



# Examining Student Conceptions of Covariation: A Focus on the Line of Best Fit

[Stephanie A. Casey](#)

Eastern Michigan University

*Journal of Statistics Education* Volume 23, Number 1 (2015),

[www.amstat.org/publications/jse/v23n1/casey.pdf](http://www.amstat.org/publications/jse/v23n1/casey.pdf)

Copyright © 2015 by Stephanie A. Casey, all rights reserved. This text may be freely shared among individuals, but it may not be republished in any medium without express written consent from the author and advance notification of the editor.

---

**Key Words:** Statistics Education; Statistical Association; Linear Regression.

## Abstract

The purpose of this research study was to learn about students' conceptions concerning the line of best fit just prior to their introduction to the topic. Task-based interviews were conducted with thirty-three students, focused on five tasks that asked them to place the line of best fit on a scatterplot and explain their reasoning throughout the process. The results include descriptions and categorizations of students' meanings, criteria and methods for placement, accuracy of placement, and interpretation of the line of best fit. The discussion addresses how students' prior study of mathematics and statistics impacted their conceptualizations of the line of best fit as well as implications for the teaching and learning of the line of best fit.

## 1. Introduction

One of the fundamental statistical ideas in school curricula is statistical association ([Burrill and Biehler 2011](#); [Garfield and Ben-Zvi 2004](#)), also known as covariation. Students are typically introduced to statistical association through the study of the line of best fit, as it is a natural extension of their concurrent study of linear equations and functions in mathematics (e.g., Australia: [Australian Curriculum, Assessment, and Reporting Authority 2012](#); England: [Qualifications and Curriculum Authority 2007](#); U.S.A.: [Common Core State Standards Initiative \(CCSSI\) 2010](#)). For example, the authors of the Common Core State Standards in Mathematics (CCSS-M) ([CCSSI 2010](#)) state that students in grade eight in the United States should begin studying this topic:

8.SP.A.2 [Students should] Know that straight lines are widely used to model relationships between two quantitative variables. For scatter plots that suggest a linear association, [students will] informally fit a straight line, and informally assess the model fit by judging the closeness of the data points to the line (p.56).

Informal fitting of the line refers to the idea that students are fitting a line to data displayed in a scatterplot by eye, without using calculations or technology to place the line for them. The line placed in this manner is noted as the informal line of best fit. The [Common Core Standards Writing Team \(2011\)](#) expanded on the purpose of this standard to include an expectation that students in grade eight determine that the informal line of best fit for data that has no association should be a horizontal line, and that a horizontal fitted line implies that there is no association between the variables.

Students' study of the line of best fit through this informal approach initiates the developmental progression of educational experiences through which students need to work in order to understand linear regression ([Bargagliotti et al. 2012](#)). Such a progression is endorsed by recent curriculum documents such as the Guidelines for Assessment and Instruction in Statistics Education (GAISE) Report ([Franklin et al. 2005](#)) and the CCSS-M ([CCSSI 2010](#)). The learning of informal regression lays the foundation for students' future study of formal regression (e.g., least-squares regression line, median-median line) in later coursework, including mathematics ([CCSSI 2010](#)), science ([Next Generation Science Standards \(NGSS\) Lead States 2013](#)), and social studies ([National Council for the Social Studies \(NCSS\) 2010](#)).

Informally finding a line of best fit, where students are naturally developing their own methods, criteria, and properties for the line of best fit, has been included in curriculum standards for more than two decades (e.g., [National Council of Teachers of Mathematics \(NCTM\) 1989](#); [NCTM 2000](#)), yet students' conceptions of line of best fit are largely undocumented. A review of the literature found only one study that examined university students' conceptions of the line of best fit ([Sorto, White, and Lesser 2011](#)). Most significantly, empirical research has not examined what conceptions are held by younger (ages 11-14) students who are studying this topic for the first time. Unanswered questions include: What do students conceive the line of best fit to be? What are students' natural inclinations regarding methods for finding the line of best fit? What criteria do students develop for the line of best fit, and do they agree with the [Common Core State Standards Initiative's \(2010\)](#) approach of "closeness of the data points to the line" (a precursor to the criterion of the least-squares regression line)? What are the sources of students' conceptions regarding the line of best fit? Consequently, the purpose of this study is to answer these questions through studying students' conceptions concerning the line of best fit just prior to their introduction to the topic.

Mathematics teachers and curriculum writers need to know how students naturally think about specific topics taught in the mathematics classroom. [Ball, Thames, and Phelps \(2008\)](#) termed this 'knowledge of content and students' and deemed it an essential component of mathematical knowledge for teaching since effective teaching helps students learn through the process of building more advanced knowledge from their prior understandings ([Bransford, Brown, and Cocking 1999](#)). Oftentimes, students' conceptions differ in fundamental ways from the concepts that we try to teach ([Confrey 1990](#)), so research regarding student conceptions can determine

standard stumbling blocks that occur during the learning of a topic. For the topic of statistical association in particular, research has found that students' conceptions are resistant to change ([Batanero, Estepa, Godino, and Green 1996](#)). This accentuates the need for research regarding students' conceptions for topics that teach statistical association, including the line of best fit. Student conceptions described by studies such as this one enable teachers' to anticipate ideas students typically bring to the table when they learn about a certain topic. Teachers' anticipation of likely student responses is a key teacher practice in leading productive mathematical discussions, an instructional practice currently encouraged in order to prepare students for demands of the twenty-first century ([Smith and Stein 2011](#)).

## 2. Related Literature

### 2.1 Student knowledge of line of best fit

The current study builds on research by [Sorto et al. \(2011\)](#) who investigated university students' methods for drawing the line that best fit the data and whether these methods agree with the method of least-squares. The most relevant part of this study for the purposes of the present study was its methodology. The assessment tasks used by Sorto et al. framed the thinking about the tasks that were designed for the task-based interview method utilized in the present study. The participants in their study were presented with three assessment tasks on a piece of paper and asked to complete them. The first two tasks centered on a scatterplot of eight data points that displayed a positive association between job satisfaction scores and salary. The participants were asked to draw a line on the scatterplot which best fits the data and to describe the criterion they used to determine the line. The third assessment task presented side-by-side graphs with identical data points plotted but different lines of best fit. The task asked "If you had to decide which of the two lines fits the data better, how would you decide? Is there anything you could compute?" (p.50). Sorto et al. found that this task forced the participants to focus on measuring the goodness of fit in order to compare the two lines. Looking across the three tasks, the participants were placing the informal line of best fit according to the following four criteria: equal number of points on both sides of the line; at the middle or average of all the points; fit the trend of the data; and consideration of the distance of the points to the line. This is the only known study to address students' knowledge of the line of best fit; however, it was done with university level students who are distinctly different in important ways (e.g., mathematical sophistication, exposure to line of best fit) from younger (e.g., ages 11-14) students who are learning this topic for the first time.

### 2.2 Students learning to reason about data

Over the last twenty years, statistics education researchers have learned a considerable amount about how young students learn to reason about data. Some of these research findings were relevant to the data analysis procedures of this study and are described below.

Students new to the study of data tend to focus their attention on individual cases in the data set and perceive data as a series of individual cases (case-oriented view) rather than considering the whole data set holistically (aggregate view) as experts do ([Bakker 2004](#)). For example, when experts with an aggregate view look at a histogram of a population's heights, they tend to focus

on its general shape, how spread out it is, and where the data are centered to inform their analysis of that population's heights. A person with a case-oriented view, however, would tend not to see these aggregate features. This person might focus on the tallest or shortest individual in the population. Another person with a case-oriented view might group the individuals into short, usual, and tall groupings then reason about each of these groupings separately without ever considering the data set collectively as a group. [Estepa and Batanero \(1996\)](#) documented the prevalence of the case-oriented view by seniors in secondary school when reading scatterplots to judge correlations between quantitative variables. In their hypothetical learning trajectory for reading scatter plots, [Cobb, McClain, and Gravemeijer \(2003\)](#) hypothesized that students need to develop an aggregate view of data in order to see the shape associated with a bivariate data set and reason about any pattern that exists within.

A robust finding that has emerged from multiple studies is that peoples' prior beliefs about the relationship between two variables has a great deal of influence on their judgments of the covariation between those variables (e.g., [Batanero et al. 1996](#); [Jennings, Amabile, and Ross 1982](#); [Kuhn, Amsel, and O'Loughlin 1988](#)). This is a specific example of a larger psychological phenomenon known as confirmation bias ([Nickerson 1998](#)). Confirmation bias describes the propensity of humans to look for or interpret evidence in a manner that will confirm their existing expectation regarding what the evidence should show. In the case of bivariate data analysis in statistics, if the data provided do not agree with a student's expectation about the relationship between the variables based on what the student knows about them from context, he or she often misinterprets or disregards the data. Students also tend to view cases that confirm their expectations as more relevant than disconfirming cases ([McGahan, McDougal, Williamson, and Pryor 2000](#)), and often infer that association means a causal relationship exists between the variables ([Batanero, Estepa, and Godino, 1997](#)).

## 2.3 Research questions

The research questions for this study on students' conceptions of line of best fit were:

1. What do students conceive the line of best fit to *be*?
2. What are students' natural inclinations regarding *criteria* the line of best fit is trying to meet and what *methods* do they use to meet them?
3. What are *sources* of students' conceptions?
4. How *accurately do students place* an informal line of best fit, relative to the least-squares regression line?
5. How do students *interpret* a placed informal line of best fit in context?

## 3. Methodology

### 3.1 Participants

All thirty-three of the participants were eighth grade students in the United States. They participated in the study in the fall of 2012 in the week prior to studying the informal line of best fit in their mathematics class, so they were prepared to learn the topic according to the authors of their curriculum. Twenty-one of the participants, eight male and thirteen female, attended a middle school in a medium-sized, Midwestern town. This school was diverse both in terms of

ethnicity (58% white, 15% Asian, 13% African-American, 6% Hispanics) and socio-economic status of its students. Of these twenty-one students, five were in the honors mathematics class and two received special education support services. Another twelve participants, eight male and four female, were students in an honors mathematics class at a predominately white, upper-middle class middle school in the northeast region of the United States.

### 3.2 Method: Task-based interview

A task-based interview was conducted individually with each participating student by the author, a trained graduate research assistant, or a research collaborator. Task-based interviews have been found effective at eliciting student conceptions regarding statistical concepts (e.g., [Konold et al. 2002](#); [Mokros and Russell 1995](#)) and have been used in much of the exemplary research that has been done on student learning ([Confrey 1990](#)). All interviews were video recorded and followed a semi-structured format, which involved a mix of predetermined tasks and questions as well as clarifying follow-up questions created at the time of the interview. Each interview took approximately 20 minutes. The interview protocol was revised based on a pilot study with two students during the summer of 2012. The final interview protocol used in this study, including the tasks, is in [Appendix A](#).

For the predetermined interview protocol, a series of tasks and questions were assembled to elicit students' thinking related to fundamental aspects of the line of best fit. The interview began by establishing that the participating students knew how to read a scatterplot. Next, the interviewer read the following statement that summarizes what the student was asked to do in the first five tasks:

For each graph, we will talk about the data displayed and then I will ask you to determine the line of best fit for the data points. I would like you to think out loud as you decide where to place the line on each graph so that I can understand how you are deciding where to place it.

