# Data visualization of between-subject and within-subject factors in JMP
## Chong Ho Yu, Ph.Ds (2012, August 12)

chonghoyu@gmail.com
http://www.creative-wisdom.com/computer/sas/sas.html

The objective of this write-up is to illustrate how data visualization techniques could be utilized to unveil insights from a data set. It is important to point out that this document reflects the personal preferences of the author only. This write-up is by no means suggesting that the following method is the best.

Figure 1 shows a typical 2X2 factorial design with one between-subject factor (gender) and one within-subject factor (two measures: Test 1 and Test 2). It is assumed that a treatment is implemented between the two tests. If the analyst wants to focus on the gender effect, she could employ "General Linear Model" in SPSS by assigning the posttest score to the dependent variable, treating gender as the fixed factor, and making the pretest score a covariate. A covariate is a variable that reflects the pre-existing difference of the subjects. The pretest is used in this modeling to take the prior knowledge possessed by the participants into account when the gender difference in terms of academic performance is inquired.
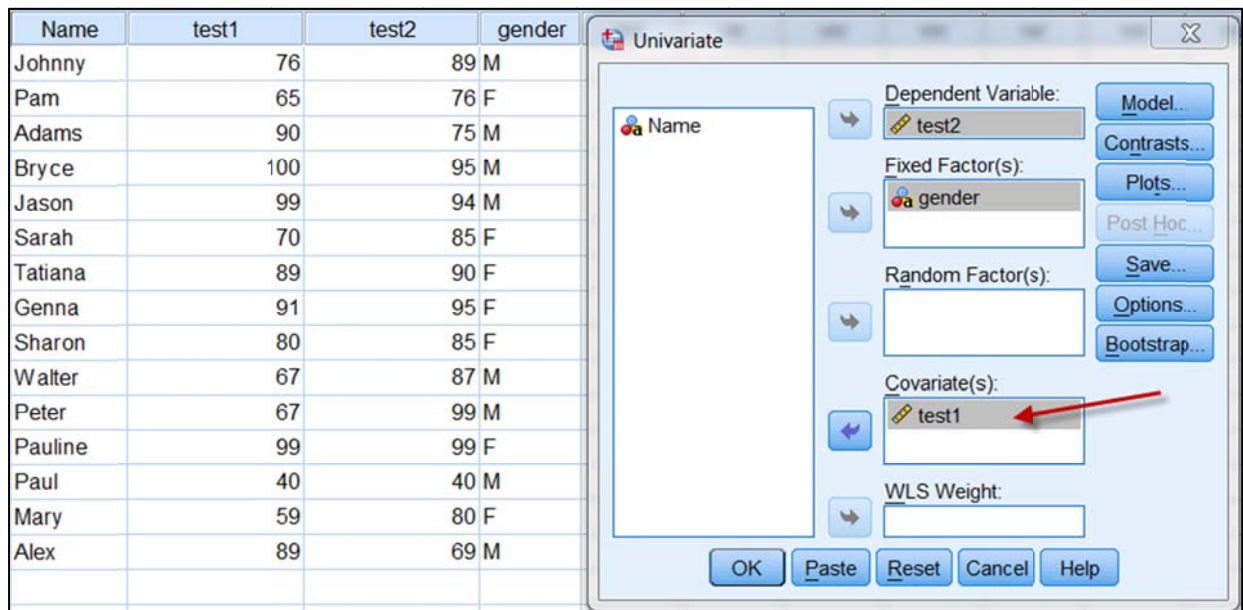


Figure 1. Pretest as a covariate in GML of SPSS.

Figure 2 shows the result of the GLM with a covariate. The *p* value (sig.) is .367, indicating a non-significant gender effect. This report is acceptable, but if the analyst stops right here, something important might be overlooked.

**Tests of Between-Subjects Effects**

Dependent Variable: test2

| Source | Type III Sum of Squares | df | Mean Square | F | Sig. |
|---|---|---|---|---|---|
| Corrected Model | 1434.581ª | 2 | 717.291 | 4.915 | .028 |
| Intercept | 1055.600 | 1 | 1055.600 | 7.234 | .020 |
| test1 | 1293.705 | 1 | 1293.705 | 8.865 | .012 |
| gender | 128.407 | 1 | 128.407 | .880 | .367 |
| Error | 1751.152 | 12 | 145.929 | | |
| Total | 108690.000 | 15 | | | |
| Corrected Total | 3185.733 | 14 | | | |

a. R Squared = .450 (Adjusted R Squared = .359)

Figure 2. GLM output with pretest as a covariate.

