

# Spatial Reasoning and Statistical Graphics

Susan VanderPlas, Heike Hofmann

October 22, 2014

## 1 Introduction

Intro about statistical graphics, lineups, ...

### 1.1 Spatial Reasoning and Statistical Graphics

Statistical graphics provide quick summaries of data, models, and results, but these displays may not be equally useful to all viewers. Mathematical ability is important, even for the simplest graphs: [Shah and Carpenter \(1995\)](#) showed that mathematical ability, but not spatial ability, was associated with accuracy on a simple two-dimensional line graph. Spatial ability becomes more important, however, when more complicated graphical displays are used in comparison tasks: [Mayer and Sims \(1994\)](#) hypothesize that those with low spatial ability perform worse on tests utilizing diagrams and graphs because more cognitive resources are required to process the visual stimuli (leaving fewer resources to make connections and draw conclusions from those stimuli). Many theories of graphical learning center around the difference between visual and verbal processing: the dual-coding theory emphasizes the utility of complementary information in both domains, while the visual argument hypothesis emphasizes that graphics are more efficient tools for providing data with spatial, temporal, or other implicit ordering, because the spatial dimension can be represented graphically in a more natural manner ([Vekiri, 2002](#)). Both of these theories suggest spatial ability would impact a viewer's use of graphics, because spatial ability either influences cognitive resource allocation or affects the processing of spatial relationships between graphical elements. Previous investigations into graphical learning and spatial ability have found relationships between spatial ability and the ability to read information from graphs ([Lowrie and Diezmann, 2007](#)).

### 1.2 The Lineup Protocol

More lineup introduction here

Statistical lineups ([Hofmann et al., 2012](#); [Buja et al., 2009](#)) depend on the ability to search for a signal amid distractors (visual search) and the ability to infer patterns from stimuli (pattern recognition). Some lineups (polar coords) also depend on the ability to mentally rotate stimuli (spatial rotation) and mentally manipulate graphs (spatial rotation and manipulation). By breaking the lineup task down into component parts, we may be able to determine which visuospatial factors most strongly correlate with lineup performance, using carefully chosen cognitive tests to assess these aspects of visuospatial ability. These assessments will then be used to account for some of the variation in lineup performance differences. In addition, we will also examine the effect that previous experiences (science-based major, research experience, Auto-CAD skills) have on lineup performance and the implications for the use of lineups to assess the graphical display of statistical information.

## 2 Methods

Use the setup in the lineup protocol section to justify each of the tests chosen.

Besides pointing out what we asked participants to do, we also need some why, and also why in particular these tests were selected. part of the validity consideration can go in here - and the classification from the naval test of the tests we picked into different classes (maybe already with a hint that for our data these weren't actually orthogonal classes, because we found a really high correlation between card rotation and paper folding.)

*I don't know if the tests are supposed to be orthogonal. It would be hard to construct a set of tests that was orthogonal, because the abilities themselves aren't easily separable - math reasoning and spatial reasoning are related, for instance, and verbal abilities are correlated with creativity too.*

Participants will complete the following tasks (sample pictures included, full stimuli set will be added to the appendix once we are sure there is no need for follow-up experiments). Tasks are designed so that participants are under time pressure; they are not expected to complete all of the problems in each section. This provides more discrimination between high scorers and prevents score compression at the top of the range.

- Visual Search Task: designed to test participants' ability to find a target stimulus in a field of distractors. An example is shown in figure 1. The visual search task is similar in concept to lineups: it tests one's ability to find the target plot. Historically, it has been used as a measure of brain damage (Goldstein et al., 1973; DeMita et al., 1981; Moerland et al., 1986); however, similar tasks, have been used to measure cognitive performance in a variety of situations (under the influence of drugs, for example, in Anderson and Revelle (1983)). The similarity to lineup protocol as well as the simplicity of the test and its' lack of color justify the slight deviation from forms of visual search tasks typically used in normal populations.
- Paper Folding Task: tests participants' ability to visualize and mentally manipulate figures in three dimensions. Associated with the ability to extrapolate symmetry and reflection over multiple steps. An example is shown in figure 2. Lineups require similar manipulations in two-dimensional space, and also require the ability to perform complex spatial manipulations mentally (for instance, comparing the interquartile range of two boxplots as well as their relative alignment to a similar set of two boxplots in another panel).
- Card Rotation Task: tests participant's ability to rotate objects in two dimensions to distinguish between left-hand and right-hand versions of the same figure. Tests spatial reasoning ability and mental rotation skills. An example is shown in figure 3. As with the paper folding test, two-dimensional comparisons are an important component of lineup performance. In some lineup situations, these comparisons involve translation, in other lineups, rotation is required.
- Figure Classification Task: tests participant's ability to extrapolate rules from provided figures. This task is associated with visual reasoning capabilities and we expect that it should correlate with the ability to pick out a signal plot from a lineup. An example is shown in figure 4. The figure classification test requires the same types of reasoning as the lineups: participants must determine the rules from the provided classes, and extrapolate from those rules to classify new figures. In lineups, participants must determine the rules based on the panels appearing in the lineup; they must then identify the plot which does not conform. As such, the figure classification test has content validity in relation to lineup performance: it is measuring similar underlying criteria.

Between cognitive tasks, participants will also complete three blocks of 20 lineups each. These lineups have been previously tested (Hofmann et al., 2012). Participants have 5 minutes to complete each block of 20 lineups. Figure 5 shows a sample lineup of box plots.

In addition to these tests, participants will complete a questionnaire which includes questions about colorblindness, mathematical background, self-perceived verbal/mathematical/artistic skills, time spent playing

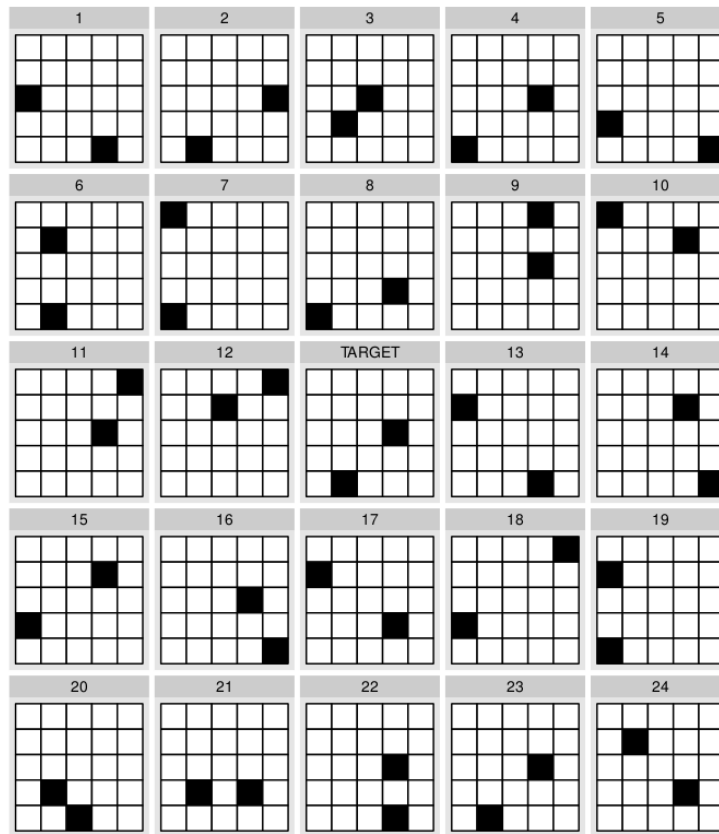


Figure 1: Visual Search Task. Participants are instructed to find the plot numbered 1-24 which matches the plot labeled "Target". Participants will complete up to 25 of these tasks in 5 minutes.