Let's use this piece of wire as the line so you can move it on the paper and place it where it best fits the data. [Place a wire on the paper in the white space, not on the graph] There are no right or wrong answers. I am interested in how you think about placing the piece of wire, so please tell me what you are thinking as you do this so I can follow your thoughts. Once you decide where you want the wire line on the scatterplot, I will use scotch tape to keep it in place. Go ahead and place the line where you think it should be.

In accordance with the work of [Sorto et al. \(2011\)](#), this type of task was deemed effective at eliciting students' conceptions regarding the line of best fit; thus the protocol for this study included posing five similar tasks to the participants. The students followed this procedure for a series of five scatterplots for contextual data sets with eight ordered pairs each. This study's protocol included five tasks to encompass varying directions and strengths of association within the tasks. Since students better understand covariation when it's presented with data from a meaningful context ([Moritz 2004](#)), all of the tasks were set in contexts familiar to students. Use of a wire or similar straight-edge manipulative was a common strategy used in textbook presentations of this topic (e.g., [Cooke, Heideman, Keene, Lin, and Reeves 2007](#)). Music wire

was used in this study because it was not bendable (hence maintaining its shape as a line) or breakable and was thin (thickness of approximately 0.05 inches) to minimize the issue of obscuring the data points when placed on the scatterplot. In addition to asking the students to talk aloud during their completion of these tasks, clarifying questions regarding the reasoning behind the student's completion of the task were asked by the interviewer as necessary. The interviewer was allowed to restate students' statements but refrained from answering any questions regarding how to complete the task correctly or if the given response was correct.

The findings from the pilot study, a task from the Core-Plus Mathematics Project's textbook "Contemporary Mathematics in Context: A Unified Approach" Course 1 (Coxford et al. 2003), and personal experience teaching this topic, informed initial ideas regarding criteria students might have for the line of best fit. These included going through the maximum number of points, having the same number of points above and below the line, connecting the first and last points, and forcing the line to go through the origin. These criteria informed the task design process, as described below in the explication of each task.

Tasks 1 and 2 displayed data relating the height a golf ball bounces to the height from which it was dropped. The plots for these tasks had positive linear associations. The data in Task 1 had a strong positive linear association ( $r = 0.96$ ) and the least-squares regression line had a y-intercept very near to the origin ( $b_0 = 0.33$ ), which was anticipated to be a comfortable situation for many students and therefore appropriate for the first task. Task 1 also presents a graph that could help identify whether students had a criteria of hitting as many points as possible because there are multiple sets of four points that are close to collinear. If a student took the approach of having equal numbers of points above and below the line then his or her line may be close to the least-squares regression line for this data. Task 2's data also had a strong positive linear association ( $r = 0.97$ ), but the pattern of the points was different than Task 1. Pairs of consecutive points formed collinear pairs with slopes distinctly different from one another. It was anticipated that students may have a criterion of equal number of points above and below the line, which could result in a poor line of best fit for this data set. Also, two of the pairs of points had approximately equal bounce heights despite different drop heights, which provided the students with an opportunity to consider such a situation when determining the line of best fit. It would be reasonable to have the best fit line go through the origin on this plot as well.

Data regarding the yearly attendance at movies in the United States at differing average ticket costs was shown in Tasks 3 and 4. These variables had a negative association with less strength in their association than the golf ball bounce and drop heights (Task 3:  $r = -0.76$ ; Task 4:  $r = -0.92$ ) but still relatively strong associations. It would not be appropriate to have the line of best fit for these tasks cross the origin, a novel scenario in this interview and one purposefully included for students to consider. In Task 3, a student who is trying to make a line that goes through as many points as possible would likely create a poor best fit line, as would a student who wants to have equal number of points above and below the line. The data in Task 4 were less linearly aligned than in the previous plots to provide students with an opportunity to perhaps play around with the placement of the line in a greater number of positions. This was purposefully done to push students to make decisions about what criteria are important to consider. Additionally, neither the "hit as many points as possible" nor "equal number of points below and above the line" methods will work well to model the entire data set.

The plot in Task 5 graphed data regarding the height and shoe size of eight elementary school students in a teacher's class. The variables are not associated ( $r = 0.05$ ), which is contrary to a common assumption of a positive association between these variables from one's personal experience. This was done to determine the level of influence context and beliefs regarding the variables' association have on the placement of the line of best fit. It was also important to determine how students would conceive of the line of best fit in a plot that displays no association.

Following the completion of these five tasks, each student was asked three summary questions. The first question, "Could you tell me what you would say to another student that asked you 'What is the line of best fit?'" was asked to have the student verbalize his or her definition of the line of best fit and better understand the meaning the term 'line of best fit' had for the student. This was followed by a question concerning the student's criteria of the line of best fit: "What would you say to another student to help them draw the line of best fit on a scatterplot?" The third query was designed to prompt the student to be reflective about his or her experience in doing the interview thus far, asking "As you completed the tasks, did your thoughts about where the line of best fit gets placed change? If so, how?"

Next the students were presented with Task 6, an adaption of the third task used by [Sorto et al. \(2011\)](#). The scenario for this task was that two students, Angelo and Barbara, were given the same task the participating students were given in Task 1: to find the line of best fit for the data on a golf ball's drop height and bounce height. Angelo and Barbara had different solutions, presented in side-by-side scatterplots with the best fit lines displayed and equations of the best fit lines given below the plots. The prompt given to each participant was: "Which student's line fits the data better and why?" The task was designed so that Angelo and Barbara's line placement would be similar and neither would go through the origin, which the pilot study indicated might be a criterion; however, Angelo's line goes through two of the points while Barbara's is closest to all of the points (it is the least-squares regression line) but does not go through any points. One purpose of this task was to force students to choose between which of these criteria they felt did a better job of fitting the data. Another was to see whether students would choose the least-squares regression line, the model statisticians have deemed to be the best fit line, as the line that better fits the data. This task was placed in the protocol after the first five tasks and the summary questions so as not to interfere with the students' responses to those queries.

Finally, the students were asked two concluding questions to see if they could broadly interpret their best fit lines in the context of the data. Each student was shown the graphs from Task 1 and Task 3 with their best fit line wires attached, and asked to talk about what the line shows about the relationship between those variables (given by name for each task).

### 3.3 Analysis

Analysis of the data began with the creation of written transcripts for the interviews. The author and research collaborator then went through an iterative process of analysis, moving back and forth between viewing of the data sources (videos, transcripts, and student work) and identifying, discussing, and classifying the student responses. This analysis mode is in accordance with a

grounded theory method ([Glaser and Strauss 1967](#)) that uses empirical data as the basis for developing themes and theories to describe a phenomenon. The core elements from each student's interview were recorded in a spreadsheet and documented the corresponding times in the video each observation was based on. The core elements included in the initial analysis were each student's criteria for fitting the line of best fit on Tasks 1 through 5, statements about the meaning of the line of best fit, including the response to summary question one, and answer to Task 6, including justification. These were initial core elements needed to answer research questions one through three. A strength of the interview protocol was having multiple responses inform each of these research questions. For example, an answer to research question one 'What do students conceive the line of best fit to be?' was predominately provided by a student's response to the first summary question, "Could you tell me what you would say to another student that asked you 'What is the line of best fit?'" but was also informed by responses to the tasks which forced the student to consider the line in different scenarios. Similarly, a student's criteria for fitting the line was informed by his or her criteria utilized on Tasks 1 through 5, response to summary question two "What would you say to another student to help them draw the line of best fit on a scatterplot?," and response to Task 6 (choosing between Angelo's and Barbara's lines as the better line of best fit). During this first round of analysis, notes of interesting or unique occurrences during the interview were made, as well as ideas for further analysis. For example, it was noted that many students referred to measures of center during the interview, sometimes in productive ways, and in other instances in ways that interfered with a correct conception of the line of best fit. It was then decided to include that comprehensively as a part of the analysis during the next round. Following this first round of analysis, the findings thus far were discussed, including themes, ideas for further analysis, and surprising responses. For example, it was evident that students generally struggled to engage with Task 5. Therefore, we decided to isolate that task's line placement and criteria from the other line placement tasks (Tasks 1-4) in the analysis.

During the second round of analysis, all of the transcripts and videotapes were revisited to add additional elements to the analysis that emerged as important from our first round of analysis. The following elements were added to the spreadsheet in the second round of analysis: previous study of scatterplots; reference to a measure of center; case-oriented or aggregate view of the data; univariate or bivariate view of the data; reference to the context of the data and whether it interfered with their approach to the task or not; interpretations of their fitted lines (to answer research question 5); and whether the student made a causal statement when interpreting the line.

In the next phase of analysis, the author synthesized the core elements documented in the spreadsheet. This began with an analysis of the commonalities of the conceptions and criteria of the line of best fit of the participants. A tag was assigned to each conception as well as criteria using the constant comparative method ([Glaser and Strauss 1967](#)) in order to minimize the number of tags created. This allowed overall trends regarding the students' conceptions and criteria to emerge, and associations with categories for other core elements were explored. Each conception and criteria was assigned a tag, not each student, so some students had multiple tags for their conceptions and/or criteria if more than one was identified in the student's responses during the interview.



Lastly, to answer research question four regarding the accuracy of placed lines of best fit, plots were made that displayed all 33 students' lines for Tasks 1 through 5 with the least-squares regression line highlighted in red. Further analysis of these plots was done to determine how criterion categories were related to the accuracy of the placed lines.

## 4. Results

### 4.1 Meaning

Students' meanings for the line of best fit (i.e., what they conceive the line of best fit to be) fell into four categories: representing where you expect the relationship between the variables to be; shows what the data looks like; something you use to get close predictions; and average of the points. Those students who conceived of the line of best fit as 'representing where you expect the relationship between the variables to be' viewed the line as representing a general relationship between the two variables of interest, or as Jackie (pseudonym, as are all student names in the manuscript) described it, the line tells the story that was defined by the graph's title (e.g., Task 5: Heights and shoe sizes of elementary students). In contrast, students with the second conception of 'shows what the data looks like' demonstrated a localized understanding of the relationship. They viewed the line as a way to indicate where the individual points were located on that particular scatterplot. As Fran explained, "if you put the line there, and then you didn't have the dots, it would show you...it would most show you what the data looks like, and like even though it's not there...like the data, um...you can still guess at what it is, and like it's pretty easy to figure out."

Student responses in the next category, 'something you use to get close predictions,' emphasized the use of the line for making predictions. Students who used this definition said things such as 'if you used it for prediction, it would be close.' Like the students in the previous category, they viewed the line as something one can use to estimate points, but they viewed the utility of the line to be for estimating points that were not in the original data set, i.e. predictions, while students in the previous category were using the line to estimate the points in the original data set. The final category is for meanings that defined the line of best fit in terms of the average of the points, using statements like 'average of the data points' and 'average in between all of the points' to define the line. The meaning of the word 'average' for these students varied widely, from a univariate perspective to a bivariate perspective as well as corresponding to different measures of center such as the mean, median, and mode. These different meanings emerge in the next section on students' criteria and methods for placing the lines.

### 4.2 Criteria and methods

[Table 1](#) lists the criteria used by multiple students for placing the line of best fit when completing Tasks 1-5 and the number of different students who used each criterion.