The analyst could switch the emphasis from the between-subject factor to the within-subject factor by conducting a dependent t-test in JMP. The right panel of Figure 3 shows that when both genders are included as a single group, the mean difference between Test 1 and Test 2 is not statistically significant ($p$ = .1743).
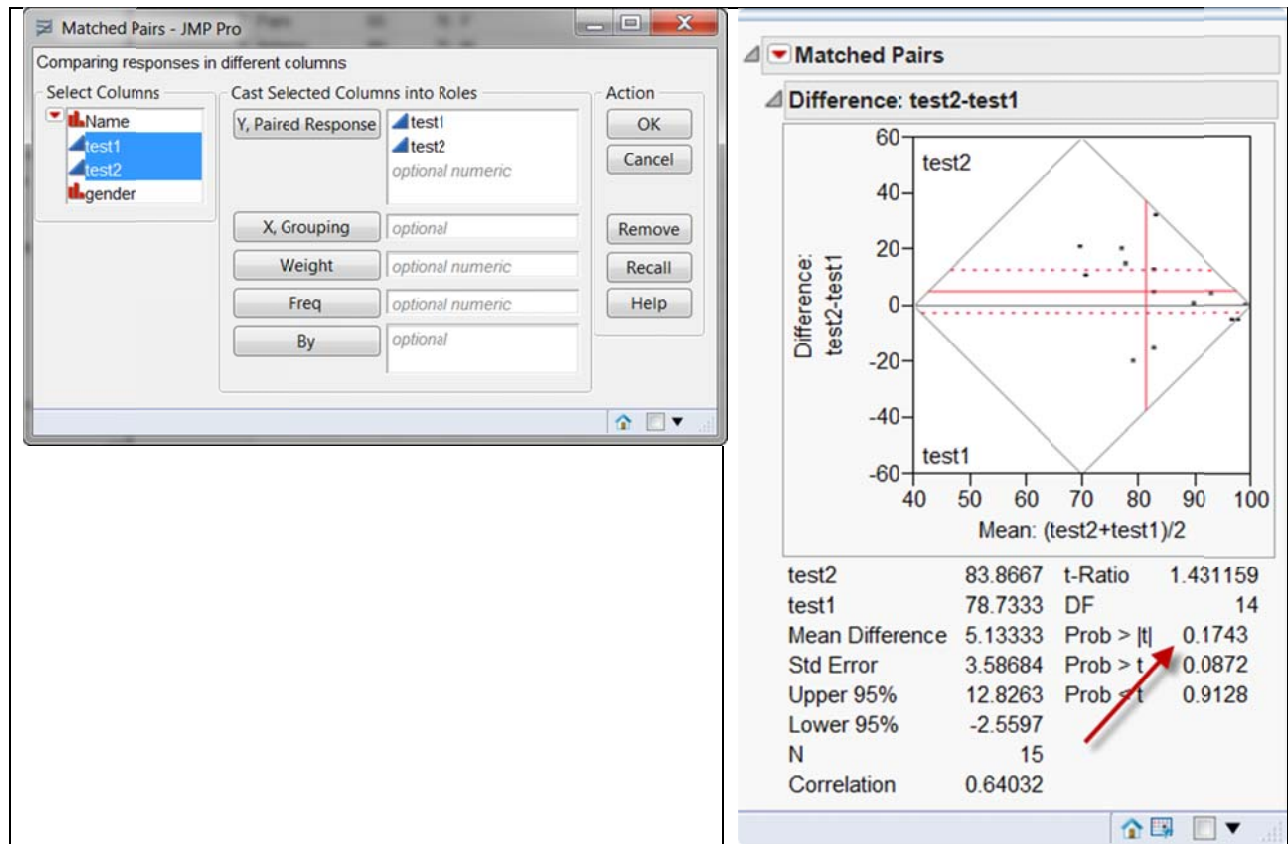


Figure 3. Dependent t-test using all data.

Interestingly enough, when the data set is partitioned by gender (F, M), a different story emerges. The paired t-test using female observations yields a significant result ($p$ = .0327) whereas the t-test using male participants reaches the opposite conclusion ($p$ = .7030) (see Figure 4). In other words, the treatment effect as measured by the change of test performance over time is definitely *moderated* by gender.

