

Usable Security and Privacy - Tutorial 3: Statistical Data Analysis

Kami Vaniea

Mohammad Tahaei

1 Introduction

In this tutorial, you will be practising quantitative data analysis on the survey you took the first couple of weeks of class. The survey was comprised of three scales, questions from the CyBok knowledge set, and demographics. In the survey, we used *Security Behavior Intentions Scale (SeBIS)* [1]. The scale, shown in Table 1, has 16-item over four sub-scales that measures attitudes towards device securement, password generation, proactive awareness, and updating. The scale we will be using today is ‘password generation’. Today we will be using the data to answer the following research question:

Research Question: Do people who have prior experience working in a ‘high tech job’ such as computer programming, IT, or computer networking have significantly different behavioural intentions towards password generation than those that do not have such experience?

The research question has two parts: 1) *high tech job* and 2) *behavioural intention towards password generation*. The high tech job part is drawn from a single question on the survey which has only two possible values identifying which group participants belongs to and is, therefore, categorical. The behavioural part is numeric which is computed using one of the scales with values from 1 to 5. Since our research question is asking if one group of people (high tech job), categorical, has better (behaviour), numeric, we can use a *t*-test.

We will use Microsoft Excel in this tutorial, but for your coursework, you can use either the R programming language or Microsoft Excel.

2 Data Analysis

We recommend you work individually on the tutorial, but you can also work in pairs. You can either use your own laptop or a DICE machine. You only need Microsoft Excel installed. The tutorial will also work with LibreOffice Calc. All files are available on USEC coursework page. We will analyse the data through a number of steps. First, you will load the data to your excel. Then, you will compute the average behavioral intention towards password generation per person. After that, we will divide the participants into two groups and compute the *t*-test and interpret the results.

2.1 Load and clean the data

1. Open the [numerics-data.csv](#) file in Excel.
2. Convert the file to Excel format (.xlsx). You can use ‘Save As’ to do this.

Normally you would also need to clean the data of inattentive respondents. Cleaning typically involves removing any incomplete responses as well as removing any participant that failed the attention check questions. We will be skipping the cleaning in this tutorial, but you are expected to do data cleaning in your coursework.

#	<i>Device Securement</i> (28.47% of variance explained; $\lambda = 4.555$)	μ	σ
F4	I set my computer screen to automatically lock if I don't use it for a prolonged period of time.	3.20	1.559
F6	I use a password/passcode to unlock my laptop or tablet.	3.78	1.525
F3	I manually lock my computer screen when I step away from it.	2.63	1.343
F5	I use a PIN or passcode to unlock my mobile phone.	3.21	1.733
#	<i>Password Generation</i> (12.95% of variance explained; $\lambda = 2.071$)	μ	σ
F12	I do not change my passwords, unless I have to. ^r	2.65	1.091
F13	I use different passwords for different accounts that I have.	3.75	1.037
F15	When I create a new online account, I try to use a password that goes beyond the site's minimum requirements.	3.31	1.096
F14	I do not include special characters in my password if it's not required. ^r	3.30	1.292
#	<i>Proactive Awareness</i> (8.36% of variance explained; $\lambda = 1.337$)	μ	σ
F8	When someone sends me a link, I open it without first verifying where it goes. ^r	4.01	1.014
F11	I know what website I'm visiting based on its look and feel, rather than by looking at the URL bar. ^r	3.17	1.077
F16	I submit information to websites without first verifying that it will be sent securely (e.g., SSL, "https://", a lock icon). ^r	3.69	1.102
F10	When browsing websites, I mouseover links to see where they go, before clicking them.	3.69	1.027
F7	If I discover a security problem, I continue what I was doing because I assume someone else will fix it. ^r	4.08	0.976
#	<i>Updating</i> (6.77% of variance explained; $\lambda = 1.082$)	μ	σ
F1	When I'm prompted about a software update, I install it right away.	3.07	1.035
F2	I try to make sure that the programs I use are up-to-date.	3.78	0.890
F9	I verify that my anti-virus software has been regularly updating itself.	3.55	1.228

Table 1: The final questions for the Security Behavior Intentions Scale and associated sub-scales [1]. The means of reverse-scored questions (denoted by r) have been recorded. The final question numbers have been enumerated from F1-F16 in order to differentiate them for easy reference. Responses were reported on the following scale: Never (1), Rarely (2), Sometimes (3), Often (4), and Always (5).

2.2 Password Generation Sub-scale

In this tutorial, we will be using the password generation sub-scale. In this step, we will compute the behavioural intention scale for each participant.