**Table 1.** Criteria: Description and number of students who used each

<i>Criteria</i>	<i>Number of students</i>
Through as many points as possible	13
Equal number of points on both sides	8
As close to all the points as possible	7
Reflect the relationship the variables have, based on context knowledge	5
Halfway between the lowest and highest points	3
Through the first and last points	3
Starting from the first point then maximizing the number of points it goes through	2

First, it should be noted that a sizeable number of students (nine), when asked to find the line of best fit, wanted to bend the wire to connect the points on the scatterplot. For example, after receiving the prompt for Task 1, Marcus asked “Wouldn’t it be like the line that starts here [the origin] and like, connects...connects all these points, right?” When statements like this occurred at the start of the interview, the interviewer explained that the student was to find the line of best fit, and that since lines are straight they were not to bend the wire. After receiving this instruction, all of these students were able to move on to complete the tasks in the interview protocol.

This desire to connect the points relates to the predominant criteria students had for the line of best fit: through as many points as possible. For many students this was related to their conceptions of the line of best fit as ‘something you use to get close predictions’ or ‘shows what the data looks like.’ Hence, these students thought if you put the line through the most points you will get the most accurate prediction or view of what the data looks like because your estimate will be exactly right for those points which the line goes through. Students who used this criterion utilized the method of looking for a subset of collinear points in the scatterplot, rather than viewing the general trend of all of the points. Some students even considered those points that were not collinear to be outliers, such as Sam who said “I decided to go through the main, average points...and these [points that the line does not go through] are like outliers.” A portion of these students required their lines to start at the origin, with rationale including ‘that is where all lines start’ as well as knowledge of the context for a particular plot (e.g., In Tasks 1 and 2, dropping a ball from a height of zero centimeters should result in it bouncing zero centimeters).

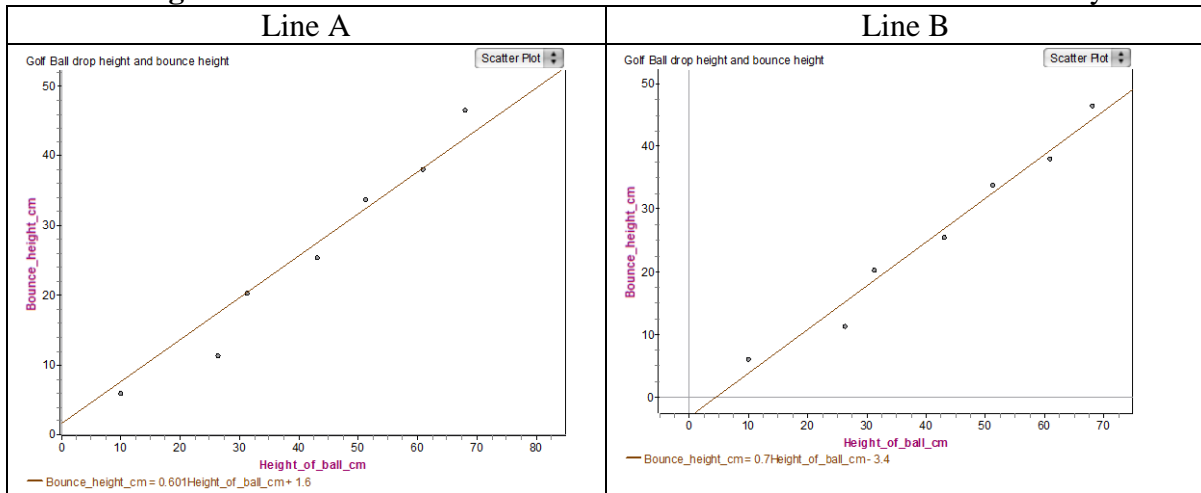
Getting an equal number of points on both sides of the line was the second most common criteria. This criterion was common for those students who viewed the line as representing where you expect the relationship between the variables to be as well as those students who thought of the line as an average. For these students, they connected the meaning of a best fit line as an average to a previous average they had used in univariate settings—the median—which has the property that there are an equal number of data points on either side of it. These students could use that same idea in this bivariate setting to find the average of the points by making a line with an equal number of points on either side of it. Interestingly, three students always placed their lines horizontally or vertically so that there were an equal number of points on either side of the

line. These students decided to focus their attention on only one of the variables of interest in the scatterplot, and placed the line of best fit to represent the middle of that univariate data set. These students were unable to coordinate the relationship between the two variables in the bivariate setting of covariation, and thus reduced it to a univariate setting.

Seven of the thirty-three students used the criteria encouraged by the authors of the CCSS-M (CCSSI 2010) and in agreement with the approach of the ordinary least-squares regression line: as close to all of the points as possible. Students with this approach justified their methods by saying things like the line was placed to be ‘relatively close to all of the points’ or ‘in between so they [the points] are all close.’ Three students who used methods to place the line through the most points on earlier tasks transitioned to use this criterion in the latter half of the interview. Students using this criterion included those whose meanings for the line were ‘average of the points,’ ‘show what the data looks like,’ and ‘something you use to get close predictions.’ These students saw the placement of the line so it is closest to all of the points as a way of actualizing these meanings.

Task 6 (see Figure 1) of the interview was designed to see whether students considered it more important for the line to go through some of the points (Line A by Angelo) or to be closest to all of the points, but not necessarily go through any of them (Line B by Barbara).

**Figure 1.** Student interview Task 6: Which line fits the data better and why?



One-third (11) of the students preferred Line A while the other two-thirds (22) preferred Line B. Seven of the 11 students who preferred Line A explicitly stated they preferred it because it went through some of the points, including three students whose dominant criterion for placing lines was through the most points. The other four students who preferred Line A did so because of concern regarding intercepts. For example, Travis explained that Line B’s x-intercept meant that when you dropped the golf ball it was not going to bounce at all, and that wouldn’t happen in reality: if you drop a golf ball, it will always bounce a little bit. Therefore, he chose Line A because it did not have that issue. Concerns with the y-intercept included that Line A starts on the y-axis, “which is what you have to do when graphing a line” (Lacey), and closer proximity of the line to the origin (Marcus, Travis, and Thurston). Arguments regarding the y-intercept also

were given by two students who chose Line B. Sam and Abby noted that Line A shows a positive drop height when you haven't dropped the ball, which isn't physically possible, and therefore they chose Line B.

A deeper analysis of the remaining 19 students who stated Line B fit the data better shows that all of these students noted that Line B was closer to all of the points than Line A and that key distinction was the impetus for their choosing of Line B. Disaggregating the students by dominant criterion, all of the students with 'closest to all the points' (with the exception of Travis and Thurston who focused on the intercepts when doing this task) and 'equal number of points on either side of the line' as their dominant criterion chose Line B, which is consistent with those criteria. Surprisingly, however, seven of the ten students whose dominant criterion was 'through the most points' chose Line B as the better line, changing to note that being closest to all of the points is what is most important for the line of best fit. As previously mentioned, for three of these students, their progression through the tasks involved a transition away from the criteria of 'through the most points' that they had utilized for the earlier tasks. As Sasha described, she "started out thinking like Angelo but now sees that Barbara's is better." For others, completing this task was an illuminating experience. To illustrate, here is an excerpt from Tom's interview:

Tom: I think...I would say Angelo's. The reason why I would say Angelo's is because two dots fit in here.

*Interviewer: Okay.*

Tom: But, this one [Barbara] looks like it has...this one looks like it goes in between them, so it kind of looks like you would...it would be more accurate for each time you do this. Like this one [Line A] is like doing the exact like...making a line using this data...like what I did.

*Interviewer: Sure.*

Tom: But, I saw her [Barbara's] idea and I started to think that hers looks more like it would be used because it's not showing...it's not going through any of the dots.

*Interviewer: Right.*

Tom: Maybe because she was thinking that this was like...if we kept on doing this, this is what it would be like. Like it wouldn't exactly be this, or that, or that [points to three consecutive dots in the scatterplot]...it would line up a little bit close to this line each time. So, I think Barbara's situation.

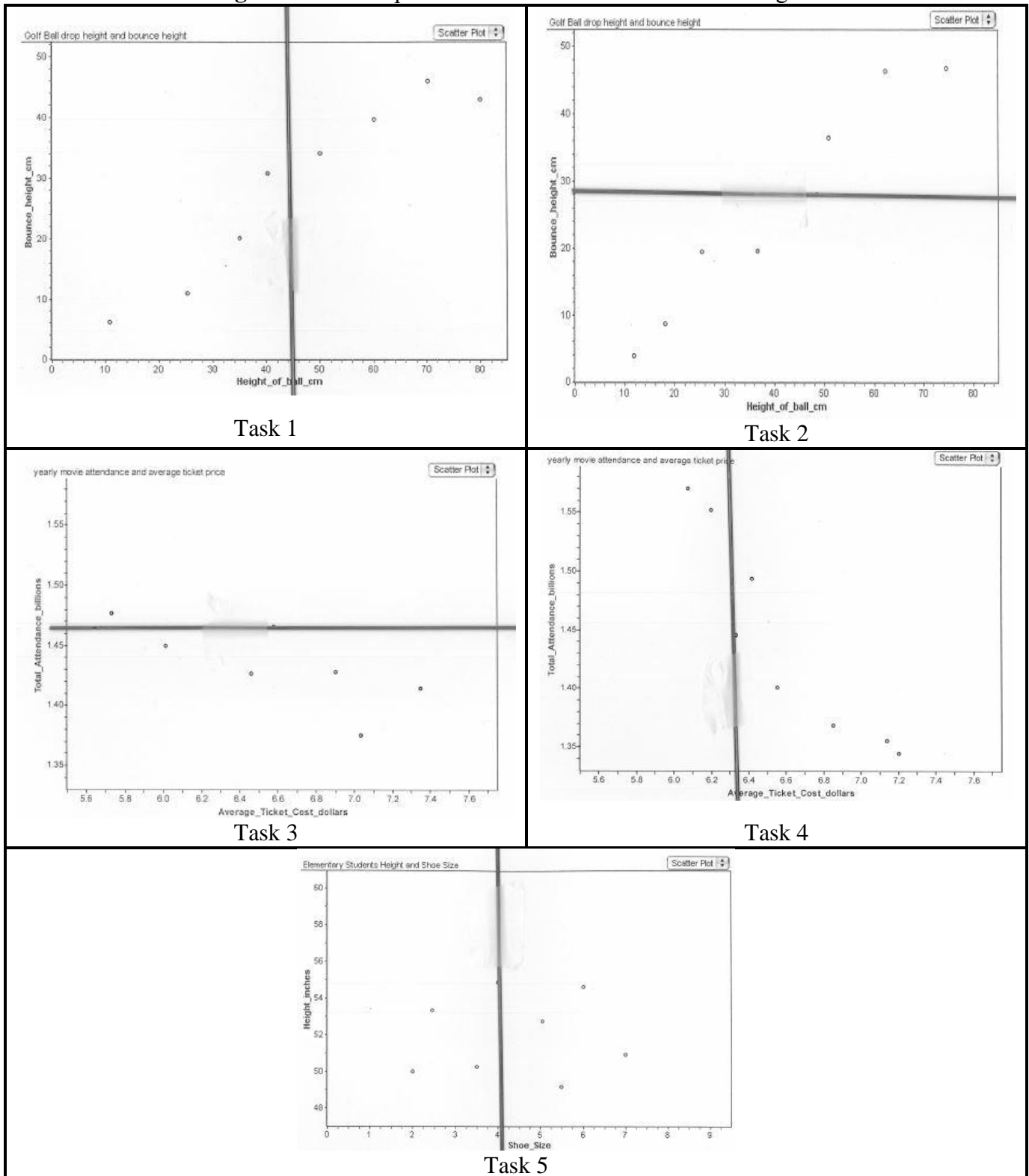
*Interviewer: Now you switched to Barbara?*

Tom: Yeah.