It is crucial to mention that in this example gender should be conceptualized as a moderator, not a mediator. Usually a mediator is a variable on the causal pathway. For example, if A affects B and B affects C, then B is considered a mediator between A and C. On the contrary, a moderator is not a causal variable. When the relationship between A and C is not consistent across all the levels of B, this moderating effect does not necessarily imply that B causes C. Put it bluntly, we could not declare that gender (or race) is a "cause" of academic performance. Rather, psychologists would examine the cognitive structure or the cultural traits of certain groups to identify the causal mechanism.
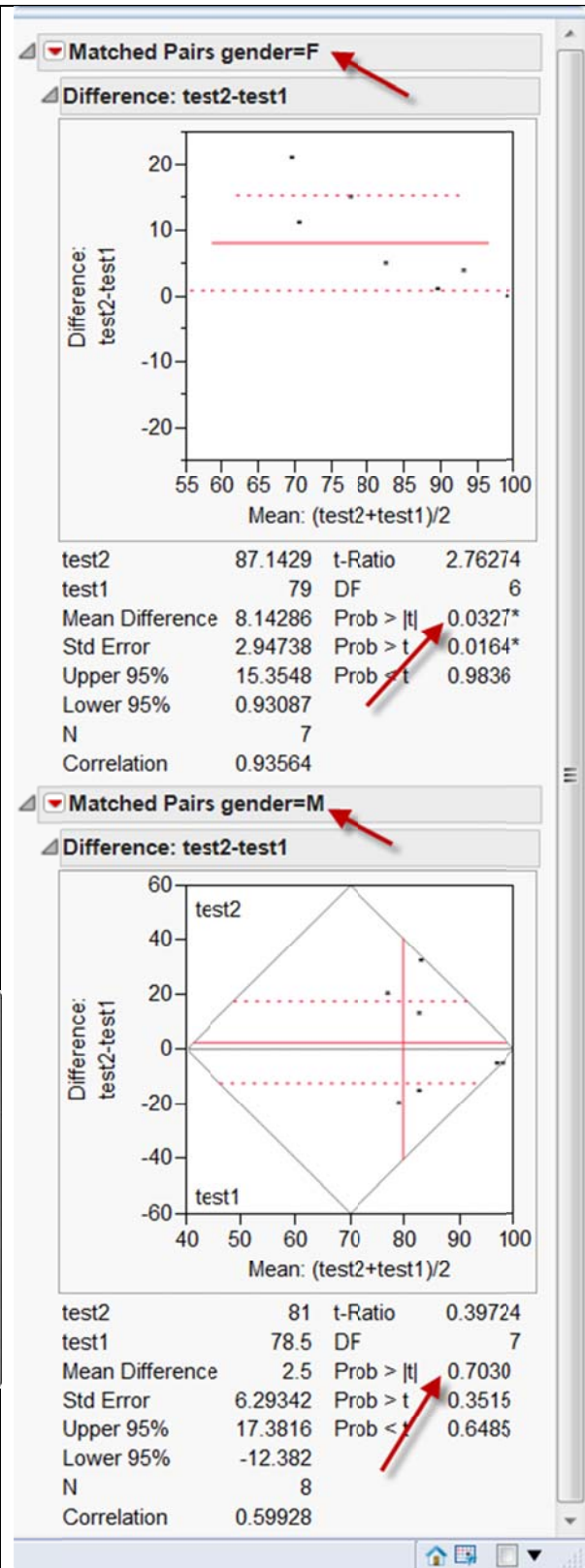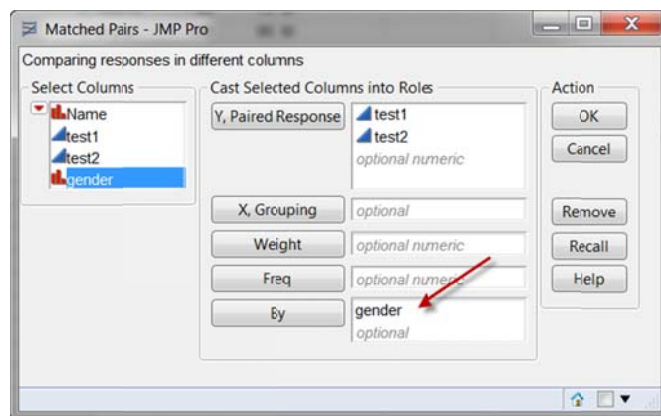


