





BAD NEWS: THEY FINALLY DID A META-ANALYSIS OF ALL OF SCIENCE, AND IT TURNS OUT IT'S NOT SIGNIFICANT. https://xkcd.com/2755/

I have the (Developer) Power! Supporting Power Analysis and Statistical Reporting in Usable Security and Privacy

Anna-Marie Ortloff¹, Christian Tiefenau², Matthew Smith^{1,2} ¹: University of Bonn, ²: Fraunhofer FKIE







A Qualitative Study on How Usable Security and HCI Researchers Judge the Size and Importance of Odds Ratio and Cohen's d Effect Sizes

Anna-Marie Ortloff University of Bonn Bonn, Germany ortloff@cs.uni-bonn.de

Simon Lenau CISPA Helmholtz Center for Information Security Saarbrücken, Germany lenau@cispa.de Julia Angelika Grohs University of Bonn Bonn, Germany s6jugroh@uni-bonn.de

Matthew Smith University of Bonn Bonn, Germany Fraunhofer FKIE Bonn, Germany smith@cs.uni-bonn.de



CHI'25





[Find recommendations here]







Background



Hypothesis Testing









UNIVERSITÄT BONN









The probability of detecting an effect, if a true effect exists.







https://rpsychologist.com/d3/nhst/





The strength of the relationship of predictor variables with outcome variables.







https://rpsychologist.com/d3/nhst/





Odds ratio

Ratio of two odds

	Adopted PWM	Did not adopt PWM
Group 1	10	30
Group 2	20	20

 OR=1 means the odds of an outcome are the same in both groups.

Cohen's d

- Normalized difference in means
- d = 0.8 is difference of 0.8 standard deviations between the means



https://rpsychologist.com/cohend/











Do we have the power?



Developer-Centered Usable Security (DCUS)









Literature Collection

- SOUPS, USENIX Security, S&P, CCS, ICSE, USP Tracks of CHI
- 2010 2021
- Include user study
- Participants: software developers, similar expert users, or proxies
- Domain of usable security and privacy



- 54 papers
- including 64 studies, 467
 hypothesis tests, 413 variables

Data Structure

 Relevant information on power and effect sizes in these studies





Power Meta-Analysis (simulated a-priori power analysis)











In SOUPS and CHI USP publications from 2020/2021 only 8 of 74 (10.8%) quantitative papers used a priori power analysis





How can we get the power?







- Use general guidelines for large, medium, small effects
- Use context specific guidelines for large, medium, small effects
- Literature research
- Do a pilot study
- Decide on the smallest effect size of interest



PowerDB

Database + Companion Website



HOME SEARCH TUTORIALS ABOUT

I have the power!

This is the companion website for the paper SoK: I Have the (Developer) Power! Sample Size Estimation for Fisher's Exact, Chi-Squ Wilcoxon Rank-Sum, Wilcoxon Signed-Rank and t-tests in Developer-Centered Usable Security.

For more information, see About.



https://powerdb.info/

Website developed by Ahmad M. Assaf • University of Bonn Copyright © University of Bonn 2023



Searching the database



PowerDB		HOME	SEARCH	TUTORIALS	ABOUT
	Click here to download the database as a SQLite file				
	Operator All selected (AND) Any selected (OR) 				
	Variable				
	- Variable category				
	Participants type				
	fishers				
	Data collection method				
	SEARCH				
	Website developed by Ahmad M. Assaf • University of Bonn Copyright © University of Bonn 2023				



Searching the database



PowerDB			HOME SEARCH TUTORIALS ABOUT
Results:			
1. Test: Fisher's Exact Test			
Paper: On Conducting Security Developer Studies with C (2020) Naiakshina et al.	S Students: Examining a Password-Storage Study with	h CS Students, Freelancers, and Company Developers	
Participants: professional software developer (N=36)			
Dependent Variable	DV Categories	Independent Variable	IV Categories
security "a binary variable secure indicating		prompting condition set by researchers "whether	
whether participants used any kind of security in their code"	Categories for security	the participant is asked to store the password securely"	Categories for prompting
1. Level: yes Any kind of security was used in the code	1. security	1. Level: true participant was prompted for security	1. study related variable
2. Level: no No kind of security was used in the code		2. Level: false participant wasn't prompted for security	
Effect sizes: odds ratio=46.33 Cohen's d=2.11			