There were some students who used unusual, and many times unique, criteria that resulted in

very inaccurate lines. One interesting case was Sheila. [Figure 2](#) displays Sheila's lines of best fit on Tasks 1 through 5.

**Figure 2.** Sheila's placed lines of best fit for Tasks 1 through 5

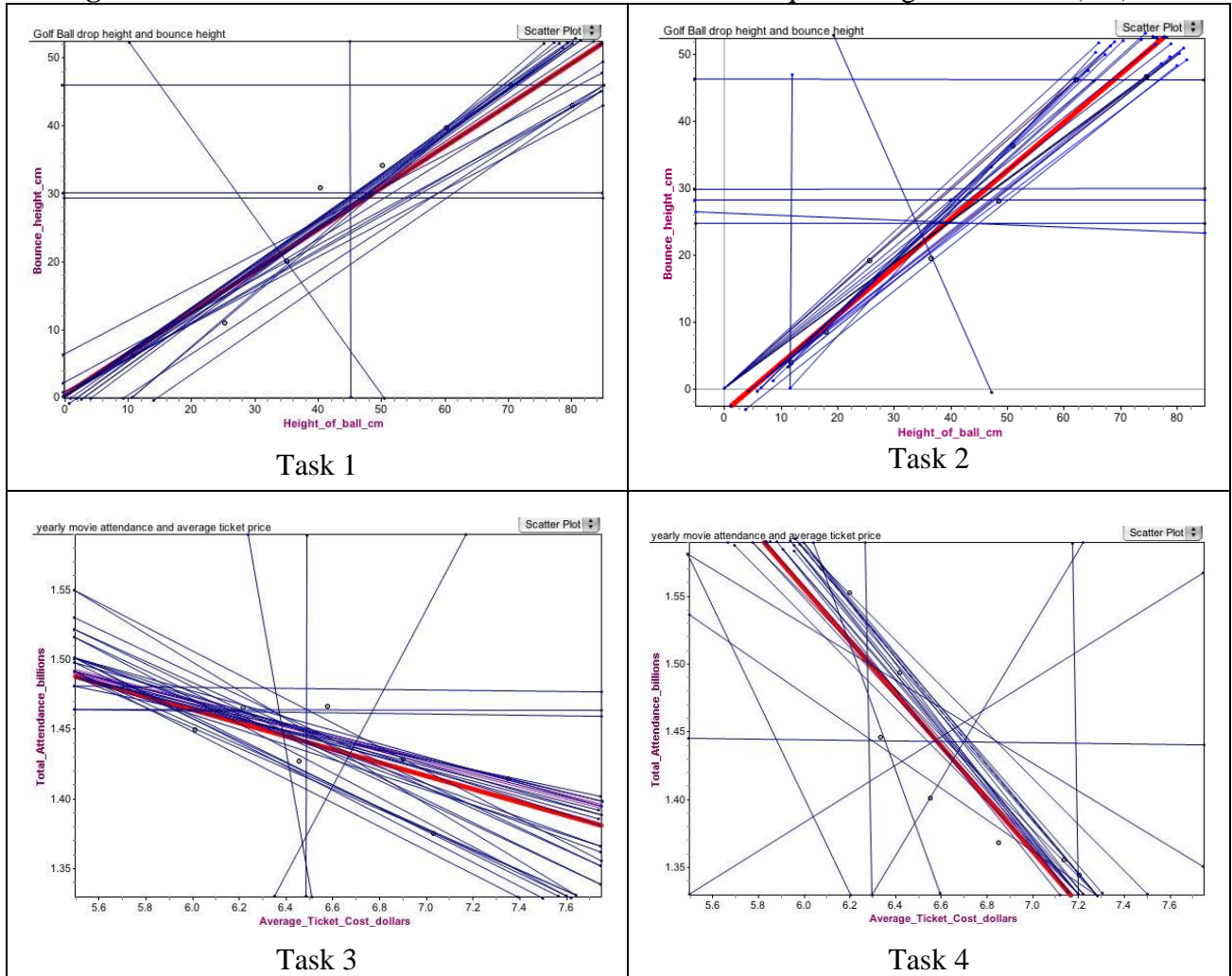


For the first task, she placed the line of best fit vertically between the 4<sup>th</sup> and 5<sup>th</sup> points in the plot because those points are in the middle of the 8 total points and it's "actually kind of like the *median*." For the second task, she placed the line of best fit horizontally at a place she deemed in the middle of the  $y$ -axis that would also enable her to place the line through a point. She thought this line was more accurate because it goes through one of the points. Her line on Task 3 was again placed horizontally, but this time she looked for points that had the same  $y$ -value because they "occur most often, kind of like *mode*; if want to describe data going to want to put what there are more of." She found two points with approximately equal  $y$ -values, and placed her line horizontally to go through those points. For Tasks 4 and 5, she went back to placing the line vertically so that there were an equal number of points on either side of the line. In explaining why she did so, she said it's "in middle of everything, like the median." She analyzed the two presented lines in Task 6 from this median and mode perspective as well, stating that Angelo's Line A used a mode approach, Barbara's Line B used a median approach, and Line A was better. She reiterated this same type of reasoning when asked to tell another student how to find the line of best fit. Sheila said her first preference was to have the line be like the mode, going through two points with the same  $y$ -coordinate. If that wasn't possible, her second preference was to place the line so it's like the median with equal number of points on either side. Sheila is a compelling case, as it seems her previous knowledge of statistics, namely median and mode, interfered with her ability to conceive of the line of best fit in standard ways and may have kept her from viewing the data from a bivariate perspective. The notion of prior knowledge interfering with one's ability to engage in the task of fitting an informal line of best fit will be revisited in the discussion section.

### 4.3 Placement

The fourth research question queried how accurately students place an informal line of best fit relative to the least-squares regression line. [Figure 3](#) presents all 33 students' lines for Tasks 1 through 4 with the least-squares regression line highlighted in red, providing a visual image regarding the accuracy of the placed lines for these tasks.

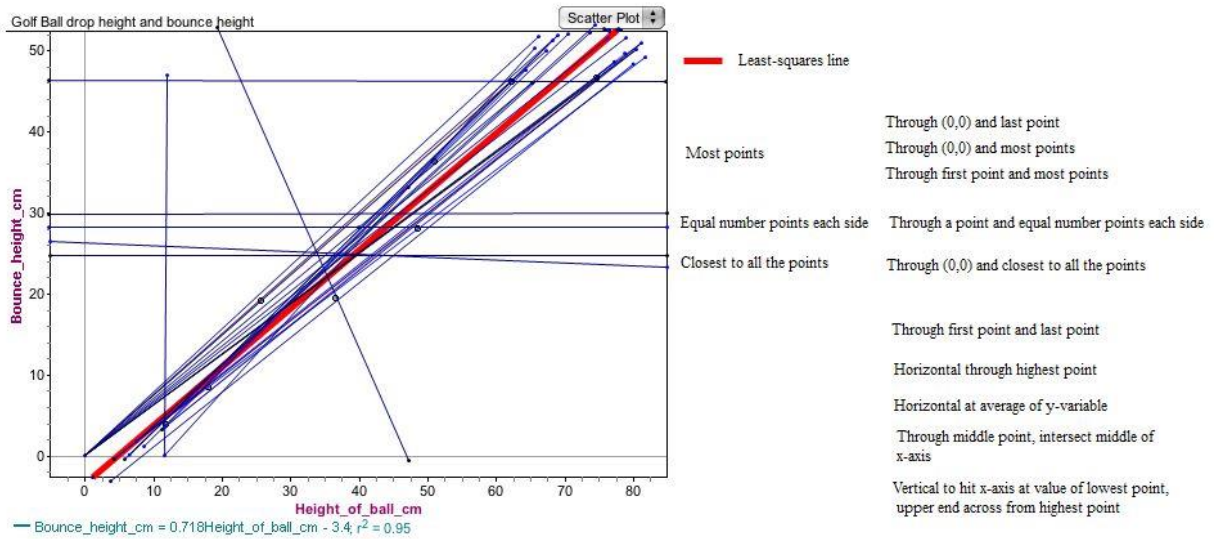
**Figure 3.** Placement of lines for Tasks 1-4 with the Least-Squares Regression Line (red)



These displays show that there is considerable variability in the placed lines' locations. The majority of the lines are reasonably accurate in that they are close to the least-squares regression line, but there are also a considerable number of lines that are placed inaccurately. A closer examination of the lines placed on Task 2 will provide insight regarding the placement of lines of best fit by the students.

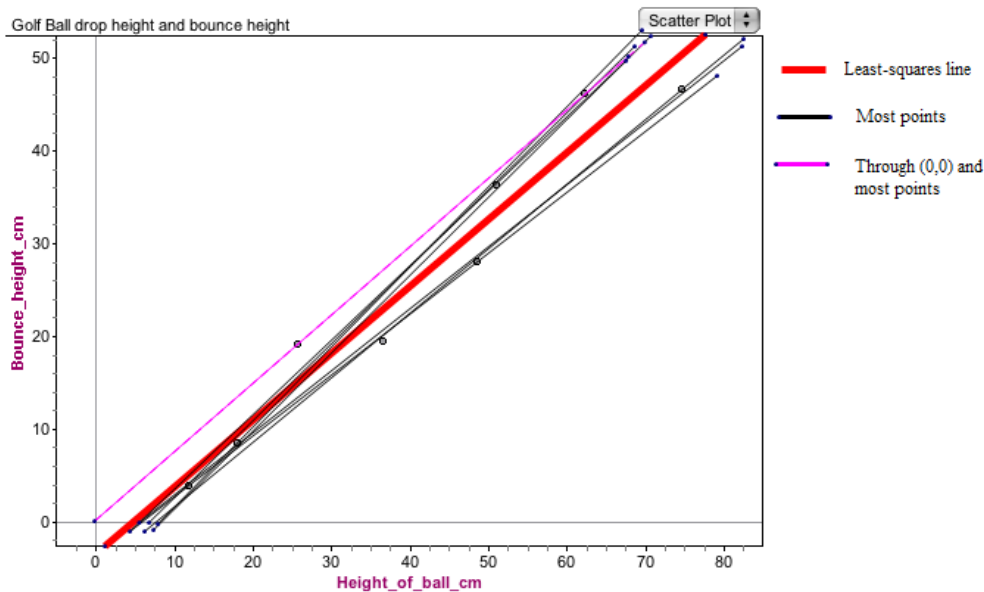
[Figure 4](#) shows the lines of best fit for all student participants, the least-squares regression line, and a list of all of the criteria that were utilized by students when placing the line on Task 2.

**Figure 4.** Task 2 placed lines with the Least-Squares Regression Line (red) and criteria



The 13 criteria utilized by the 33 students when placing the line on this task resulted in a large number of lines placed near the least-squares regression line. However, it is noticeable that these lines generally run parallel to or split the least-squares regression line, with very few following it. This is largely related to the predominance of the most points and equal number criteria. [Figure 5](#) displays the lines placed by those students who used the criterion of making the line go through the most points (including one student who added the additional restriction that the line go through (0,0)).