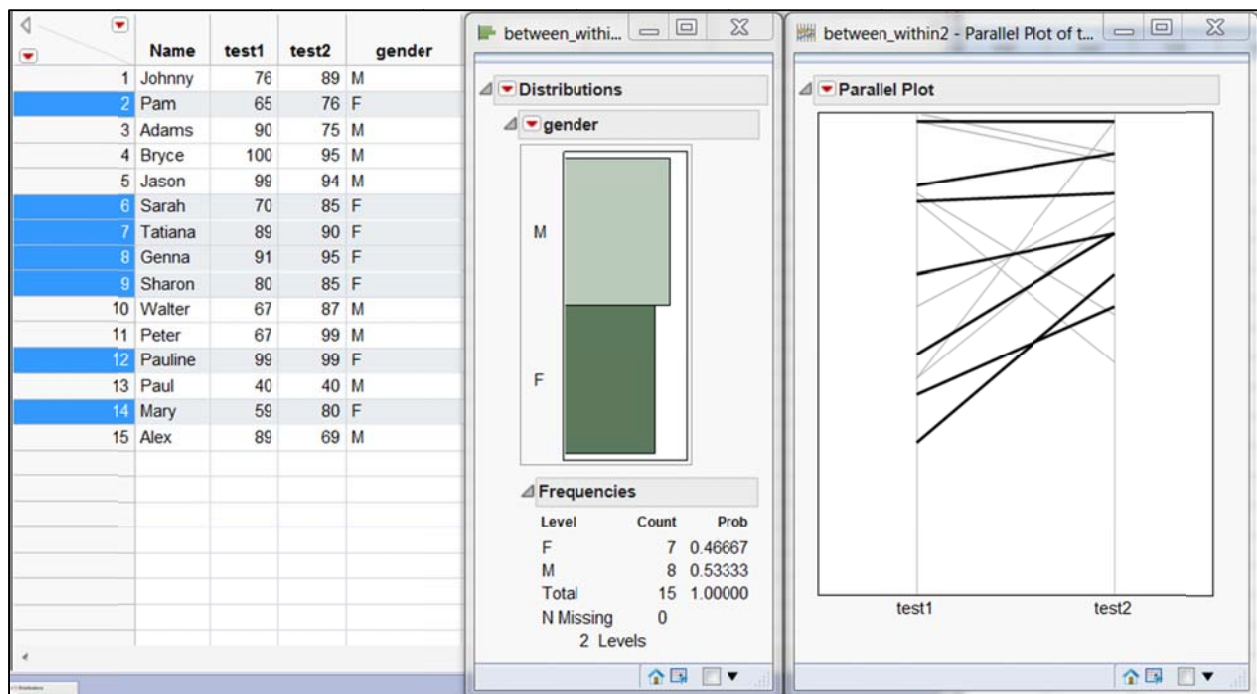Figure 4. Dependent t-tests by gender.

Figure 5. Parallel plot showing female subjects.