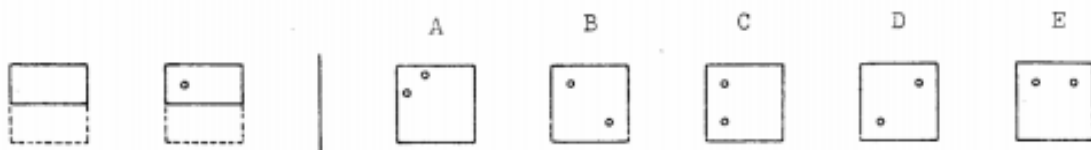


Figure 2: Paper Folding Task. Participants are instructed to pick the figure matching the sequence of steps shown in the left-hand figure. Participants will complete up to 20 of these tasks in 6 minutes.



Figure 3: Card Rotation Task. Participants mark each figure on the right hand side as either the same or different than the figure on the left hand side of the dividing line. Participants will complete up to 20 of these tasks (each consisting of 8 figures) in 6 minutes.

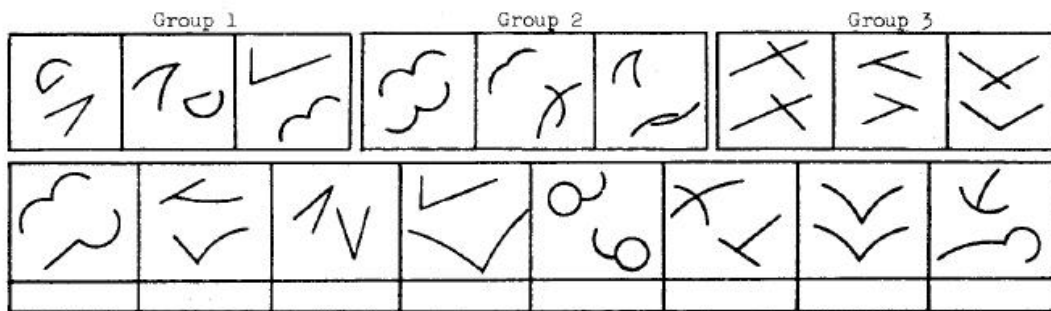


Figure 4: Figure Classification Task. Participants classify each figure in the second row as belonging to group 1, 2, or 3 (if applicable). Participants will complete up to 14 of these tasks (each consisting of 8 figures to classify) in 8 minutes.

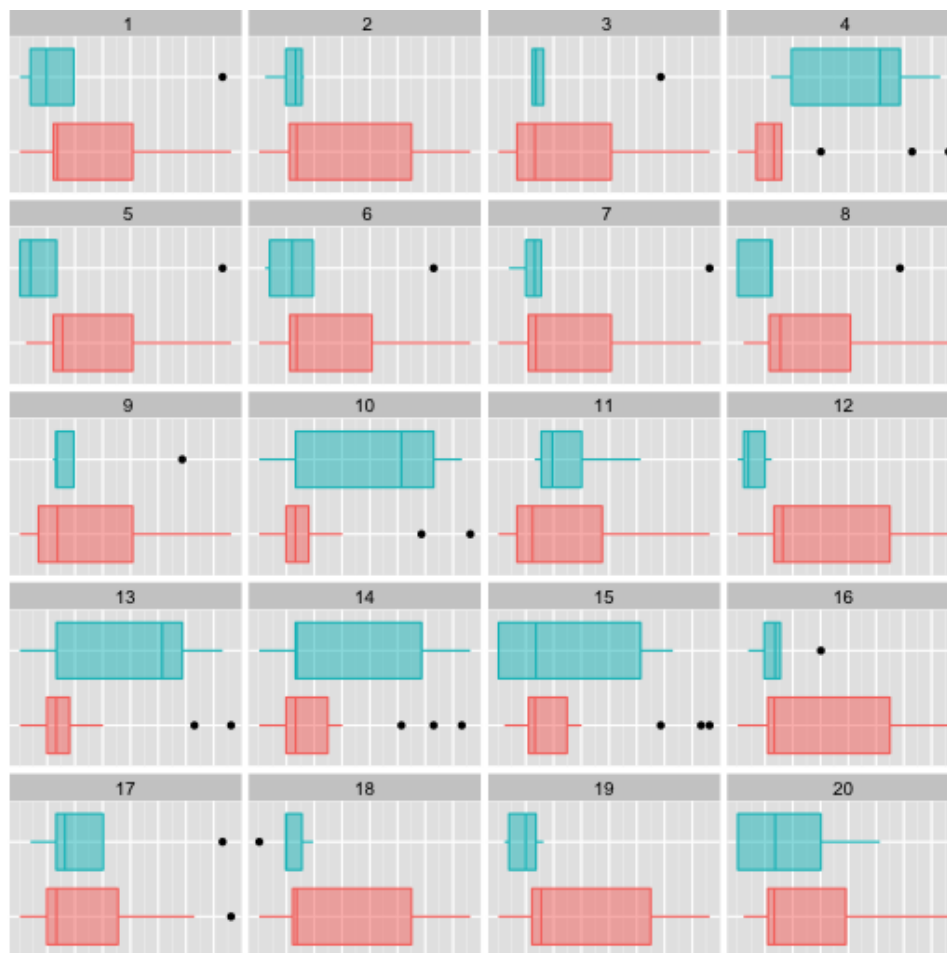


Figure 5: A sample lineup. Participants are instructed to choose the plot which appears most different from the others. In this lineup, plot 13 is the target.

video games, and undergraduate major. These questions are designed to assess different factors which may influence a participant’s skill at reading graphs and performing spatial tasks.

## 2.1 Test Scoring

Scoring of all test results was done such that random guessing leads to an expected value of 0; therefore each question answered correctly contributes to the score by 1, while a wrong answer is scored by  $-1/(k - 1)$ , where  $k$  is the total number of possible answers to the question. Thus, for a test consisting of multiple choice questions with  $k$  suggested answers with a single correct answer each, the score is calculated as

$$\# \text{correct answers} - 1/(k - 1) \cdot \# \text{wrong answers.} \quad (1)$$

This allows us to compare each participant’s score in light of how many problems were attempted as well as the number of correct responses. Combining accuracy and speed into a single number does not only make a comparison of test scores easier, this scoring mechanism is also used on many standardized tests, such as the SAT and the battery of psychological tests (Diamond and Evans, 1973; Ekstrom et al., 1976) from which parts of this test are drawn.

could we also refer to the fact that these tests have been long established (maybe with references to earlier tests or updates?) *Yes, definitely. I'll pull up some references for that too...*

The advantage of using tests from the Kit of Factor Referenced Cognitive tests (Ekstrom et al., 1976) is that the tests are extremely well studied (including an extensive meta-analysis by Voyer et al. (1995) of the spatial tests we are using in this study) and comparison data are available from the validation of these factors (Schaie et al., 1998; Hampson, 1990; Mayer and Sims, 1994) and previous versions of the kit (French et al., 1963).

the well studied reference list needs to be a bit more extensive ... do you know this book by Michel Hersen?  
*I was working on it yesterday along with trying to validate the general idea that reading graphics requires certain spatial abilities. It's definitely a work in progress. I haven't seen that book before, but I'm not at all surprised it exists.*

## 3 Results

Results are based on an evaluation of 38 undergraduate students at Iowa State University. About one-half of the participants were in STEM fields, the others were distributed relatively evenly between agriculture, business, education, and liberal arts. Students were evenly distributed by gender, and were between 18 and 24 years of age with only one exception. This is reasonably representative of the university as a whole; ISU graduates are about 20% engineering, 16% business, and 11% agriculture students, however, these numbers are for completed degrees, not current majors.