**Figure 5.** Task 2 lines placed to go through the most points

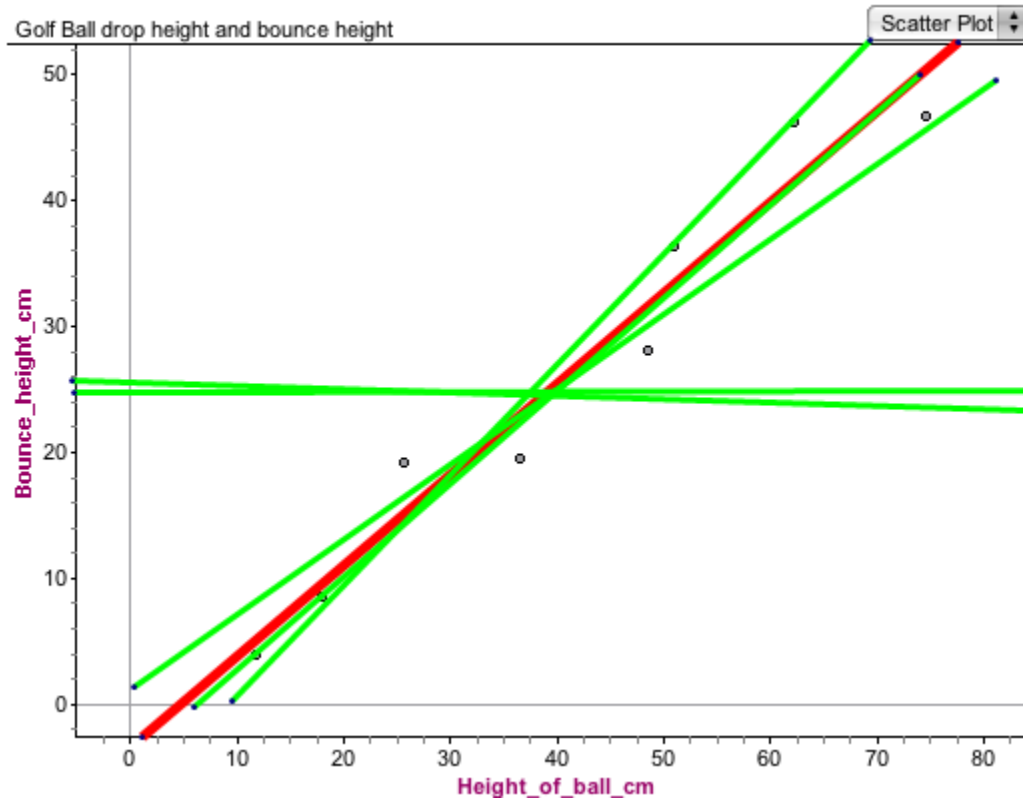


This figure shows that the lines placed by students with the most points criterion split the least-squares regression line, largely because these students forced their lines to go through one of the



last two points. [Figure 6](#) shows, in green, the lines other students placed so that an equal number of points would be on each side of the line.

**Figure 6.** Task 2 lines placed to have equal number of points on each side



Two of these lines split the least-squares regression line, placed in locations almost identical to lines shown in [Figure 5](#) which were placed using the most dots criterion. Another line follows the least-squares regression line nearly exactly. However, the other two lines placed by students utilizing the equal numbers criterion were inaccurate as they were placed horizontally. What was interesting about these students was their explanations regarding why they placed the line in that location sounded appropriate (“I’m putting it in the middle,” “It’s at the average”) and a teacher would be inclined to think these students understood best fit line; however, their lines were inaccurately placed. These students’ application of average or middle was actually done in a univariate sense which is why their lines were horizontal, not in agreement with a best fit line following the trend of bivariate data.

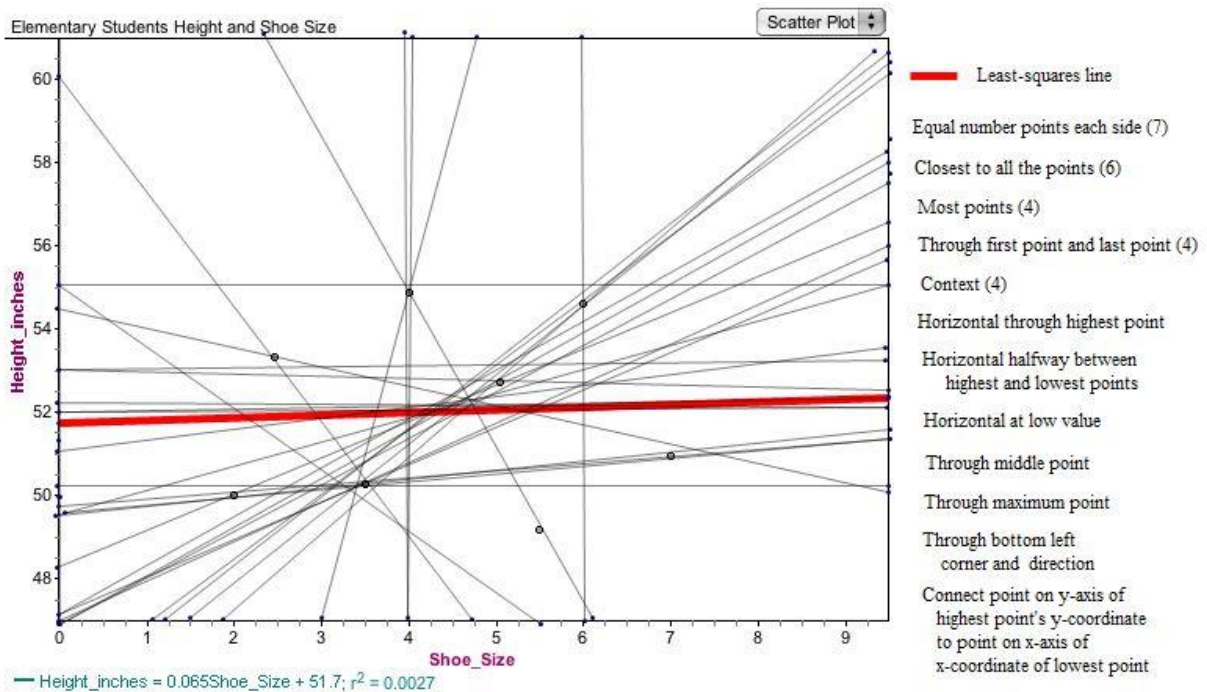
Lastly, it is noted that 8 of the 33 students forced their line to go through the origin on this task. Two of the students gave the explanation that the origin is where graphs start – ‘all graphs start there’-as their reason for this requirement. Consideration of the context, with statements like ‘if you didn’t bounce the ball there wouldn’t be a height’ and ‘if you don’t drop it then no bounce,’ were provided by three of the students.

### 4.4 Task 5: No Association

The presentation of Task 5’s scatterplot that displayed no association evoked responses and approaches from the students that were different than the previous four tasks which presented scatterplots with associations. The time it took students to complete this task was considerably longer than the other tasks, and many students studied the plot in silence for a substantial time (around twenty seconds) before responding. Six of the students initially commented that they did not see a general trend or direction in the plot and were confused about what to do. Jackie captured the sentiment of these students when she said “they [the points] were so spread out. And there wasn’t like a real thing that clicked in my mind for it to make sense. Like there’s dots everywhere, there’s no particular place or motion that it’s going.” There was one student, however, who commented that she did not see a general trend in the plot but (correctly) used that observation to decide to put the line horizontally halfway between the lowest and highest points because “it’s not decreasing or increasing” (Lacey). Lacey was the only student who was able to conceptualize that a horizontal line of best fit is how one indicates that there is no association between the variables.

There was a wider variety of locations for the placed lines on this task than the previous four tasks. [Figure 7](#) displays all of the lines placed by the students for this task (Sasha said “I have no idea” and did not place a line), along with the least-squares regression line. A list of the criteria utilized by students on this task ordered by frequency of use is provided on the right side of the figure, with the number of students utilizing a criterion shown in parentheses if it was used by multiple students.

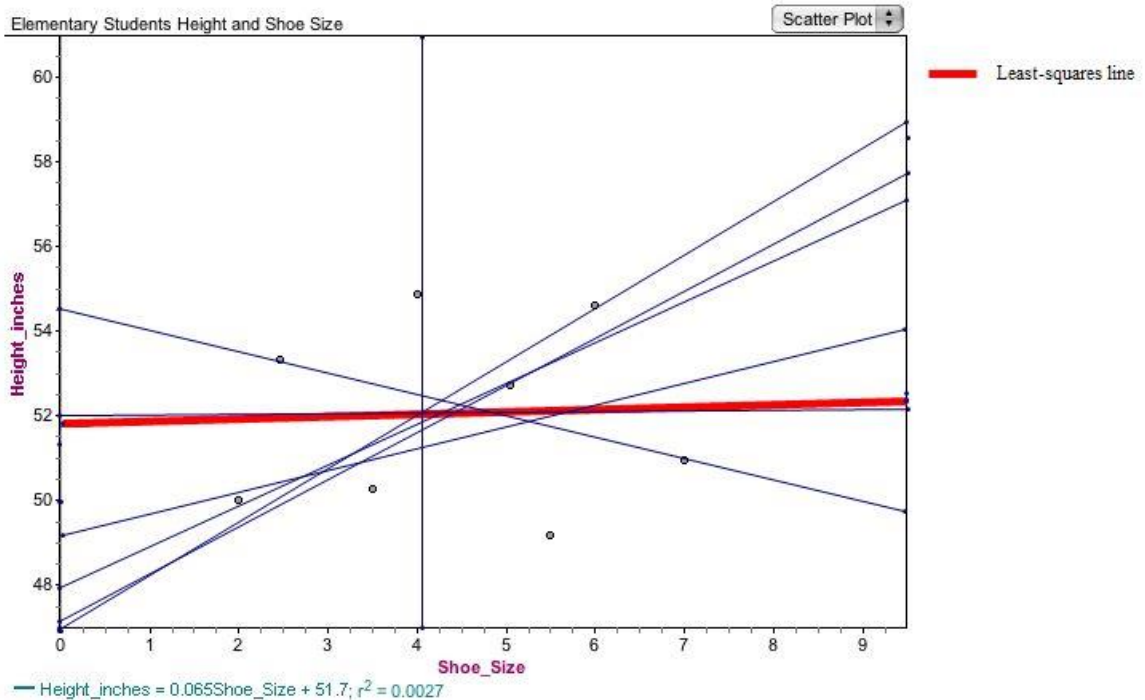
**Figure 7.** Task 5 placed lines with the Least-Squares Regression Line (red) and criteria



Notably, there are a large number of placed lines with positive slopes on this task. Those were produced by students who utilized criteria of ‘through first point and last point,’ ‘equal number points each side,’ ‘most points,’ ‘through bottom left corner and direction,’ and ‘context.’ In addition to those students who used their knowledge of the context exclusively when placing their lines, many of the students who utilized other criterion explained that the context was a secondary consideration when choosing where to place the line. Abby, who utilized the ‘equal number points each side’ criterion, explained that “it makes sense that the taller they are the bigger the shoe size.” Therefore, she placed her line with a strong positive slope and an equal number of points on each side. Looking at the placed lines that are horizontal, there were a few lines that are very close to the least-squares regression line. One was placed by Lacey, who, as already described, was the only student to purposefully place a line horizontally to represent no association between the variables, while the others were placed by students utilizing the ‘equal number of points each side’ and ‘closest’ criterion. Other students placed their lines horizontally but inaccurately. Marta placed the line horizontally through the highest point, as she had done previously on Tasks 1 and 3, to show a key point in the data. Jackie placed her line low on the plot (around height = 50.5 inches) because this plot was of elementary school students who are younger and smaller, so she “was trying to go lower...to make...so it would be more of like an average number” for these younger individuals.