PowerDB		HOME SEARCH TUTORIALS ABO			
		Introduction			
		A Cuide te Dewer Analysia			
Introduction		A Guide to Power Analysis			
Fisher's Exact Test	~	for Hypothesis Tests with One Categorical Independent Variable with Two Groups			
Chi-Squared Test	~				
	Ť	Anna-Marie Ortloff Christian Tiefenau Matthew Smith University of Bonn University of Bonn University of Bonn, Fraunhofer FKIE			
McNemar's Test	~				
Independent t-Test	~	What can I expect?			
Paired t-Test	~	What is Power Analysis?			
Wilcoxon Rank-Sum Test	~	Four parameters are relevant to power analysis: power, the significance criterion (i.e. the \(\alpha\) error level), the reliability of the sample results or sensitivity of the test, and the effect size [2]. These four parameters are			
Wilcoxon Signed-Rank Te	st 🗸	interdependent, such that when three of them are available, it is possible to calculate the fourth. Such calculations are referred to as power analysis. In general, there are four different kinds of power analysis, each used to determine both \(\alpha\) and power if a ratio for \(\alpha\) and \(\beta\) is given together with the other two parameters - this is termed			
		compromise power analysis [5]. The other four flavors are summarized, e.g. by Cohen [2] in Chapter 1.5.			
		The Four Parameters Explained			
		All of the mentioned parameters are explained here. Click the cards below to see the definition.			
		Power V			
		Significance Criterion			
		Reliability (Sample Size)			
	the second s				





Example Power Analysis



Simulation for Power Analysis









"Not the Right Question?" A Study on Attitudes Toward Client-Side Scanning with Security and Privacy Researchers and a U.S. Population Sample

Lisa Geierhaas^{*}, Florin Martius^{*}, Arthi Arumugam[†], Matthew Smith[‡] ^{*}University of Bonn, {geierhaa,martius}@cs.uni-bonn.de [†]University of Bonn, arumugam@uni-bonn.de [‡]University of Bonn, Fraunhofer FKIE, smith@cs.uni-bonn.de

S&P'25



Results from Simulation









- Currently unpublished work
 - Comparing two types of intervention in a fully crossed design
- Effect size estimate based on
 - small-N pilot
 - prior work with raw data/frequencies published in the paper





Results from Simulation









Simulation

- One way: 14 / group
- Multi-way:
 - 48 / group
 - 11 / group
- Post-hoc comparisons



Direct

- One-way: 15 / group
- Multi-way
 - 55 / group
 14 / group
- Post-hoc comparisons:
 - 159 / group
 29 / group
 + 9 / group





More Findings from DCUS



Many Statistical Tests






Too? Many Statistical Tests











Specify if you are using confirmatory or exploratory analysis methods.





Problems with Reporting











Report complete and appropriate descriptive statistics.









Make fully anonymized data sets available when needed.





Effect sizes in DCUS









Extension: Effect sizes at CHI





Small, Medium, Large? A Meta-Study of Effect Sizes at CHI to Aid Interpretation of Effect Sizes and Power Calculation

Anna-Marie Ortloff University of Bonn Bonn, Germany ortloff@cs.uni-bonn.de

Theo Raimbault University of Bonn Bonn, Germany s6thraim@uni-bonn.de Florin Martius University of Bonn Bonn, Germany martius@cs.uni-bonn.de

Lisa Geierhaas University of Bonn Bonn, Germany geierhaa@cs.uni-bonn.de Mischa Meier Fraunhofer FKIE Bonn, Germany mischa@mmisc.de

Matthew Smith University of Bonn Bonn, Germany Fraunhofer FKIE Bonn, Germany smith@cs.uni-bonn.de

CHI'25





Find a way to standardize and make machine-readable the research output / statistical results.