Heike, can you get better enrollment numbers for current students? My guess is there are more people enrolled in engineering than actually graduate :).

The testing was conducted by two research assistants; the first tested 18 participants, the second tested 20 participants. A comparison of results from each researcher is available in appendix A.

**Comparison of Spatial Tests with Previously Validated Results** The card rotation, paper folding, and figure classification tests have been validated using different populations, many of which are demographically similar to Iowa State students (naval recruits, college students, late high-school students, and 9th grade students). We compare Iowa State students’ unscaled scores in table 1, adjusting data from other populations to account for subpopulation structure and test length.

<sup>1</sup>ISU students took only Part I due to time constraints.

|                 | Card Rotation                      | Paper Folding                         | Figure Classification   | Visual Search |
|-----------------|------------------------------------|---------------------------------------|---|---------------|
| ISU Students    | 83.4 (24.1)                        | 12.4 (3.7)                            | 57 (23.8) <sup>1</sup>  | 21.9 (2.3)    |
| Scaled Scores   | 88.0 (34.8)                        | 13.8 (4.5)                            | 58.7 (14.4) <sup>2</sup>  | N/A           |
| Unscaled Scores | 44.0 (24.6) <sup>3</sup>           | 13.8 (4.5)                            | M: 120.0 (30.0), F: 114.9 (27.8)  | N/A           |
| Population      | approx. 550 male<br>naval recruits | 46 college students<br>(1963 version) | suburban 11th & 12th grade students<br>(288-300 males, 317-329 females) | N/A           |

Table 1: Comparison of scores from Iowa State students and scores reported in [Ekstrom et al. \(1976\)](#). Scaled scores are calculated based on information reported in the manual, scaled to account for differences in the number of questions answered during this experiment. Data shown are from the population most similar to ISU students, out of the data available. The Visual Search task ([Goldstein et al., 1973](#); [DeMita et al., 1981](#); [Moerland et al., 1986](#)) is not part of the Kit of Factor Referenced Cognitive Test data, and thus we do not have comparison data for the form used in this experiment.

Table 1 shows mean scores and standard deviation for ISU students and other populations. Values have been adjusted to accommodate differences in test procedures and sub-population structure; for instance, some data is reported for a single part of a two-part test, or results are reported for each gender separately (adjustment procedure is described in more detail in Appendix B). Once these adjustments have been completed, it is evident that Iowa State undergraduates scored at about the same level as other similar demographics. In fact, both means and standard deviations of ISU students’ scores are similar to the comparison groups, which were chosen from available demographic groups based on population similarity.

Comparison population data was chosen to most closely match ISU undergraduate population demographics. Thus, if comparison data was available for 9th and 12th grade students, we have compared Iowa State students’ scores with the 12th grade students, as they are closer in age to college students. When data was available from college students and Army enlistees, we have compared ISU students to other college students, as college students are more likely to have similar gender distribution to ISU students.

*We may want to include a table of scaled scores here, so that we don’t just bust right into lineup discussion without contextualizing everything properly.*

Applying the grading protocol discussed in section 2.1, we see that the ranges of lineup and visuospatial test scores do not include zero; this indicates that we do not see random guessing from participants in any task. Figure 6 shows the range of possible scores and the observed score distribution.

<sup>2</sup>Averages calculated assuming 294 males and 323 females.

<sup>3</sup>Data from Part I only.

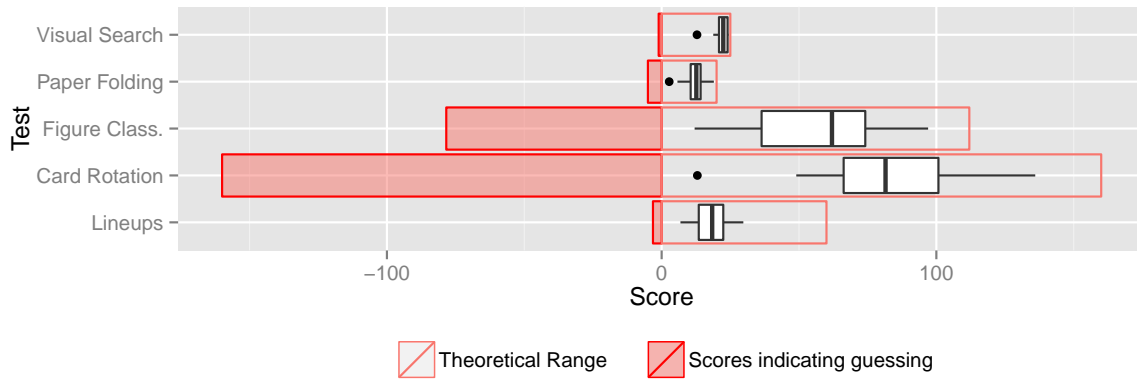


Figure 6: Test scores for lineups and visuospatial tests. As none of the participants scored at or below zero, we can conclude that there is little evidence of random guessing. We also note the score compression that occurs on the Visual Search test; this indicates that most participants scored extremely high, and thus, participants' scores are not entirely representative of their ability.

OK, so on to the next steps ... we would like to figure out, how the lineups play into the established test scores. I think we need to talk about how to best talk through the visual search score. We have a lot of information on the other three scores - is there an established test in the naval test that is 'like' the visual search? - i.e. could we use those scores to compare to?

so outline for the next few steps:

1. some general words on behavior of lineup scores by themselves - I think the whole discussion that is at the moment at the end of the paper should move up here. The levels of the video gaming are a bit messed up - they should be ordered according to number of games played rather than alphabetically :) - and I think that then the relationship actually is a monotonic one ... In any case, sort the figures in the table according to their  $p$ -values in the table. Gender shows up as significantly different in my table of  $p$  values, but in the text you say that it is not. The table with the t-tests should probably go into an appendix ....
2. figure 8 with either three or four scatterplots depending on how the visual search scores fit in
3. a paragraph on discussing the multicollinearity between the test scores based on correlations (either, again, including visual search or not) - this should lead to the conclusion that we cannot directly use a linear model of lineup in test scores.
4. PCA of test scores, with a bit of discussion
5. linear model of lineup scores in PCA with discussion and tying it back to raw scores (figure classification really is the best predictor for lineups, but it only explains some of the variance).