For this task, the dominant criterion were ‘equal number of points each side’ (7 students), ‘closest to all the points’ (6 students), through the ‘most points’ (4 students), and knowledge of the ‘context’ (4 students). Looking closer at the lines placed by the seven students who used the ‘equal number of points each side’ criterion (Figure 8) illustrates that a weakness of that criterion is that it can result in a multitude of different lines.

**Figure 8.** Task 5 lines placed to have equal number of points on each side with the Least-Squares Regression Line (red)



Four of these students exhibited confirmation bias ([Nickerson 1998](#)). Their expectation that students with bigger shoe sizes have greater heights influenced the placement of their lines; all of their lines had strong positive slopes. The student who placed the line with a negative slope had a secondary consideration of going through points, placing her line to have an equal number on each side as well as to go through the 2<sup>nd</sup> and 8<sup>th</sup> points. The vertical line was placed by Sheila, who as previously mentioned viewed these tasks from a univariate perspective tied to median and mode, so in this task she placed the line to be like the median for the shoe sizes. Lastly, one student was able to place a relatively accurate, horizontal line using this criterion.

## 4.5 Interpretation

Three-fourths of the students (25 out of 33) were able to broadly interpret their best fit lines from Tasks 1 and 3 to construct correct statements regarding what the lines show about the relationship between the variables for those tasks. A statement like this one, given by Travis, typifies these students' responses to interpreting their lines from Task 1: "Shows that the higher the ball is dropped from the higher it bounces." Likewise, this quote from Talia is representative of these students' interpretations of their lines of best fit for Task 3: "As the price went up you could see that the amount of people went down." Examining the responses of the remaining eight students, three did not respond appropriately to the prompts despite numerous redirects from the interviewer, focusing instead on how they placed their line or explaining the physics behind ball bounces, and another student was unable to respond to the prompt saying she did not know what these lines showed about the relationship. Marta, whose method for placing the lines was to go through a key point such as the maximum, struggled to use her placed lines to say anything about the relationship and came to the conclusion that the line "doesn't show anything" about the relationship between the displayed variables. What stood out, however, were the remaining three students responses' that the lines do not show anything about the relationship between the variables because they do not match the variability exhibited in the data. For example, Evan's response was "Well I don't know 'cause it's [the line she placed] not consistent [with the data points]. Like it doesn't go up evenly [gestures to the points in the scatterplot] so like I wouldn't be able to tell you." These students do not yet understand the line of best fit as a model for summarizing data that exhibit variability.

Four of the students took the additional leap to interpret their line as evidence of a causal relationship between the variables. Jamal explained his line on task three as showing that "Less people are going [to the movies] but the prices are increasing since they have to make money so the price goes up", so a decrease in movie attendance is causing the film industry to increase the average price of a movie ticket. Similar responses by the other three students confirm [Batanero et al.'s \(1997\)](#) finding that persons often believe association implies causation.

## 5. Discussion

### 5.1 Discussion of the results

The purpose of this study was to learn what conceptions students have about the line of best fit just prior to formal instruction on the topic. Their conceptions regarding what they conceive the

line of best fit to be fell into four categories: representing where you expect the relationship between the variables to be; shows what the data looks like; something you use to get close predictions; and average of the points. These were related to the criteria and methods students utilized to place a line of best fit as described previously. For example, students whose meaning for the line of best fit was ‘something you use to get close predictions’ or ‘shows what the data looks like’ constituted most of the students whose criterion for the line of best fit was to maximize the number of points that it went through. Reflecting upon the criteria students used to place their lines of best fit, it is clear that they fall along the continuum of case-oriented to aggregate views of data ([Bakker 2004](#)). Students whose criterion for the line of best fit was halfway between the highest and lowest points, through the first and last point, or through a middle point had case-oriented views in approaching the task of fitting a best fit line. They focused their attention on specific points they deemed significant (i.e., highest, lowest, first, last, or middle point) and considered those points exclusively when placing the line, ignoring the rest of the points in the data set. Other students used approaches that were in the middle of the continuum. Examples include going through the most points and going through points with approximately the same  $y$ -value. To apply these criteria, the students studied the entire data set until they found a set of points that enabled them to meet the criteria, but they were really fitting the line to a subset of points rather than the data set as a whole and their insistence on going through points showed they were still concerned with the line displaying characteristics of the data set at the individual case level. Therefore, these students were considered to be in transition from a case-oriented view to an aggregate view of data. An aggregate view of data was identified in students whose criteria required analysis of the data as a whole entity with characteristics not necessarily visible in any of the individual cases. Students who placed the line to have an equal number of points on each side as well as those who wanted the line as close to all of the points as possible met this description and thus were classified as having an aggregate view of the data.

Comparing the results of the present study with grade eight students to that of [Sorto et al.’s \(2011\)](#) with university students, both groups of students included members whose conception for the line of best fit was at the middle or average of all the points, and some students accomplished that by having an equal number of points on either side of the line. While seven of the eighth-grade students felt it important that the line be placed closest to all of the points, the university students were more sophisticated in their responses and better able to explain what that meant in terms of the distance of the points to the line. This is not unexpected due to the difference in ages and mathematical sophistication of the students.

It was evident from this study that students’ prior study of mathematics and statistics impacted their conceptualizations of the line of best fit, in both productive and unproductive ways. Looking at the influence of their prior learning in mathematics, all of the students were able to ‘read’ the scatterplots in that they knew and understood what each point represented, largely due to their previous work in mathematics with the coordinate grid and plotting of points. They were also generally successful at interpreting their lines of best fit as showing increasing or decreasing relationships depending on which way their line was positioned (going up from left to right for increasing and going down from left to right for decreasing), which can be attributed to their previous work with lines with positive and negative slopes in mathematics. However, there were some ways in which their prior learning of mathematics interfered, including their ability to make sense of the task of fitting a best fit line. Many students struggled to conceive of the line of

best fit as a line that does not necessarily go through all of the points, likely because this differs from graphs of linear functions that these students have studied in mathematics that necessarily go through every plotted point. Their previous study of lines in mathematics was also the source of their undue concern with the intercept(s) of the line. As noted, many students required that their line go through the origin—some for reasons that tied to the context (if you don't drop a ball, it will not have a bounce height), but others because they had a preconceived notion that all lines start there. Concern with the intercepts of the line also arose in Task 6 where students were asked which line was the better fit for the data set. Six of the students based their decision exclusively on the intercept(s) of the lines rather than examining how well the line fit the data. This concern with the intercepts of a line originates from students' work in mathematics where the  $y$ -intercept is traditionally used as the starting point for graphing a line and questions about the  $x$ - and  $y$ -intercepts are commonly posed to students. In the statistical setting of the line of best fit, however, the intercepts are generally not of interest; the primary and sole criteria is getting the line as close to all of the points as possible. The  $y$ -intercept is not necessarily to be used as the starting point for determining the placement of the line, and it is entirely possible that the  $y$ -intercept of the line of best fit will have no meaning in the context of the problem. This is acceptable to statisticians because that is not a goal of the line of best fit, particularly for data sets which do not have values of the explanatory variable near zero. However, as seen in the results of this study, these differences between mathematical and statistical approaches to lines are harder for students to accept and can create obstacles to their learning of the topic of line of best fit.

Students' prior learning of statistics to analyze univariate data sets also was a source of students' conceptions about the line of best fit. Sixteen of the students referenced the term 'average' when describing the line of best fit during their interviews, and many of them were able to leverage their meaning of that term to help them understand the line of best fit correctly as a model depicting the relationship between the variables and place accurate lines. For others, the notion of fitting a line to data evoked conceptions of the line as an average in the sense of the median, so that students placed their line to have an equal number of points above and below the line. As demonstrated, sometimes this method can result in an accurately placed line but it can also result in a line which is completely inaccurate. For some students who placed inaccurate lines via an 'equal number of points on each side' approach, they were unable to make the transition to understanding average in a bivariate sense, continuing to think of average from a univariate perspective and placing horizontal and/or vertical lines that would be at the average of one of the variables. One student, Sheila, referenced average as mode in addition to median as she completed the interview tasks, placing inaccurate lines for all tasks and choosing Line A as the better line. Still, others used the line to highlight an average point in the plot, such as forcing the line to go through a middle point or a point with an average  $y$ -value. It is clear from these examples that it is common and often natural for students to evoke the notion of 'average' when thinking about the line of best fit. However, many students need to learn how to utilize the meaning of average in productive ways to accurately place a line of best fit and understand its meaning.

## 5.2 Implications for the field

These findings provide teachers and teacher educators with descriptions of students' initial conceptions regarding the line of best fit, which can be utilized to develop teachers' knowledge of content and students ([Ball et al. 2008](#)) for teaching the line of best fit. It is hoped that through attaining this essential knowledge for teaching, teachers can better design effective instruction that begins with the conceptions students have and work to move students from those initial conceptions to more advanced knowledge ([Bransford et al. 1999](#)) about the line of best fit. For example, by knowing that it is likely that a sizeable number of students in a class will think going through points is more important than being closest to all points, a teacher can design instructional activities to address this idea and transition students to understand why statisticians consider it most important for the line to be closest (measured vertically) to all points. Teachers can also use this knowledge to anticipate what students typically think about and do when considering the meaning and placement of a line of best fit. This will better prepare them to orchestrate productive discussions in their instruction on the topic ([Smith and Stein 2011](#)). It should be noted that these findings are not only applicable to mathematics teachers and teacher educators, as has been emphasized, but also to science and social studies teachers and teacher educators as the topic of line of best fit is included in those curricula as well ([NCSS 2010](#); [NGSS Lead States 2013](#)).

The findings from this study also have implications for school curriculum writers, such as textbook authors. Three specific examples will be described. First, due to students' undue concern with starting lines of best fit at the origin, it is recommended that students not begin to work with lines of best fit in settings where placing the line through the origin is reasonable (either in that setting or due to the pattern of points in the scatterplot). This way, the students will be able to focus their attention on the pattern of points in the scatterplot and will not be distracted by concerning themselves with the  $y$ -intercept of the line being the origin. Another recommendation is that students be presented with a task similar to Task 6 from the interview, where they are forced to choose which of two lines is the better fit and one of the lines goes through points while the other is closest to all of the points. Engaging in this task was a productive learning experience for some students and could be an effective way to transition students away from the idea that going through points is the criterion of the line of best fit. Lastly, there is value to engaging students in the task of fitting a line of best fit to data that has no association (Task 5 in this study). It was apparent from this study's findings that students do not naturally conceive of a horizontal best fit line as the appropriate line of best fit for data sets with no association, as only one of the 33 students in this study was able to determine this. However, the horizontal placement of the line of best fit in these situations is an important concept, as it is the basis of future statistical tests of association between quantitative variables which determine whether the slope of the best fit line is statistically significantly different from zero in order to conclude that there is an association between the variables. Inclusion of such a task at this point in students' education is important, as it lays the foundation for students' learning of related, advanced statistical approaches in the future and is linked to the determination of the location of a line of best fit.