Although the preceding analysis by group has revealed some insight, the analyst could further drill down into the individual level. The middle panel of Figure 5 is a bar chart created from the pull down menu "Analyze→Distribution" whereas the right panel is generated from "Graph→Parallel Plot." The parallel plot is a simple line plot connecting the "dots" in Test 1 and Test 2. When the researcher clicks on the bin representing female, the corresponding observations on the parallel plot are highlighted. Obviously, four of them (the bottom four) achieved substantive gains in test performance after the treatment, but three of them (the top three) had minimal gains only. A plausible explanation is that the top three students already know a lot about the subject matter and thus the treatment could not add much to their existing knowledge base. This phenomenon is known as the *ceiling effect*.
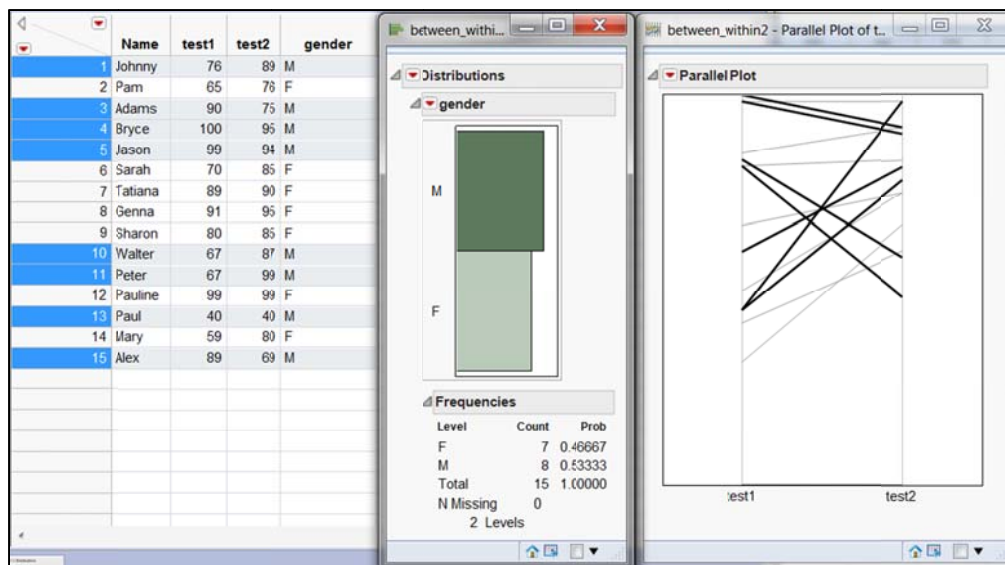
Figure 6. Parallel plot showing male subjects.

The same data visualization approach is applied to the male subjects (Figure 6). Unlike their female counterparts, the test results of the boys are very diverse. Specifically, four male students got worse after the treatment while three of them improved their performance in Test 2. It is noteworthy that those who scored low in Test 2 were the top-performing students in Test 1. Conversely, those who made progress in Test 2 did not do well in Test 1. This is a puzzling phenomenon that awaits further investigation. For example, the steep slope on the parallel plot is very eye-catching and by clicking on the line the researcher found that this student is Peter (Figure 7). To go beyond this data set, she might talk to Peter to find out why he is benefited so much from the treatment program.
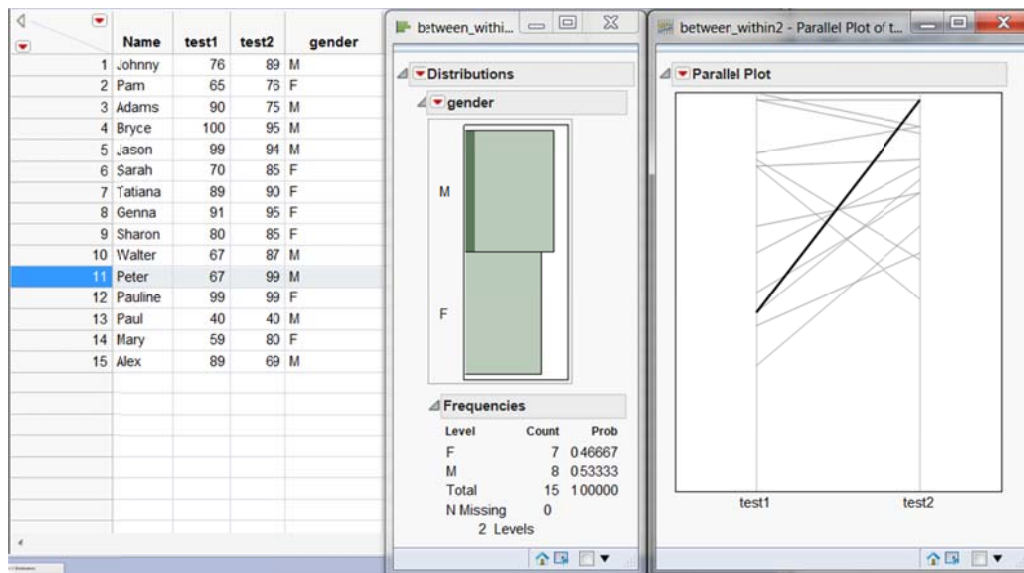