**Lineup Performance and Demographic Characteristics** Completion of Calculus I is associated with increased performance on lineups. This may be related to general math education level, or it may be that success in both lineups and calculus requires certain visual skills. This association is consistent with (Shah and Carpenter, 1995), who found an association between mathematical ability and performance on simple graph description tasks. There is also a significant association between hours of video games played per week and score on lineups, however, this association is not monotonic and the groups do not have equal sample size, so the conclusion may be suspect. There is a (nearly) significant difference between male and female

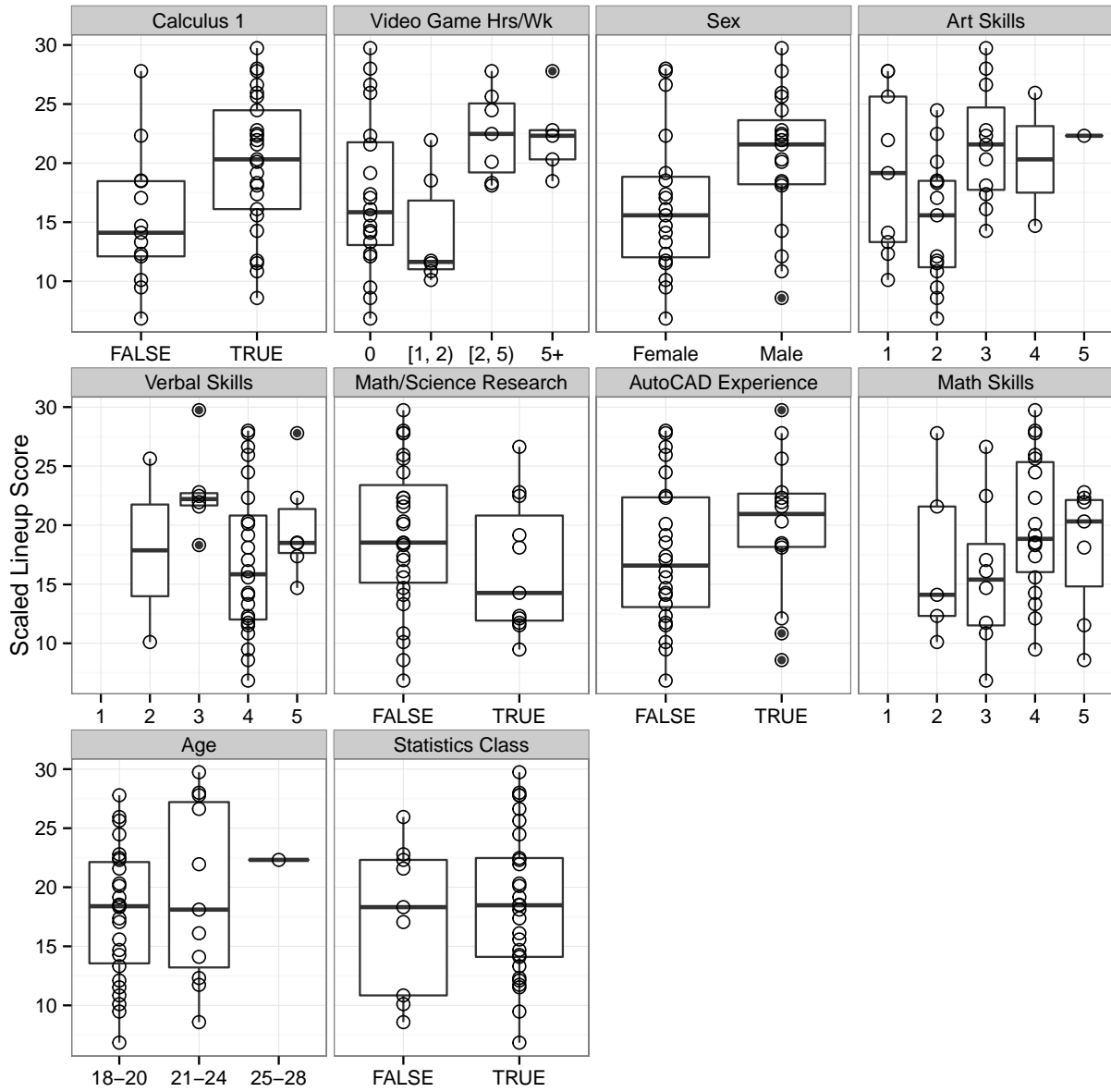


Figure 7: Sample demographic characteristics compared with lineup score. Categories are ordered by effect size; calculus completion, hours spent playing video games per week, and sex are all associated with a difference in lineup score.



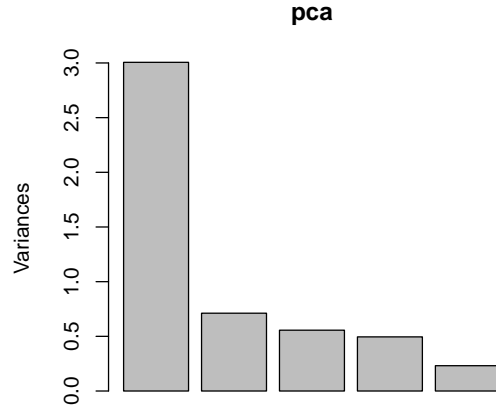


Figure 8: Scree plot of principle component analysis of performance on the different test batteries.

performance on lineups; this is not particularly surprising, since men perform better on many spatial tests (Voyer et al., 1995) and performance on spatial tests is correlated with phase of the menstrual cycle in women (Hausmann et al., 2000). There is no significant difference in lineup performance for participants of different age, self-assessed skills in various domains, previous participation in math or science research, completion of a statistics class, or experience with AutoCAD. These demographic characteristics were chosen to account for life experience and personal skills which may have influenced the results. Statistical test results are available in appendix C.

|            | lineup | card.rot | fig.class | folding | vis.search |
|------------|--------|----------|-----------|---------|------------|
| lineup     | 1.000  | 0.505    | 0.512     | 0.471   | 0.363      |
| card.rot   | 0.505  | 1.000    | 0.474     | 0.705   | 0.609      |
| fig.class  | 0.512  | 0.474    | 1.000     | 0.539   | 0.397      |
| folding    | 0.471  | 0.705    | 0.539     | 1.000   | 0.405      |
| vis.search | 0.363  | 0.609    | 0.397     | 0.405   | 1.000      |

Table 2: Correlation matrix for the five tests.

|                        | PC1    | PC2    | PC3    | PC4    | PC5    |
|------------------------|--------|--------|--------|--------|--------|
| Standard deviation     | 1.7337 | 0.8435 | 0.7458 | 0.7036 | 0.4811 |
| Proportion of Variance | 0.6011 | 0.1423 | 0.1112 | 0.0990 | 0.0463 |
| Cumulative Proportion  | 0.6011 | 0.7435 | 0.8547 | 0.9537 | 1.0000 |

Table 3: Importance of components in principal component analysis of the five tests.

```
biplot(pca, choices = 1:2, pc.biplot = T, cex = c(0.5, 1), adj = 0.75)
biplot(pca, choices = 3:4, pc.biplot = T, cex = c(0.5, 1), adj = 0.75)
```

In figure 11, we see that participant performance on lineups is positively correlated with performance on card rotation, figure classification, and paper folding tasks. This suggests that skills associated with visual

|            | PC1  | PC2   | PC3   | PC4   | PC5   |
|------------|------|-------|-------|-------|-------|
| lineup     | 0.42 | 0.49  | -0.46 | 0.60  | -0.10 |
| card.rot   | 0.50 | -0.30 | 0.28  | 0.23  | 0.73  |
| fig.class  | 0.43 | 0.45  | -0.15 | -0.75 | 0.18  |
| folding    | 0.47 | 0.07  | 0.68  | 0.04  | -0.56 |
| vis.search | 0.41 | -0.69 | -0.48 | -0.15 | -0.33 |

Table 4: PCA Rotation matrix for all five tests.

