop</b>	-0.33***	0.50***	-0.48***	0.55***	0.17***	0.31***	0.38***	0.09	-0.28***
<b>Essay value</b>	1.00	0.07	-0.21***	0.18***	0.23***	-0.20***	-0.20***	-0.12**	0.11*
<b># of Assignments</b>		1.00	-0.59***	0.56***	0.13**	0.22***	0.66***	0.20***	-0.25***
<b>LGs</b>			1.00	-0.38***	-0.29***	0.25***	0.00	-0.12**	-0.04
<b>Slide dummy</b>				1.00	0.30***	0.34***	0.40***	-0.01	-0.23***
<b># of Pages</b>					1.00	0.08	-0.06	-0.13**	0.12**
<b>Exams</b>						1.00	0.69***	0.29***	-0.33***
<b>CTEs</b>							1.00	0.34***	-0.39***
<b>Prev. Essays</b>								1.00	0.47***
<b>CREs.</b>									1.00

p<.1 \*\*: p<.05 \*\*\*p<.01

**Table 3**  
**Stepwise Backward Elimination Regression Results: Evidence, Conclusion, Objections**

	Evidence			Conclusion			Objections		
	<i>B</i>	<i>SE</i>	$(\eta^2/r^2)$	<i>B</i>	<i>SE</i>	$(\eta^2/r^2)$	<i>B</i>	<i>SE</i>	$(\eta^2/r^2)$
<b>Intercept</b>	-1.72	16.87		10.20	18.78		-69.04***	17.55	
<b>TopicID (Prob&gt;F)</b>	.100*		0.085	.137		0.083	<.001***		0.157
<b>ACTs</b>	1.20***	0.29	0.044	1.05***	0.33	0.027	0.77**	0.36	0.011
<b>GPA</b>	6.27***	2.14	0.022	4.12*	2.40	0.008	15.98***	2.53	0.097
<b>Essay value</b>	4.90***	1.68	0.022	5.47***	1.77	0.026	3.59*	1.84	0.009
<b>Reading %</b>	0.13***	0.05	0.020	0.15***	0.05	0.022			
<b>CREs</b>	5.22**	2.48	0.012	7.02***	2.38	0.023	6.38**	2.61	0.014
<b>Slide dummy[1]</b>	7.97**	3.68	0.012						
<b>Exams</b>	18.42**	7.38	0.016						
<b># of Assignments</b>	-6.20***	1.43	0.049	-5.83***	1.70	0.032			
<b>LG Dummy[1]</b>	-7.21*	4.35	0.007	-10.81**	4.38	0.016			
<b>CTEs</b>				13.82***	3.73	0.037	6.64**	2.95	0.012
<b>Slide %</b>				12.65**	5.79	0.013			
<b>Gender [F]</b>							-4.11***	1.47	0.019
<b>Age</b>									
<b>Prev. Essays</b>									
<b>Essay Drop</b>									
<b># of Pages</b>									
<b>Attendance</b>									
<b>Hours</b>									
<b>Major</b>									
<b>Adj. R<sup>2</sup></b>		.256			.232			.305	
<b>R<sup>2</sup></b>		.340			.319			.375	
<b>N</b>		286			286			289	

\*: p<.1; \*\*p<.05; \*\*\*:p<.01

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 4**  
**Stepwise Backward Elimination Regression Results: Thesis, Focus, Interpretation**

	Thesis			Focus			Interpretation		
	<i>B</i>	<i>SE</i>	( $\eta^2/r^2$ )	<i>B</i>	<i>SE</i>	( $\eta^2/r^2$ )	<i>B</i>	<i>SE</i>	( $\eta^2/r^2$ )
<b>Intercept</b>	-10.07	16.75		-52.25***	19.42		-19.63	19.12	
<b>TopicID (Prob&gt;F)</b>	.012**		0.128	.010***		0.127	.100*		0.092
<b>ACTs</b>	0.97***	0.35	0.022	0.89**	0.35	0.019			
<b>GPA</b>	10.14***	2.50	0.048	7.87***	2.67	0.025	11.84***	2.33	0.069
<b>Essay value</b>	3.08*	1.76	0.009				4.19**	1.95	0.012
<b>Reading %</b>				0.14***	0.05	0.020	0.20***	0.05	0.040
<b>CREs</b>									
<b>Slide dummy[1]</b>				-7.07**	3.26	0.013			
<b>Exams</b>							10.72**	4.97	0.012
<b># of Assignments</b>									
<b>LG Dummy[1]</b>									
<b>CTEs</b>									
<b>Slide %</b>									
<b>Gender [F]</b>									
<b>Age</b>				2.30***	0.76	0.026			
<b>Prev.Essays</b>				6.81**	3.16	0.013			
<b>Essay Drop</b>							-5.58**	2.79	0.011
<b># of Pages</b>									
<b>Attendance</b>									
<b>Hours</b>									
<b>Major</b>									
<b>Adj. R<sup>2</sup></b>		.151			.183			.232	
<b>R<sup>2</sup></b>		.229			.268			.310	
<b>N</b>		292			289			289	

