

TPS REPORT 05

Dispersion , Scaling, and Simple Linear Correlations in Excel

1. Download the TPS 05 Excel file from www.terevaka.net/nau/ant568/tps01.html
2. In this report, you'll be analyzing a variety of different hypothetical datasets. You'll begin with the Depression sheet data, which tracks the short-term benefits of an experimental drug. Our ultimate goal will be to determine which of the eight tests is most closely correlated with random variation in the experimental drug.
3. Although we might be interested in the distributions of data in Columns **A** through **I**, creating a histogram for each Test would be a time-consuming process for the sake of visual analysis. Instead, we can calculate skewness and kurtosis measures to get an idea of what our distributions look like. Remember:

$$\text{Skewness} = \frac{3(\bar{X}-M)}{\sigma} \quad \text{and} \quad \text{Kurtosis} = \frac{\sum\left(\frac{X-\bar{X}}{\sigma}\right)^4}{n} - 3$$

where:

- \bar{X} is the mean
- M is the median
- σ is the standard deviation
- X is the individual score
- n is the sample size

Write a formula in Cell **A33** to calculate the skewness of the sample of doses in the experiment between Cell **A2** and Cell **A31**. (Hint: since we're dealing with a sample, remember to use the unbiased estimate for standard deviation.)

➤ Write your formula in your TPS Report.

Now Copy your formula and Paste it to Cells **B33** through **I33**.

➤ In your TPS Report, indicate which column of data is the least skewed (i.e., most symmetrical) of all.

4. Based on your measures of skewness, which Test do you think will have the strongest correlation coefficient with the Dose data?
➤ Write your answer and explain your reasoning in your TPS Report.
5. In Cell **B34**, write a formula to calculate the standard deviation for Test 1 scores. Copy your formula and Paste it to Cells **C34** through **I34**.

- In your TPS Report, indicate which column of data shows test scores most tightly bunched together (i.e., which Test shows the lowest average distance between data points and the mean)?

Add a label to Row 34 in Cell A34,

6. In Cell B37, write a single formula that can be Copied and Pasted to all Cells between C37 and I37 to calculate the Pearson product-moment correlation coefficient between each column of test scores and Doses in Column A. (Hint: make sure that when you Copy and Paste your formula from B37 to other columns that you are always comparing the test data to the Dose data in Column A.)

- Write the formula from B37 in your TPS Report.

Copy your formula from B37 and Paste to Cells C37 through I37.

Add a label for Row 37 in Cell A37.

7. Suppose your job is to decide which of the tests provides the most valid reflection of behavioral changes associated with differences in dosage of the experimental drug.

- In your TPS Report, explain which test you'd choose and explain your justification.

8. Now you'll focus on the Test-Retest sheet of the same Excel file. Suppose that the test scores in Column B have been selected as the most valid reflection of behavioral changes associated with differences in dosage of the experimental drug. Before we get too confident in that particular test, you'll need to determine if this very expensive test is reliable over multiple applications.

- In your TPS Report, explain what level of test-retest reliability this test provides. Make sure to cite at least one specific statistic to support your explanation.

9. Now turn your attention to the Obama sheet in the same Excel file. Columns A through F represent different polls taken repeatedly in Los Angeles, aimed to get an idea whether the population that was sampled supports Obama's position(s) on minority-rights issues that have emerged over two terms of presidency. Each cell in Column A represents the mean value of poll scores when Poll 1 was administered to 100 randomly-selected residents of L.A., each cell in Column B represents the mean value of poll scores when Poll 2 was administered to 100 randomly-selected residents of L.A., etc.

In Columns A, B, C, E, and F, higher numbers represent stronger support for Obama's actions, but clearly each poll calculated support on a different scale.

In order to convert Poll 4 to a numeric scale, select all of Column **D** and use Ctrl+H (or Command+H for Mac) to replace values based on the following scale:

Strongly support = 5
Support = 4
Neutral = 3
Disagree = 2
Strongly disagree = 1

10. Our next goal is to convert numeric values on various scales from Columns **A** through **F** to **z**-scores so that all measurements are comparable on a single scale centered around a zero value. In general, we use the formula:

$$z = \frac{(X - \bar{X})}{\sigma}$$

where:

- \bar{X} is the mean
- σ is the standard deviation
- X is the individual score

In Cell **I2**, write a formula that can be Copied and Pasted to Cells **I3** through **I51** that will calculate **z**-scores for scores in Cells **A2** through **A51**.

➤ Write your formula in your TPS Report.

Copy your formula from Cell **I2** and Paste to Cells **I3** through **I51**.

In Columns **J** through **N**, calculate **z**-scores for all data in Columns **B** through **F**, respectively.

11. In Cells **A53** and **B53**, calculate the median values of raw mean scores for Poll 1 and Poll 2.
- In your TPS Report, provide a detailed explanation as to which poll's median value, after 50 repeated implementations, indicates a stronger level of support for Obama's actions.
12. The **z**-table at the end of this document offers a useful way to convert **z**-scores to probability values when we are dealing with data that we suspect approximates a normal distribution. Data won't always be distributed normally. However, when we're calculating the mean of many ($n > 30$) random variables (like 100 random poll-participants), and we repeat this random sampling process within the same

population (like Los Angeles), the mean value of our many means should approximate the shape of a normal distribution (remember the Central Limit Theorem?)

As an example of how to use the z -table, Cell J34 should contain a z -score of approximately 1.26. If you look up the corresponding area on the z -table (look at the cell on the table where the z value on the left-most column 1.2 intersects with the z value on the top row of 0.06; i.e., $1.2+0.06=1.26$), you'll see that the area under a normal curve between the mean and a z -score of 1.26 is 0.3962. Or in other words, approximately 40% of the total area under the normal curve is contained in the space between the mean and the z -score of 1.26.

This also tells you, from a predictive standpoint, that if our poll data truly approximates a normal distribution, that we would only expect to see z -scores higher than 1.26 (or equivalently raw mean scores higher than 9 in Poll 2) approximately 10% of the time. That's because under the normal curve, 50% of the area will be to the left of the mean, and we already figured that 40% of the area under the curve will be between the mean and the z -score of 1.26. That leaves only 10% of the area in the tail to the right of our z -score.

Using the z -table, calculate the predicted percentage of raw mean scores that will fall between 0.76 and 0.82 in Poll 6.

➤ In your TPS Report, write your answer and explain how you reached that answer.

13. Finally, turn your attention to the Emissions sheet in the same Excel file. Suppose that the scores represent emission scores for two prototypes and 100 randomly-selected vehicles as they performed in two different areas (highway and city).

Specifically, you are asked to determine whether the new Mitsubishi Volanz or the new Ford Ninny performed better when it came to emissions scores.

➤ In your TPS Report, write your answer and explain how you reached that answer. (Hint: to assess the Volanz and Ninny on the same scale for Highway, City, and Average driving, calculate z -scores).

{ mean, median, stdev.s, correl, pearson }