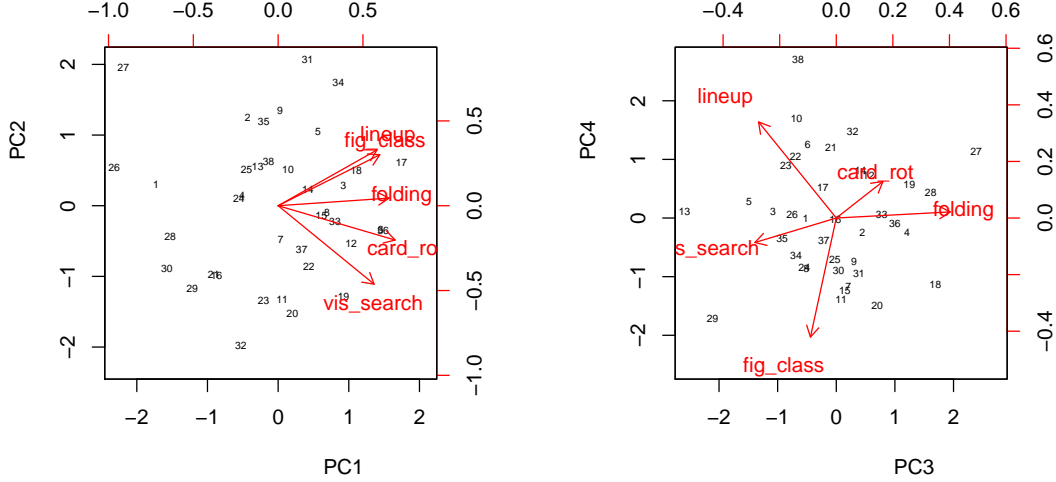


Figure 9: Plots of the principal components with observations. PC1 seems to represent a general “visual intelligence” factor. Figure classification is strongly related to lineups (PC2, PC3), and visual search is related to lineup performance as well (PC3). Paper folding does not seem to be strongly related to lineup performance outside of the first principal component; card rotation is only positively associated with lineup performance in PC4.

reasoning ability are related to lineup performance. As participants must use the same skills in lineups (mental rotation, classification and determining categorization schemes, and multi-step spatial reasoning) as in the factor-referenced tests, this is not particularly surprising. In addition, there seems to be some positive relationship between a participant’s score on the visual search task and their score on lineups: the visual search task represents a baseline of a participant’s ability to find a matching pattern, while lineups require that task as well as the ability to determine what the pattern is for a particular graph. Even excluding the one low visual search score that is a high-leverage point, there seems to be a positive relationship between a participant’s score on lineups and their score for visual search.

Figure 7 shows participants’ responses to the questionnaire given at the beginning of the study; these demographic questions allow us to compare the participants in our study to the undergraduate population of Iowa State as well as to explore relationships between demographic characteristics (major, research experience, etc.) and score on various sections of this test.

All results and data shown here are done in accordance with IRB # 13-581.

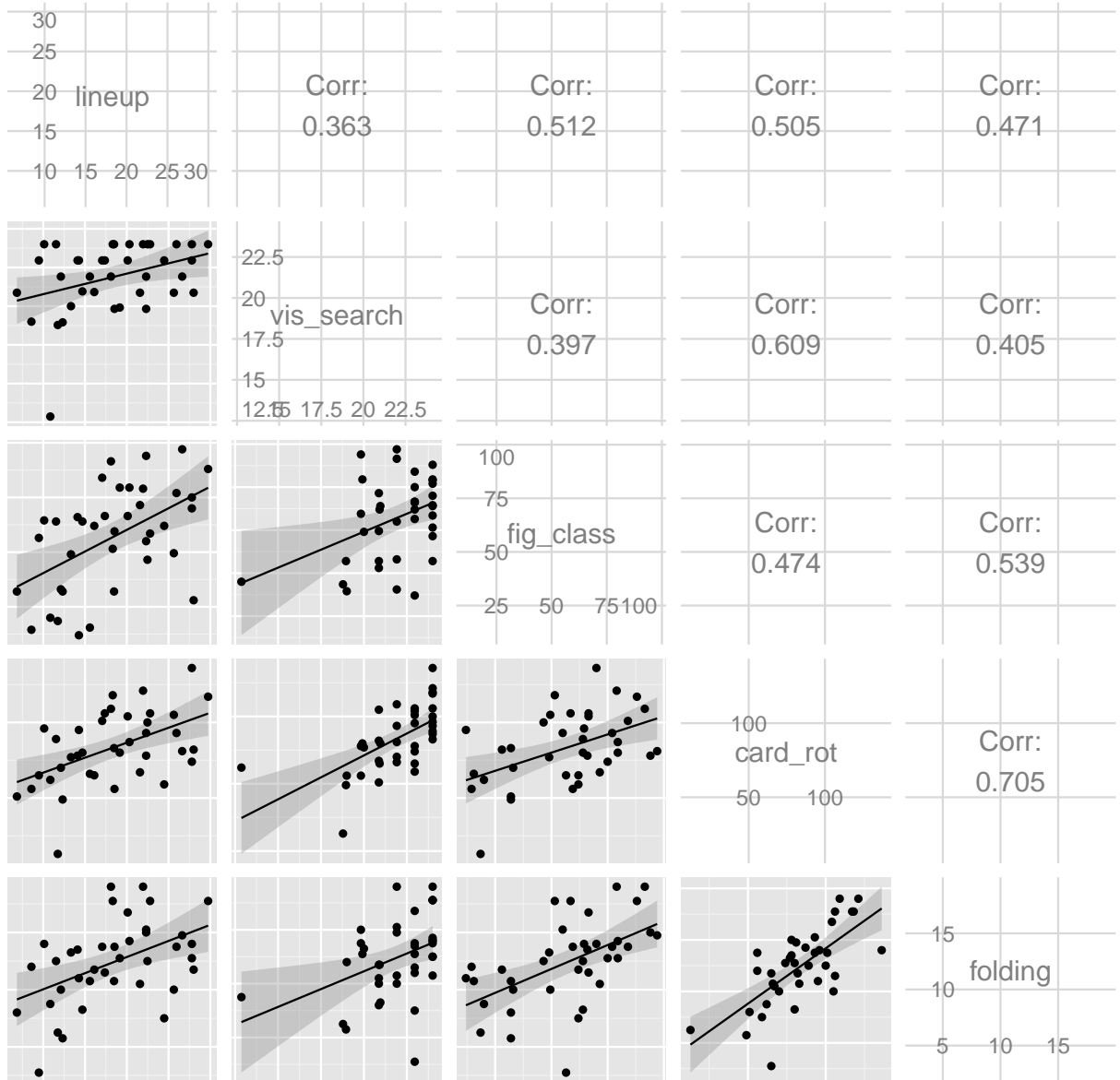


Figure 10: Pairwise scatterplots of test scores. Lineup scores are most highly correlated with figure classification scores, and are also highly correlated with card rotation scores. Paper folding and card rotation scores are also highly correlated.