1. Adjust for inverted questions.

In their survey scale, Elegeman et al. [1] used positive and negative statements to minimise acquiescence bias where participants agree on all the statements. In Table 1, statements that are marked with 'r' at the end have been "inverted". The Password Generation scale has two inverted statements 'F12' and 'F14'. To get the correct behavioral intention, you should invert the values for those questions as following:

- Find the column containing 'F12' in the data set.
- Add a new column named 'sebisLikert.sebisF12.-inverted' to the right of the original column.
- To invert all the values automatically, you need to use a formula. Start by entering '= 6 -' and then selecting a cell from 'sebisLikert.sebisF12.'. We subtract from 6 so that a 1 becomes a 5 and a 5 becomes a 1.
- Copy the formula into all cells in the column by dragging the AutoFill Handle of the new cell to the bottom of the column, and the formula will auto apply to all cells in that column.
- Search for 'F14'. Do the same as stated in the above. Call the new column 'sebisLikert.sebisF14.-inverted'.

2. Calculate the average

Next calculate the average score each participant gave the four password generation statements ((F12 inverted, F13, F14 inverted, F15) as follows:

- Add a new column to the spreadsheet named 'password-average'.
- Write this formula in the new column '=average()' (syntax from Excel help: *AVERAGE([number1], [number2], ...)*)
- The average arguments should be the columns from the password generation sub-scale (F12 inverted, F13, F14 inverted, F15)
- Drag the AutoFill handle down to the bottom in the 'password-average' column. This applies the same formula to all cells by dragging the cell to the bottom.

2.3 Student's t -test

The t -test is used in the simplest experimental situation; that is, where there are only two groups to be compared. The test statistic produced by the test is, unsurprisingly, called t and it is the ratio of the difference between means (i.e. the experimental effect) divided by an estimate of the standard error of the difference between those two sample means [2].

To determine if there is a statistically significant difference in password intention between those with and without prior high tech job experience we will be using the t -test. We are using this test because our research question has only two groups (prior high tech job, no prior high tech job) and our dependent variable is a continuous numeric variable (average of password questions).

1. Calculate the t -test:

- To keep the numbers neat, create a new sheet in Excel. Copy 'Password-average' and 'ITbg' columns and paste them in the new sheet next to each other.
- The 'ITbg' column contains our two groups, '1' for participants with high technical jobs and '0' for participants with no experience in highly technical jobs. To split the two groups, select the two columns, 'Password-average' and 'ITbg', and sort them based on the 'ITbg' column. Now you have one group at the top of the table and the other at the bottom.
- Write a new formula in an empty cell: '=t.test()' (syntax from Excel help: *T.TEST(array1,array2,tails,type)*).
 - 'array1' is the part of 'Password-average' where the 'ITbg' equals to '0' (you should select 24 cells).
 - 'array2' is the remaining part of 'Password-average' where the 'ITbg' equals to '1' (you should select 35 cells).
 - 'tails' could be one of two values, '1' for 1-tail and '2' for 2-tails. The 1-tailed t -test is used when we want to see whether technical participants are better than non-technical participants and not the opposite while, in the 2-tailed t -test, we check for either a positive or a negative difference (High tech are better than non-tech or non-tech are better than high-tech). To answer our research question, we will use 2-tailed t -test. Write '2' in the third argument for 2-tail.
 - 'type' should be set to "unequal variance". This setting means that we are not certain that the arrays have equal variation around the mean. To formally prove that they do have equal variance we would use an F-test. For the purpose of this tutorial, you should just assume the default of unequal variance as this is always a safe assumption. The only negative consequence is that we are being too cautious and may therefore have something show up as not significant when it actually is. Write '3' in the last argument for unequal variance.

2. Analyse the result

- The t -test result is called the p-value which is the value that indicates whether there is statically significant difference or not.
- Compare the p-value you get to our selected α of 0.05. If the p-value $< \alpha$, it means the behavioural intention for tech participants is significantly different than the non-technical participants. Otherwise, there is no significance between the two groups.
- The t -test will tell us if the two populations are different, but it won't tell us in what direction the difference is. That is if those respondents with more technical skill have a higher or lower behavioural intention towards passwords. To find the directionality, you need to compute the means of both populations and compare them.

3. If there is time: try testing against another variable

- Research Question: Does prior knowledge about authentication impact behavioural intention towards password generation?
- cybok16 – "Do you know anything about authentication"

3 Discussion (10 minutes)

If time allows, we will do a class-wide discussion on interpreting the results. Think about the following questions individually and we will discuss it if possible.

- Is there any correlation between the tested columns? Is it statistically significant?
- What other research questions might be interesting to try and answer using the survey data?

References

- [1] Serge Egelman and Eyal Peer. Scaling the security wall: Developing a security behavior intentions scale (sebis). In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems*, pages 2873–2882. ACM, 2015.
- [2] Andy Field and Graham Hole. *How to Design and Report Experiments*. SAGE, 2011.