\*: p<.1; \*\*p<.05; \*\*\*:p<.01



1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**Table 5**  
**Stepwise Backward Elimination Regression Results: Aggregate Critical Thinking**

	All Regressors			Selected Model		
	B	SE	( $\eta^2/r^2$ )	B	SE	( $\eta^2/r^2$ )
<b>Intercept</b>	-35.56*	20.18		-19.71	12.48	
<b>TopicID (Prob&gt;F)</b>	.682			.453		0.056
<b>ACTs</b>	0.82***	0.24	0.028	0.84***	0.22	0.075
<b>GPA</b>	9.67***	1.97	0.056	9.07***	1.62	0.045
<b>Essay value</b>	4.91***	1.57	0.023	4.16***	1.10	0.035
<b>Reading %</b>	0.13***	0.04	0.029	0.15***	0.03	0.016
<b>Gender[F]</b>	-2.35**	0.92	0.015	-2.31**	0.89	0.011
<b># of Assignments</b>	-4.17***	1.46	0.019	-1.55**	0.73	0.035
<b>CREs</b>	3.49	2.23	0.006			
<b>Slide dummy[1]</b>	4.22	4.14	0.002			
<b>LG Dummy[1]</b>	-10.11*	5.29	0.009			
<b>Hours</b>	-0.02	0.04	0.000			
<b>Age</b>	1.15*	0.69	0.007			
<b>Prev. Essays</b>	0.45	3.32	0.000			
<b>Dropped Essay[1]</b>	-3.68	4.06	0.002			
<b>Exams</b>	10.19	6.65	0.005			
<b>Attendance</b>	0.02	0.04	0.001			
<b>CTEs</b>	5.08*	2.85	0.007			
<b>Total Slide %</b>	5.31	4.14	0.004			
<b>Pages Since Last Essay</b>	-0.03	0.04	0.001			
<b>Major[1]</b>	-1.56	1.21	0.004			
<b>Adj. R<sup>2</sup></b>		.337			.311	
<b>R<sup>2</sup></b>		.434			.381	
<b>N</b>		286			288	

\*: p<.1; \*\*p<.05; \*\*\*:p<.01

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49

**Table A.1**  
**Dimensions of Critical Thinking, Aggregate Critical Thinking and Descriptive Statistics**

<b>Dimension</b>	<b>Description</b>	<b>Mean</b>	<b>Std. Dev.</b>	<b>N</b>
<b>Issue Identification and Focus</b>	How well did the student show that they understood the question that was being asked? Did the student adequately examine all the concepts put into play by the prompt? Did the student stay on task and answer the question asked?	57.74	23.85	447
<b>Interpretation</b>	Did the student fairly represent different possible positions on the issue to be found in the literature? Did the student more generally draw plausible inferences from the sources they used?	63.60	24.57	446
<b>Evidence</b>	Did the student produce appropriate evidence for claims made in the essay? Did the argument build from plausible assumptions/broadly shared intuitions? Did the student show an awareness of the reliability of different sources for empirical evidence (if such were used)?	53.80	21.04	447
<b>Objections</b>	Did the student consider prominent objections to their thesis or to their assumptions? Were such objections plausibly dealt with?	35.78	26.37	447
<b>Thesis</b>	Does the paper have a clear thesis? Is the position the student's (or is it a simple restatement of a position found in the literature)? Is the argument developed internally consistent – given its premises and assumptions?	66.76	24.59	447
<b>Conclusion</b>	Is the conclusion consistent with the argument presented in the paper? Has the student identified implications for other areas of inquiry or public policy?	73.49	22.69	447
<b>Aggregate CT</b>	An additive measure combining 6 dimensions of critical thinking (each being weighted equally).	58.39	16.58	446