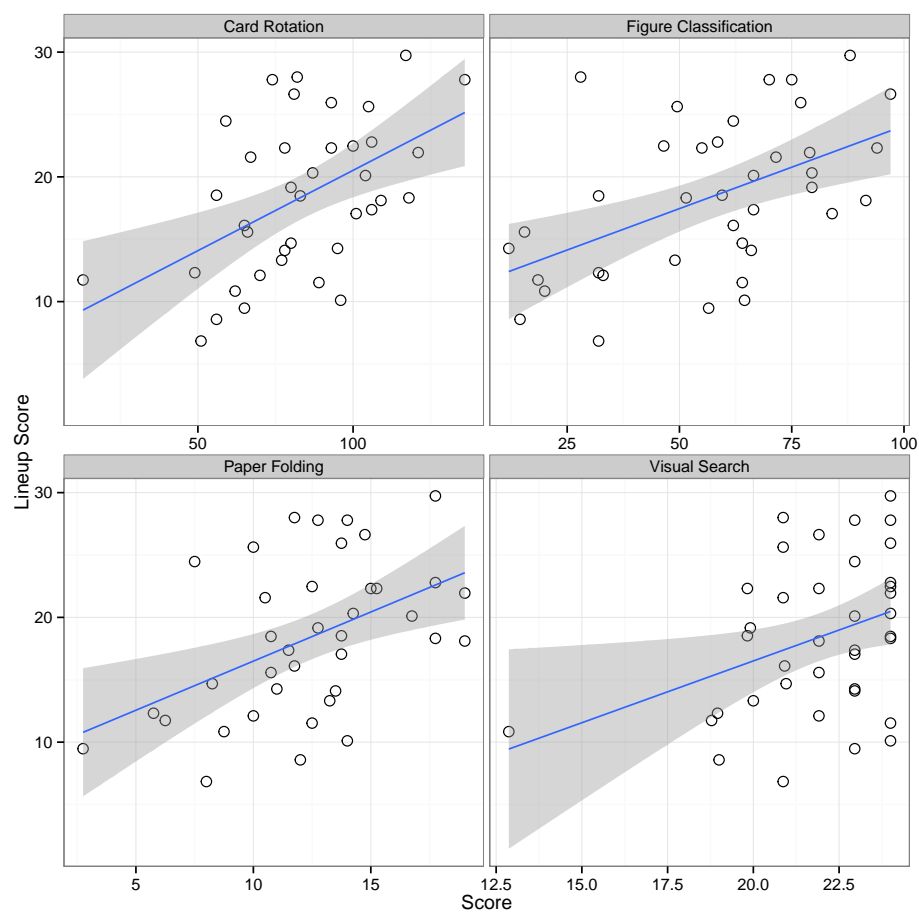


Figure 11: Scatterplots of all test scores compared to participants' scores in the lineup tests. There is a relatively strong positive correlation between lineup score and scores on visuospatial reasoning tests.

## References

- Anderson, K. J. and Revelle, W. (1983). The interactive effects of caffeine, impulsivity and task demands on a visual search task. Personality and Individual Differences, 4(2):127–134.
- Buja, A., Cook, D., Hofmann, H., Lawrence, M., Lee, E.-K., Swayne, D. F., and Wickham, H. (2009). Statistical inference for exploratory data analysis and model diagnostics. Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 367(1906):4361–4383.
- DeMita, M. A., Johnson, J. H., and Hansen, K. E. (1981). The validity of a computerized visual searching task as an indicator of brain damage. Behavior Research Methods & Instrumentation, 13(4):592–594.
- Diamond, J. and Evans, W. (1973). The correction for guessing. Review of Educational Research, pages 181–191.
- Ekstrom, R. B., French, J. W., Harman, H. H., and Dermen, D. (1976). Manual for kit of factor-referenced cognitive tests. Educational Testing Service, Princeton, NJ.
- French, J. W., Ekstrom, R. B., and Price, L. A. (1963). Kit of reference tests for cognitive factors. Educational Testing Service, Princeton, NJ.
- Goldstein, G., Welch, R. B., Rennick, P. M., and Shelly, C. H. (1973). The validity of a visual searching task as an indicator of brain damage. Journal of consulting and clinical psychology, 41(3):434.
- Hampson, E. (1990). Variations in sex-related cognitive abilities across the menstrual cycle. Brain and cognition, 14(1):26–43.
- Hausmann, M., Slabbekoorn, D., Van Goozen, S. H., Cohen-Kettenis, P. T., and Güntürkün, O. (2000). Sex hormones affect spatial abilities during the menstrual cycle. Behavioral neuroscience, 114(6):1245.
- Hofmann, H., Follett, L., Majumder, M., and Cook, D. (2012). Graphical tests for power comparison of competing designs. Visualization and Computer Graphics, IEEE Transactions on, 18(12):2441–2448.
- Just, M. A. and Carpenter, P. A. (1985). Cognitive coordinate systems: accounts of mental rotation and individual differences in spatial ability. Psychological review, 92(2):137.
- Larkin, J. H. and Simon, H. A. (1987). Why a diagram is (sometimes) worth ten thousand words. Cognitive science, 11(1):65–100.
- Lowrie, T. and Diezmann, C. M. (2007). Solving graphics problems: Student performance in junior grades. The Journal of Educational Research, 100(6):369–378.
- Mayer, R. E. and Sims, V. K. (1994). For whom is a picture worth a thousand words? extensions of a dual-coding theory of multimedia learning. Journal of educational psychology, 86(3):389.
- Moerland, M., Aldenkamp, A., and Alpherts, W. (1986). A neuropsychological test battery for the apple II-e. International journal of man-machine studies, 25(4):453–467.
- Scaife, M. and Rogers, Y. (1996). External cognition: how do graphical representations work? International journal of human-computer studies, 45(2):185–213.
- Schaie, K. W., Maitland, S. B., Willis, S. L., and Intrieri, R. C. (1998). Longitudinal invariance of adult psychometric ability factor structures across 7 years. Psychology and aging, 13(1):8.
- Shah, P. and Carpenter, P. A. (1995). Conceptual limitations in comprehending line graphs. Journal of Experimental Psychology: General, 124(1):43.
- Shah, P., Mayer, R. E., and Hegarty, M. (1999). Graphs as aids to knowledge construction: Signaling techniques for guiding the process of graph comprehension. Journal of Educational Psychology, 91(4):690.

- Vekiri, I. (2002). What is the value of graphical displays in learning? *Educational Psychology Review*, 14(3):261–312.
- Voyer, D., Voyer, S., and Bryden, M. P. (1995). Magnitude of sex differences in spatial abilities: a meta-analysis and consideration of critical variables. *Psychological bulletin*, 117(2):250.
- Zhang, J. (1997). The nature of external representations in problem solving. *Cognitive science*, 21(2):179–217.
- Zhang, J. and Norman, D. A. (1994). Representations in distributed cognitive tasks. *Cognitive science*, 18(1):87–122.