### 5.3 Future directions

This study utilized a convenience sample of 33 students from the United States. Thus, the sample has some limitations in that the students were not randomly selected from the population of all students, so the results are not necessarily generalizable to all students. However, given that every student who agreed to participate in the study was included and there was diversity in these students' mathematics backgrounds and geographic locations, it seems likely that this study's results are indicative of a broader population of United States students' conceptions. In addition, the protocol had some inherent limitations in that it was limited to five tasks where students placed lines and thus is limited in scope. Given the difficulties students had with Task 5, it would have been interesting to explore other task situations with non-associated data where the context and variables might be considered less related (e.g., height and number of family members) to reduce the influence of confirmation bias.

There are many directions that future work in this area could take. Regarding further research, similar studies could be done in other countries to describe student conceptions of line of best fit from a more global perspective. It would also be helpful to the field for an expert level study to be added to the current student level study and [Casey and Wasserman's \(2015\)](#) teacher level study. Informally fitting a line of best fit, by eye, in many ways is more of an art than a science. It would be informative to know how experts approach the task of informally fitting a best fit line and use the findings to improve the teaching of this topic to students. Another direction future research could take is to determine how and in what sequence the learning of line of best fit should be integrated with the learning of lines in mathematics. As seen in the present study, students' previous knowledge of lines learned in mathematical settings often interfered with their ability to accurately conceptualize the line of best fit. However, the study of the line of best fit is often linked to the learning of lines in the curriculum; in fact, it was placed in the same grade as the learning of mathematical lines in the United States ([CCSSI 2010](#)). One avenue of inquiry researchers could address is: What is the optimal way to teach line of best fit and mathematical lines? Is it better to teach lines of best fit first, then mathematical lines? Or is the reverse order better? How or could the learning of these topics be synthesized to maximize student understanding? As mathematical lines and lines of best fit are common to school curricula worldwide, the results of such work would have broad applicability and impact.

Another direction future work could take concerns teacher education. Development of curricula that incorporates the findings of this study and other relevant studies for use with teachers to develop their knowledge for teaching line of best fit is a needed endeavor as such curriculum does not currently exist. In particular, curriculum that helps teachers understand differences between mathematics and statistics as disciplines as well as in specific instances like mathematical lines compared to statistical lines of best fit is sorely needed, for the responsibility of teaching statistics is placed upon the shoulders of mathematics teachers in school (grades K-12) settings.

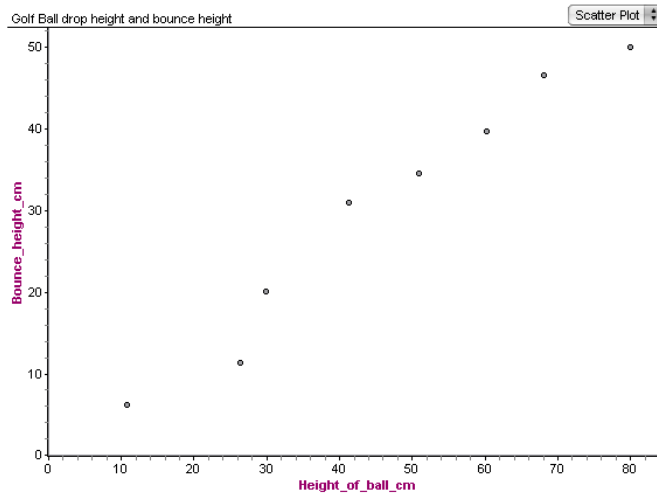
---



## Appendix A: Student interview protocol

**Introduction:** Today we're going to look at a few graphs and I'm going to ask you some questions about them. The graphs are called scatterplots. Let's take a look at one and see if it looks familiar. [Place task 1 in front of student].

**Task 1.** A golf ball is dropped from different heights and the bounce height is measured.



Here is the first scatterplot which shows data about how high a ball bounces when it is dropped from different heights. Students were given set heights to drop the ball from. Then they dropped a golf ball from each of those heights and measured how high the ball bounced back up.

For example, this point [point to the lower left point] has coordinates (10.8, 6.1). So that means the golf ball was dropped from 10.8 centimeters above the ground [move finger or pen from the point to the x-axis], and bounced back up to 6.1 centimeters above the ground [move finger or pen from the point to the y-axis]. Do you understand?

How many times was the ball dropped for the data on this graph? [Correct answer: 8; if student gives incorrect answer, explain how each data point represents another drop of the ball from a different height. As there are 8 points total, that means the ball was dropped 8 times.]

[If the student hasn't already said] Have you worked with scatterplots before?

[If student replies no], No problem. I am going to explain each graph to you so please ask me about anything you are unsure of.

For each graph, we will talk about the data displayed and then I will ask you to determine the line of best fit for the data points. I would like you to think out loud as you decide where to place the line on each graph so that I can understand how you are deciding where to place it.

Task 1: Let's use this piece of wire as the line so you can move it on the paper and place it where it best fits the data. [Place a wire on the paper in the white space, not on the graph] There are no right or wrong answers. I am interested in how you think about placing the piece of wire, so please tell me what you are thinking as you do this so I can follow your thoughts. Once you

decide where you want the wire line on the scatterplot, I will use scotch tape to keep it in place. Go ahead and place the line where you think it should be.

[If student asks you questions about what to do or if his/her answer is right, reply “It is up to you” or “It is your decision”.]

After each line of best fit is completed (if student did not mention when talking aloud), ask follow-up questions where applicable such as:

Why did you choose to put the line there? or

What criteria did you use in deciding where to put the line?

Did your criteria change for this graph? (for tasks after the first).

How did you decide on that criteria? Or

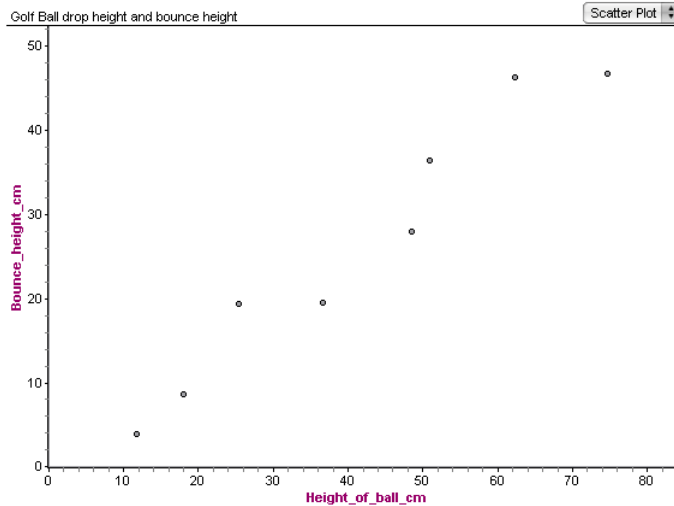
Why do you think that was important to consider?

Could you say more about how you decided at the end that it should be angled a bit more?

For Tasks #1 and #2 (if applicable): I noticed you moved the line so that it went through here (point to the intersection of the axes). Can you talk about why you thought that was important?

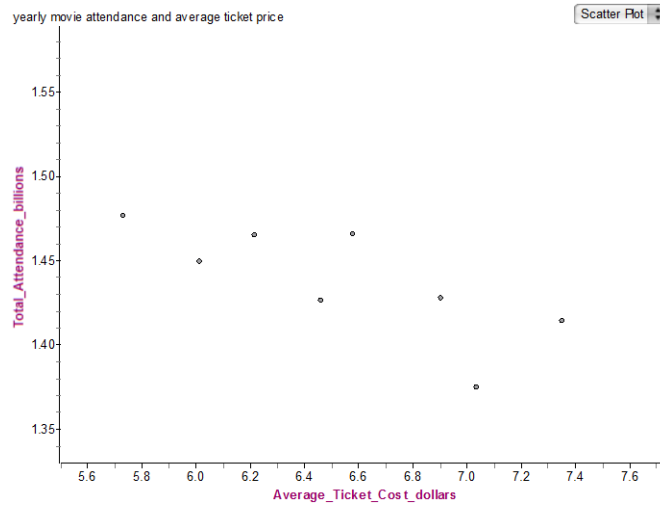
Task 2: The next scatterplot displays drop height and bounce height data for a different type of golf ball ... [repeat directions and repeat follow-up questions]

**Task 2.** A second golf ball is dropped from different heights and the bounce height is measured.



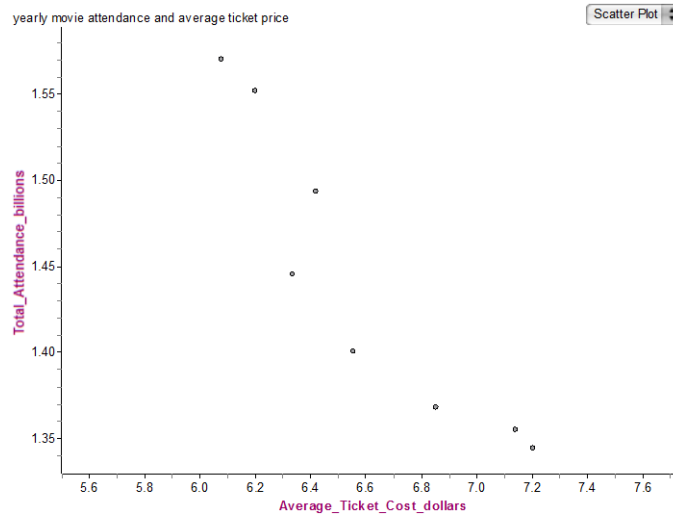
Task 3: The next two scatterplots display data from a different setting. The data on this [show Task 3 graph] graph are collected from movie theaters across the country. We can estimate the total number of movie tickets sold over an entire year and also estimate the average price of a ticket for that year based on sales in a sample of cities. For example, this point [point to second point] shows that during one year when the average price of a ticket was about \$6 the total number of people attending movies was about 1.45 billion people. Where do you think the line of best fit would be for this data?

**Task 3.** The estimated total number of people going to the movies in the United States during one year and the average price of a ticket during that year.



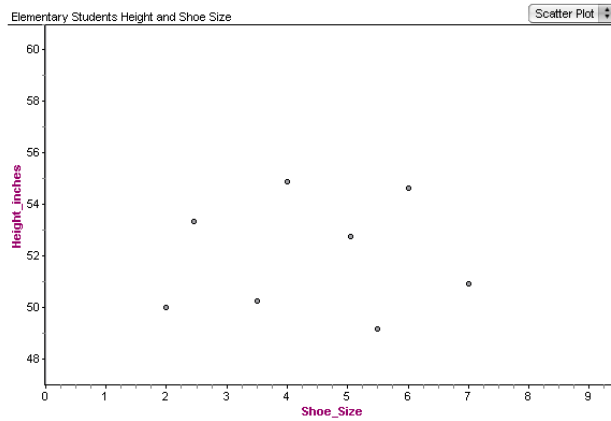
Task 4: The next graph shows similar data that is calculated using movie sales data from a different set of cities during the same 8 years.

**Task 4.** A different estimate of the total number of people going to the movies in the United States during one year and the average price of a ticket during that year.



Task 5: An elementary teacher wrote down the height in inches and shoe size for each of 8 students in his class. He plotted the data in this scatterplot.