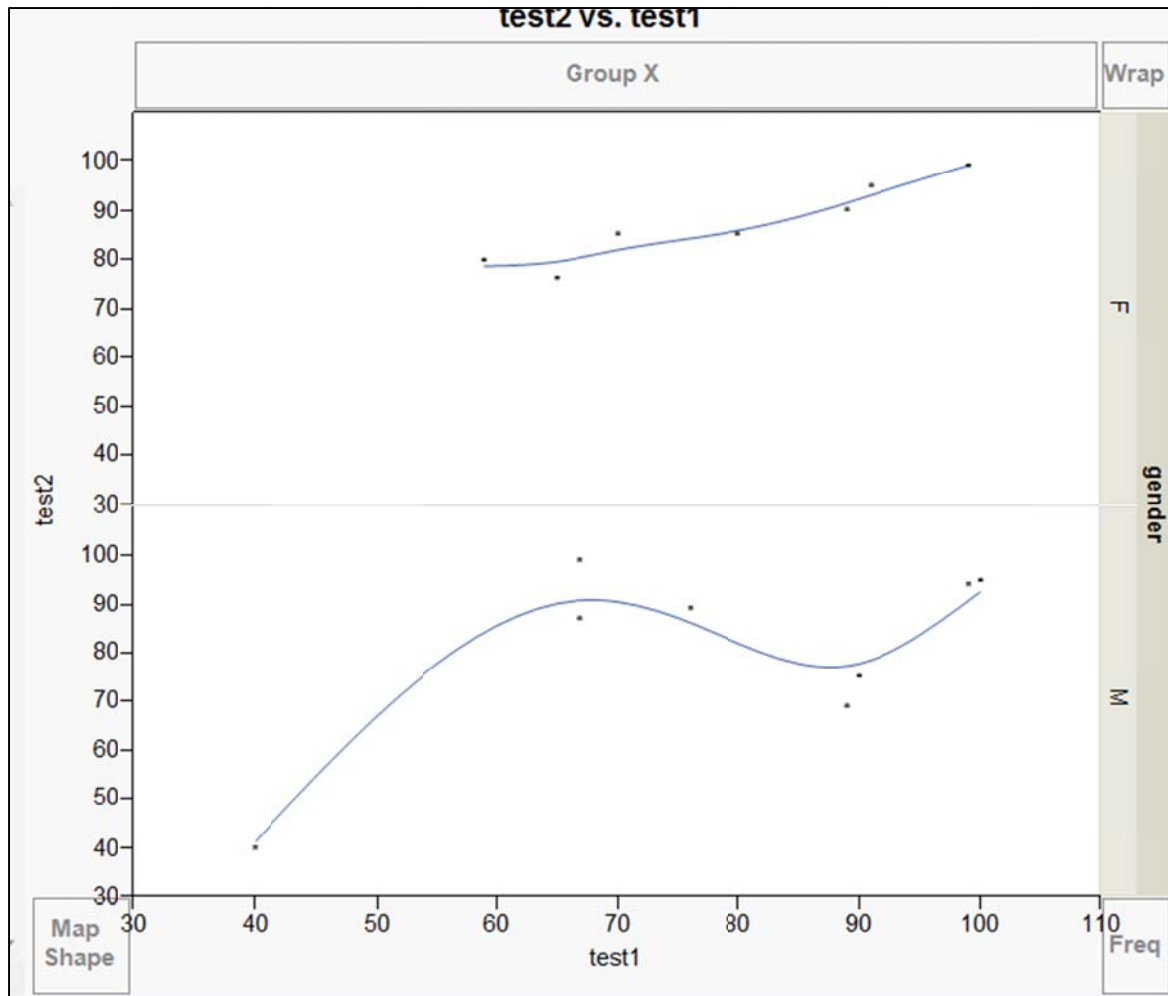Figure 7. Parallel plot highlighting one student.

Figure 8. Scatterplots of pretest and posttest by gender.

Alternately, the analyst could also scrutinize the data using "Graph→Graph Builder." Initially, a regular scatterplot is made by putting Test 2 on the Y-axis and Test 1 on the X-axis. Next, gender is dragged to the right to divide the scatterplot into two panels. It is clear that the top panel (female) shows a linear fit while the bottom one (male) depicts a non-linear fit. In other words, the change of performance over time is consistent among girls, but this varies from student to student in the male group.

As mentioned at the beginning, there is no single best approach to data analysis. Nonetheless, this author strongly recommends interactive data visualization, because sometimes statistics alone could not unearth the hidden patterns of the data.