## Appendix

### A Examining potential testing biases

T-tests of results for Hillary and Stephanie:

| Variable              | Mean.Hillary | Mean.Stephanie | t    | df    | p-value | Diff.LB | Diff.UB |
|-----------------------|--------------|----------------|------|-------|---------|---------|---------|
| Card Rotation         | 87.72        | 79.45          | 1.06 | 35.99 | 0.29    | -7.51   | 24.05   |
| Paper Folding         | 13.18        | 11.69          | 1.27 | 35.93 | 0.21    | -0.90   | 3.88    |
| Figure Classification | 65.17        | 49.60          | 2.16 | 32.86 | 0.04    | 0.87    | 30.26   |
| Lineups               | 20.17        | 16.77          | 1.74 | 35.66 | 0.09    | -0.55   | 7.34    |

### B Scaling Scores

To calculate “scaled” comparison scores between tests which included different numbers of sections, we scaled the mean in direct proportion to the number of questions (thus, if there were two sections of equivalent size, and the reference score included only one of those sections, we multiplied the reported mean score by two). The variance calculation is a bit more complicated: In the case described above, where the reference section contained half of the questions, the variance is multiplied by two, causing the standard deviation to be multiplied by approximately 1.41.

|                 | Card Rotation                      | Paper Folding                         | Figure Classification   |
|-----------------|------------------------------------|---------------------------------------|---|
| ISU Students    | 83.4 (24.1)                        | 12.4 (3.7)                            | 57 (23.8) <sup>4</sup>  |
| Scaled Scores   | 88.0 (34.8)                        | 13.8 (4.5)                            | 58.7 (14.4) <sup>5</sup>  |
| Unscaled Scores | 44.0 (24.6) <sup>6</sup>           | 13.8 (4.5)                            | M: 120.0 (30.0), F: 114.9 (27.8)  |
| Population      | approx. 550 male<br>naval recruits | 46 college students<br>(1963 version) | suburban 11th & 12th grade students<br>(288-300 males, 317-329 females) |

Table 5: Comparison of scores from Iowa State students and scores reported in Ekstrom et al. (1976). Scaled scores are calculated based on information reported in the manual, scaled to account for differences in the number of questions answered during this experiment. Data shown are from the population most similar to ISU students, out of the data available.

This scaling gets slightly more complicated for scores which have two sub-groups, as with the Figure Classification test. To get a single unified score with standard deviation, we did the following calculations:

$$\mu_{all} = (\mu_F N_F + \mu_M N_M) / (N_F + N_M) \quad (2)$$

$$\sigma_{all} = \sqrt{(N_F \sigma_F^2 + N_M \sigma_M^2) / (N_F + N_M)} \quad (3)$$

<sup>4</sup>ISU students took only Part I due to time constraints.

<sup>5</sup>Averages calculated assuming 294 males and 323 females.

<sup>6</sup>Data from Part I only.

Then, in order to account for the fact that ISU students took only part I of two parts to the Figure Classification test (and thus completed half of the questions), we adjusted the transformation as follows

$$\mu_{all} = 1/2(\mu_F N_F + \mu_M N_M)/(N_F + N_M) \quad (4)$$

$$\sigma_{all} = \sqrt{(N_F \sigma_F^2 + N_M \sigma_M^2)/(2(N_F + N_M))} \quad (5)$$

## C Lineup Performance and Demographic Characteristics

| Variable              | DF | Sum.of.Squares | Mean.Squared | F.value | p.value |
|-----------------------|----|----------------|--------------|---------|---------|
| Calculus 1            | 1  | 204.569        | 204.569      | 6.15    | 0.018   |
| Video Game Hrs/Wk     | 3  | 326.542        | 108.847      | 3.44    | 0.028   |
| Sex                   | 1  | 140.844        | 140.844      | 4.02    | 0.053   |
| Art Skills            | 4  | 303.563        | 75.891       | 2.28    | 0.082   |
| Verbal Skills         | 3  | 180.660        | 60.220       | 1.68    | 0.191   |
| Math/Science Research | 1  | 59.670         | 59.670       | 1.60    | 0.214   |
| AutoCAD Experience    | 1  | 50.893         | 50.893       | 1.36    | 0.252   |
| Math Skills           | 3  | 111.117        | 37.039       | 0.98    | 0.416   |
| Age                   | 2  | 41.445         | 20.723       | 0.53    | 0.592   |
| Statistics Class      | 1  | 9.062          | 9.062        | 0.23    | 0.631   |

Table 6: Model results of each demographic variable compared with lineup score. Multiple testing issues aside, it appears that very few demographic variables (if any) are significantly associated with score on lineups among Iowa State undergraduate students.

## D Lineup Performance: Accounting for “Visual Intelligence”

As figure 9(a) shows, the first principal component seems to primarily describe “general intelligence”, or “visual intelligence” in a generic sense (as we have only tested visual skills). As the visual search task (VST) is the most generic of visual tasks (find the matching pattern) and requires no spatial manipulation, rotation, or inference, we might use it as an indicator of general visual intelligence (or cognitive tempo and attention to detail).

```
model1 <- lm(lineup ~ vis_search, data = ans.summary)
```

|             | Estimate | Std. Error | t value | Pr(>  t ) |
|-------------|----------|------------|---------|-----------|
| (Intercept) | -3.3051  | 9.3322     | -0.35   | 0.7253    |
| vis.search  | 0.9907   | 0.4242     | 2.34    | 0.0252    |

Table 7: Results of the linear regression of visual search score on lineup score. .

```
# partial regression data set
ans.part <- ans.summary
ans.part$lineup <- resid(lm(lineup ~ vis_search, data = ans.summary))
ans.part$card_rot <- resid(lm(card_rot ~ vis_search, data = ans.summary))
ans.part$fig_class <- resid(lm(fig_class ~ vis_search, data = ans.summary))
ans.part$folding <- resid(lm(folding ~ vis_search, data = ans.summary))
```

Removing paper folding from the analysis due to its’ obvious nonsignificance and collinearity with card rotation, we get a slightly better regression result.

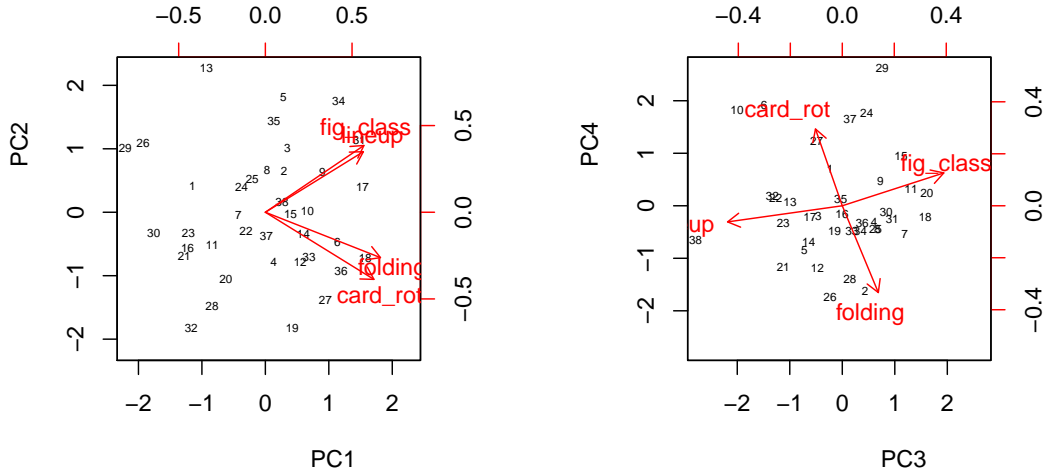


Figure 12: Biplots of the adjusted lineup, card rotation, figure classification, and paper folding scores principal components analysis. Adjusted lineup scores are most associated with figure classification in the second principal component. In the third and fourth components, adjusted scores are almost entirely orthogonal.