**Task 5.** Heights and shoe sizes of elementary students.

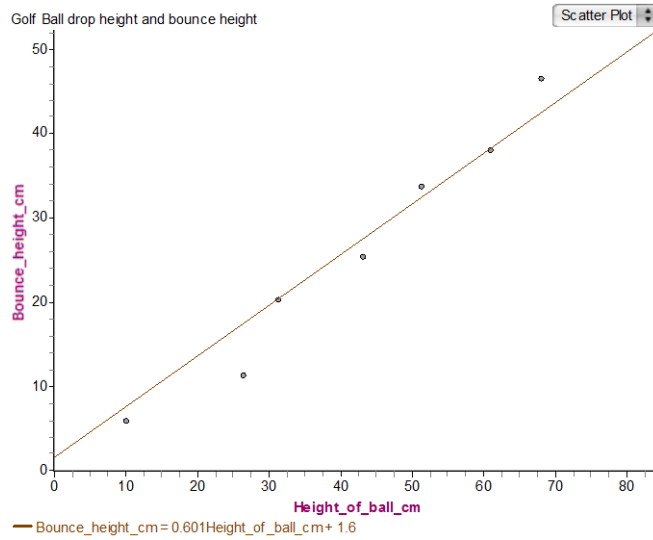


Now I have some general questions to ask you about the line of best fit.

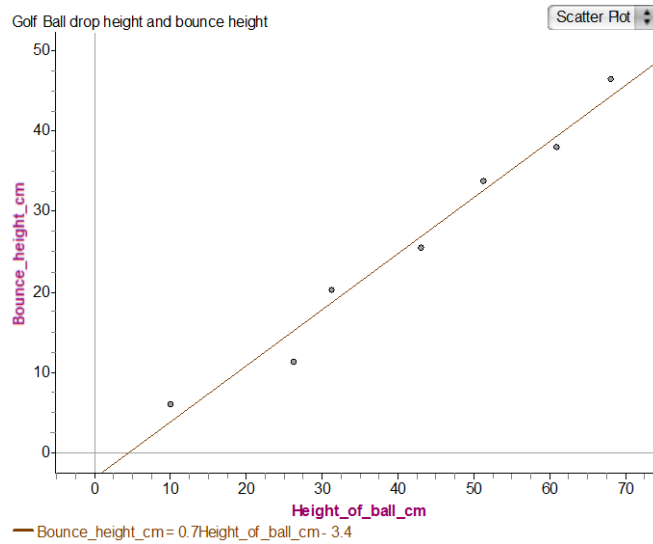
- 1) Could you tell me what you would say to another student that asked you “What is the line of best fit?”
- 2) What would you say to another student to help them draw the line of best fit on a scatterplot?
- 3) As you completed the tasks, did your thoughts about where the line of best fit gets placed change? If so, how?

Task 6: Two students, Angelo and Barbara, were given the same task you were given in task 1: to find the line of best fit for the data on a golf ball’s drop height and bounce height. [Place task 6 in front of student]. They had two different solutions. Which student’s line fits the data better and why?

**Task 6. Angelo's solution:**



**Barbara's solution:**



I have a couple more questions for you.

4) (Show student task #1 again) Can you talk about what the line shows about the relationship between the bounce height of a ball and the height it was dropped from?

(Show student task #3 again) Can you talk about what the line shows about the relationship between attendance totals for movies and the price of a movie ticket?

Thank you very much for helping me understand how you approach these kinds of tasks. It will help us teach students better if we understand how they think about problems like these. Have a good day!

## Acknowledgements

I acknowledge the collaborative work of David Wilson on this study and the remarks of Mike Shaughnessy, reviewers, and editors on earlier drafts of the manuscript.

---

## References

Australian Curriculum, Assessment, and Reporting Authority (2012), *The Australian curriculum: Mathematics*, Sydney, Australia: Author.

Ball, D. L., Thames, M. H., and Phelps, G. (2008), "Content Knowledge for Teaching: What Makes It Special?," *Journal of Teacher Education*, 59(5), 389-407.

Bakker, A. (2004), "Reasoning about Shape as a Pattern in Variability," *Statistics Education Research Journal*, 3(2), 64 – 83.

Bargagliotti, A., Anderson, C., Casey, S., Everson, M., Franklin, C., Gould, R., Groth, R., Haddock, J., and Watkins, A. (2012), "Project-SET Linear Regression Learning Trajectory," available at <http://project-set.com/presentations/121712-regressionlp-final-released/>

Batanero, C., Estepa, A., and Godino, J. (1997), "Evolution of Students' Understanding of Statistical Association in a Computer-Based Teaching Environment," in *Research on Teaching Statistics and New Technologies*, eds. J. Garfield and G. Burrill, Voorburg, The Netherlands: International Statistical Institute, pp. 191-206.

Batanero, C., Estepa, A., Godino, J., and Green, D. (1996), "Intuitive Strategies and Preconceptions about Association in Contingency Tables," *Journal for Research in Mathematics Education*, 27(2), 151-169.

Bransford, J. D., Brown, A. L., and Cocking, R. R. (eds.) (1999), *How people learn: Brain, Mind, Experience, and School*, Washington, D.C.: National Academy Press.

Burrill, G., and Biehler, R. (2011), "Fundamental Statistical Ideas in the School Curriculum and in Training Teachers," in *Teaching Statistics in School Mathematics-Challenges for Teaching and Teacher Education, A Joint ICMI/IASE Study: The 18<sup>th</sup> ICMI Study*, eds. C. Batanero, G. Burrill, and C. Reading, New York, NY: Springer, pp. 57-69.

Casey, S., & Wasserman, N. (2015), "Teachers' Knowledge about Informal Line of Best Fit," *Statistics Education Research Journal*, 14(1).

Cobb, P., McClain, K., and Gravejeijer, K. (2003), "Learning about Statistical Covariation," *Cognition and Instruction*, 21(1), 1-78.

Common Core Standards Writing Team (2011), "Progressions for the Common Core State

Standards in Mathematics (Draft),” Available at [http://commoncoretools.files.wordpress.com/2011/12/ccss\\_progression\\_sp\\_68\\_2011\\_12\\_26\\_bis.pdf](http://commoncoretools.files.wordpress.com/2011/12/ccss_progression_sp_68_2011_12_26_bis.pdf)

Common Core State Standards Initiative (CCSSI) (2010), *Common Core State Standards for Mathematics*, Washington, D.C.: Author.

Confrey, J. (1990), “A Review of the Research on Student Conceptions in Mathematics, Science, and Programming,” *Review of Research in Education*, 16, 3-55.

Cooke, G., Heideman, C., Keene, A. J., Lin, A., and Reeves, A. (2007), *Pearson Math 9*, Don Mills, Ontario: Pearson Education Canada.

Coxford, A. F., Fey, J. T., Hirsch, C. R., Schoen, H. L., Burrill, G., Hart, E. W., Watkins, A. E. with Messenger, M. J., Ritsema, B. E., Walker, R. K. (2003), *Contemporary Mathematics in Context: A Unified Approach* (Course 1), New York, NY: Glencoe McGraw-Hill.

Estepa, A., and Batanero, C. (1996), “Judgments of Correlation in Scatter Plots: An Empirical Study of Students’ Intuitive Strategies and Preconceptions,” *Hiroshima Journal of Mathematics Education*, 4, 25-41.

Franklin, C., Kader, G., Mewborn, D., Moreno, J., Peck, R., Perry, M., and Schaeffer, R. (2005), *Guidelines for Assessment and Instruction in Statistics Education (GAISE) Report: A PreK-12 Curriculum Framework*, Alexandria, VA: American Statistical Association.

Garfield, J., and Ben-Zvi, D. (2004), “Research on Statistical Literacy, Reasoning, and Thinking: Issues, Challenges, and Implications,” in *The Challenge of Developing Statistical Literacy, Reasoning, and Thinking*, eds. D. Ben-Zvi and J. Garfield, Dordrecht, The Netherlands: Kluwer Academic Publishers, pp. 397–409.

Glaser, B. G., and Strauss, A. L. (1967), *The Discovery of Grounded Theory*, Chicago, IL: Aldine.

Jennings, D., Amabile, T., and Ross, L. (1982), “Informal Covariation Assessment: Data-based Versus Theory-based Judgments,” in *Judgment under Uncertainty: Heuristics and Biases*, eds. D. Kahneman, P. Slovic, and A. Tversky, Cambridge, England: Cambridge University Press, pp. 211-230.

Konold, C., Robinson, A., Khalil, K., Pollatsek, A., Well, A., Wing, R., and Mayr, S. (2002), “Students’ Use of Modal Clumps to Summarize Data,” in *Proceedings of the Sixth International Conference On Teaching Statistics*, ed. B. Phillips [CD-ROM] Cape Town, South Africa.

Kuhn, D., Amsel, E., and O’Loughlin, M. (1988), *The Development of Scientific Thinking Skills*, Orlando, FL: Academic Press.

McGahan, J. R., McDougal, B., Williamson, J. D., and Pryor, P. L. (2000), “The Equivalence of

Contingency Structure for Intuitive Covariation Judgments about Height, Weight, and Body Fat,” *Journal of Psychology*, 134, 325-335.

Mokros, J., and Russell, S. J. (1995), “Children’s Concepts of Average and Representativeness,” *Journal for Research in Mathematics Education*, 26(1), 20-39.

Moritz, J. B. (2004), “Reasoning about Covariation,” in *The Challenge of Developing Statistical Literacy, Reasoning, and Thinking*, eds. D. Ben-Zvi and J. Garfield, Dordrecht, The Netherlands: Kluwer Academic Publishers, pp. 227-256.

National Council for the Social Studies (NCSS) (2010), *National Curriculum Standards for Social Studies: A Framework for Teaching, Learning, and Assessment*, Silver Spring, Maryland: Author.

Next Generation Science Standards (NGSS) Lead States (2013), *Next Generation Science Standards: For States, By States*, Washington, D.C.: The National Academies Press.

National Council of Teachers of Mathematics (NCTM) (1989), *Curriculum and Evaluation Standards for School Mathematics*, Reston, VA: Author.

National Council of Teachers of Mathematics (NCTM) (2000), *Principles and Standards for School Mathematics*, Reston, VA: Author.

Nickerson, R. S. (1998), “Confirmation bias: A Ubiquitous Phenomenon in Many Guises,” *Review of General Psychology*, 2(2), 175-220.

Qualifications and Curriculum Authority (2007), *The National Curriculum 2007*, Earlsdon Park, Coventry: Author.

Smith, M. S., and Stein, M. K. (2011), *5 Practices for Orchestrating Productive Mathematics Discussions*, Reston, VA: NCTM.

Sorto, M. A., White, A., and Lesser, L. (2011), “Understanding Student Attempts to Find a Line of Fit,” *Teaching Statistics*, 11(2), 49-52.

---

Stephanie A. Casey  
Eastern Michigan University  
515 Pray-Harrold  
Eastern Michigan University  
Ypsilanti, MI 48197  
[scasey1@emich.edu](mailto:scasey1@emich.edu)

---



[Volume 23 \(2015\)](#) | [Archive](#) | [Index](#) | [Data Archive](#) | [Resources](#) | [Editorial Board](#) | [Guidelines for Authors](#) | [Guidelines for Data Contributors](#) | [Guidelines for Readers/Data Users](#) | [Home Page](#) | [Contact JSE](#) | [ASA Publications](#)