**Table B1****Measures of Intercoder Reliability: Mean Absolute Deviations (MAD) and Correlation Coefficients (Pearson's r)**

Focus	Undergraduate Grader #1		Undergraduate Grader #2		Faculty Regrade		Mean Overall	
	MAD	Pearson's r	MAD	Pearson's r	MAD	Pearson's r	MAD	Pearson's r
<b>Focus</b>	1.27 (18.22%)	0.55***	1.27 (18.10%)	0.35	1.09 (15.62%)	0.60***	1.21 (17.31%)	0.50***
<b>Interpretation</b>	0.97 (13.89%)	0.59***	1.27 (18.10%)	0.57***	1.30 (18.59%)	0.55***	1.06 (15.22%)	0.57***
<b>Evidence</b>	0.78 (11.22%)	0.20	0.92 (13.17%)	0.25	0.89 (12.78%)	0.30	0.87 (12.41%)	0.25
<b>Objections</b>	1.47 (21.07%)	0.42**	1.20 (17.15%)	0.65***	1.07 (15.31%)	0.70***	1.25 (17.84%)	0.59***
<b>Thesis</b>	1.33 (19.02%)	0.29	1.63 (23.32%)	0.15	1.41 (20.11%)	0.30	1.46 (20.82%)	0.25
<b>Conclusion</b>	1.22 (17.51%)	0.39*	0.96 (13.73%)	0.53***	0.50 (7.09%)	0.81***	0.89 (12.78%)	0.57***
<b>Overall</b>	5.38 (12.81%)	0.56***	4.57 (10.87%)	0.57***	3.18 (7.58%)	0.70***	4.38 (10.42%)	0.61***

\*: p&lt;.1; \*\*p&lt;.05; \*\*\*:p&lt;.01

**Table C1****Descriptive Statistics: Continuous Variables**

	GPA	ACT Avg	Age	Slide %	Reading%	Attend	Hours	Essay value	# of Assignments	# of Pages	Exams	CTEs	CREs
<b>Mean</b>	3.04	22.97	23.26	0.25	74.18	83.99	79.76	8.20	8.24	163.91	0.19	0.51	0.76
<b>Std. Dev</b>	0.64	4.15	5.51	0.34	28.60	22.14	28.80	1.02	2.11	80.51	0.49	1.13	0.94
<b>Min</b>	1	14	18	0	0	0	16	7.14	5	72	0	0	0
<b>Max</b>	4	32	62	1	100	113	190	10	14	477	2	5	3
<b>Collection Method</b>	UR	UR	UR	D2L	D2L	D2L	UR	SYL	SYL & D2L	SYL	SYL	SYL/ D2L	SYL/ D2L
<b>N</b>	443	294	450	445	444	447	446	450	449	447	451	447	451

**Descriptive Statistics: Dummy Variables**

	Slide Dummy (1=Yes)	Logic Game Dummy (1=Yes)	Major (1=Yes)	Gender (1=Fem)	Dropped Essay (1=Yes)
<b>0</b>	198	239	67	268	175
<b>1</b>	253	208	381	182	275
<b>Collection Method</b>	D2L	SYL	UR	OBS	SYL
<b>N</b>	451	447	448	450	450

**UR=University Records; D2L=Online Course Management System; SYL=Syllabus; OBS=Observation**

**Table D1**  
**Tests of Model Assumptions: Autocorrelation (Durbin Watson), Normal**  
**Distribution of Residuals (Shapiro-Wilk), and Associated Probabilities of Accepting**  
**Model Assumptions**

	Durbin Watson	Prob<DW	Shapiro- Wilk	Prob<W	N
<b>Aggregate CT</b>	2.14	.44	.997	.95	288
<b>Evidence</b>	2.20	.55	.992	.10	286
<b>Conclusion</b>	2.21	.61	.995	.001***	286
<b>Objections</b>	2.06	.19	.995	.041	289
<b>Thesis</b>	2.28	.88	.982	.001***	292
<b>Focus</b>	2.31	.92	.991	.11	289
<b>Interpretation</b>	2.07	.22	.989	.017**	289