|             | Estimate | Std. Error | t value | Pr(>  t ) |
|-------------|----------|------------|---------|-----------|
| (Intercept) | 0.0000   | 0.8362     | 0.00    | 1.0000    |
| fig.class   | 0.0831   | 0.0436     | 1.91    | 0.0652    |
| card.rot    | 0.0665   | 0.0572     | 1.16    | 0.2526    |
| folding     | 0.1664   | 0.3451     | 0.48    | 0.6328    |

Table 8: Results of the linear regression of (the ETS tests | VST) on (lineup score | VST)

Comparing the full and reduced models, we see that the reduction does not significantly reduce the explained sum of squares.

Even using VST as a “general intelligence” factor, the first PC still does not discriminate between the tests. This approach is probably not worth pursuing further.

## E Lit Review Rejects

- [Just and Carpenter \(1985\)](#) showed that high-spatial-ability viewers used different rotation strategies than low-spatial-ability viewers when asked to whether three-dimensional alphabet cubes were the same.
- [Voyer et al. \(1995\)](#) completed a meta-analysis of spatial tasks and gender differences. They concluded that “rotation” tasks, such as the card rotation test, have robust sex differences across cohort, and while the sex differences seem to be declining for some tests, such as the card rotation test, the differences still exist for the time being. Test administration differences may account for the change. They also concluded that the paper-folding test showed no such sex difference, though this may be dependent on test scoring (a lower guessing penalty decreases the score difference).

*In the current analysis, we have scored the paper folding test out of 20. Alternate scoring procedures still need to be investigated... they cite a poster, but it's not available anywhere online as far as I can tell.*



|             | Estimate | Std. Error | t value | Pr(>  t ) |
|-------------|----------|------------|---------|-----------|
| (Intercept) | 0.0000   | 0.8270     | 0.00    | 1.0000    |
| fig.class   | 0.0902   | 0.0405     | 2.23    | 0.0326    |
| card.rot    | 0.0824   | 0.0462     | 1.79    | 0.0829    |

Table 9: Results of the reduced linear regression of (the ETS tests | VST) on (lineup score | VST)

|   | Res.Df | RSS    | Df | Sum of Sq | F    | Pr(> F) |
|---|--------|--------|----|-----------|------|---------|
| 1 | 35     | 909.52 |    |           |      |         |
| 2 | 34     | 903.34 | 1  | 6.18      | 0.23 | 0.6328  |

Table 10: Model comparison between full model (including paper folding) and reduced model (excluding paper folding), when accounting for scores on the VST in both independent and dependent variables.

- [Lowrie and Diezmann \(2007\)](#) investigated students (9-10 yrs) ability to read mathematical graphs and their visual-spatial ability. They partition graphs into 6 “graphical languages” (some of these languages aren’t statistical graphs - networks, venn diagrams, maps), and use multivariate analysis to associate spatial ability with some of the types of graphs studied.

*I can’t understand their model results very well, so I’m not entirely sure what’s going on here. They also cite a bunch of studies (including [Voyer et al. \(1995\)](#)) that I haven’t had time to completely explore.*

- [Vekiri \(2002\)](#) is lit review discussing the different theories about the utility of graphics in learning (dual coding theory, visual argument hypothesis, and conjoint retention hypothesis). The article contains a nice definition of graphics (though it might be too domain specific for statistical graphics, since it claims graphics have only one meaning and aren’t prone to interpretation.) Depending on the context, “maps”, “charts”, and “graphs” here could fall into a loose classification as “statistical graphics”, but in general “graphs” is probably accurate for most of the lineups we’re testing: “[Referents are ] quantitative data ... that enable viewers to compare and observe relations among variables”.

- **Dual Coding Theory:** people encode both words and pictures separately, so both are encoded for a graph.

Supports the use of graphics, and there is support for dual-coding in working memory research (pictures and words are processed separately but in parallel). In addition, since graphics depend on other abilities (spatial abilities), graphics that aren’t properly designed create higher cognitive load and lower performance on graphical tasks. The review also specifically discusses visuospatial ability, citing a speculation from [Mayer and Sims \(1994\)](#) that students with low spatial ability perform worse with diagrams and graphs because more cognitive resources are required.

*Lots of studies by Sweller cited through the cognitive load discussion - I haven’t gotten to that literature yet.  
That said, that research has some potential to justify the use of lineups from a psychological perspective - if you can claim that types of graphs that have poor performance do so because they create higher cognitive load (cue psych research into why that would be true, cleveland & mcgill, etc.), you can justify the suckiness of those awful boxplot-jittered-things.*

- **Visual Argument Hypothesis:** Graphics require less processing because their visuospatial properties encode some of the information through the arrangement of elements as well as the elements themselves. That is, graphics are a form of “external cognition” ([Scaife and Rogers, 1996](#)) that guide, constrain, and facilitate cognitive behavior([Zhang, 1997](#)). This constraint reduces memory load and makes more cognitive resources available for other tasks; displays that clearly present information without deep processing requirements are more effective([Zhang and Norman,](#)

1994). Graphical displays have higher search and computational efficiency than text (Larkin and Simon, 1987), and when graphics are designed according to gestalt principles of organization, where grouping lines up with visual chunks, gestalt heuristics can be used to more quickly encode information (Shah et al., 1999).

*I am working on making my way through all of this literature - right now, I'm re-citing stuff; I will print the papers out and add to this ASAP once I get toner into the printer.*

- **Conjoint Retention Hypothesis** - based on dual coding theory, but claims that there are separate but interconnected memory codes for verbal and visual information. Assumes maps/graphs are encoded intact, with visuospatial properties preserved. This is pretty unlikely, and the literature is not so relevant to our cause. Spatial information is encoded, but words can actually interfere with this encoding; it is not “intact”
- Mayer and Sims (1994) used the card rotation test (10 figures, 80 questions total) and the paper folding test (10 items) to evaluate spatial performance, along with a problem solving test. They used the same penalty structure on the card rotation test, and don't seem to have used any penalty on the paper folding test; they then scaled the two tests. Their problem solving test was more practical (engineering-based), so isn't applicable to our study. They were testing instruction using visual + audio in sync, visual + audio out of sync, and then no instruction as a control. High-ability students could make more use of the in-sync presentation of visual and auditory information than low-ability students, but did not differ significantly in the out-of-sync information condition. Thus, visualization ability moderates the results of the information presented by allowing students to process information concurrently more efficiently, but doesn't help students process the information separately or help students come up with solutions to the problem-solving task.
- Hofmann et al. (2012) for lineup stimuli and general lineup performance