Extracting Statistical Values from CHI publications











Report all effect sizes, even non-significant ones.





Overview of Quantitative CHI papers 2019 - 2023







Context-specific effect size guidelines











Take into account context when interpreting effect sizes.







What about interpretation of effects?



P-VALLE INTERPRETATION 0.001 0.01 -HIGHLY SIGNIFICANT 0.02 0.03 0.04 -5IGNIFICANT 0.049 OH CRAP. REDO 0.050] CALCULATIONS. 0.051 ON THE EDGE 0.06 OF SIGNIFICANCE 0.07 HIGHLY SUGGESTIVE, 0.08 SIGNIFICANT AT THE 0.09 P<0.10 LEVEL 0.099 HEY, LOOK AT -THIS INTERESTING ≥0.1 ¯ SUBGROUP ANALYSIS

https://www.xkcd.com/1478/











"[...] significantly increased the number of words in a sentence after which suggestions were requested – **by about 1.5 words**"

– Dang et al.: Choice over Control (<u>https://dl.acm.org/doi/full/10.1145</u> /3544548.3580969)



UNIVERSITÄT BONN





 Clarke et al.: FakeForward (<u>https://dl.acm.org/doi/abs/1</u> 0.1145/3544548.3581100)







INTERPRETATION

-HIGHLY SIGNIFICANT

SIGNIFICANT

OH CRAP. REDO

CALCULATIONS.

ON THE EDGE OF SIGNIFICANCE

P<0.10 LEVEL

HEY, LOOK, AT THIS INTERESTING

SUBGROUP ANALYSIS

HIGHLY SUGGESTIVE, SIGNIFICANT AT THE







"We could not draw a definite conclusion about which algorithm was more accurate because GlanceWriter had a lower error rate but it also had a higher number of error correction actions."

Cui et al.: GlanceWriter (<u>https://dl.acm.org/doi/abs/1</u> 0.1145/3544548.3581100)







"While eye-tracking seems to be the favourite visual attention cue, not all VR headsets have eye-tracking capabilities, and perhaps future cheap VR headsets will never incorporate such capabilities. Our results show that [...], cheap VR headsets using bi-directional CoV can still lead to the same amount of joint attention, therefore, being effective." - Bovo et al.: Speech-Augmented **Cone-of Vision** (https://dl.acm.org/doi/full/10.1 145/3544548.3581283







Interpret and Discuss Effect Sizes.







A Qualitative Study on How Usable Security and HCI Researchers Judge the Size and Importance of Odds Ratio and Cohen's d Effect Sizes

Anna-Marie Ortloff University of Bonn Bonn, Germany ortloff@cs.uni-bonn.de

Simon Lenau CISPA Helmholtz Center for Information Security Saarbrücken, Germany lenau@cispa.de Julia Angelika Grohs University of Bonn Bonn, Germany s6jugroh@uni-bonn.de

Matthew Smith University of Bonn Bonn, Germany Fraunhofer FKIE Bonn, Germany smith@cs.uni-bonn.de

CHI'25



Vignette – Password Manager Example



A study investigated the difference between the adoption of a password manager in a baseline group that received a general introduction to the password manager and an intervention group that additionally was informed that using a password manager was the top recommendation made by security experts. 2100 Participants used the password manager during the study. Two weeks after the end of the study, they were asked whether they were still using the password manager or not.	RQ and study design
511 out of 1050 (48.7%) of those who received the intervention were still using it. 224 out of 1050 (21.3%) of those who did not see the intervention were still using it	Descriptive Stats
Fisher's exact test showed that this difference was statistically significant (p<0.001, odds ratio = 3.50, 95% CI≡[2.88, 4.25]). The effect size (Odds ratio) is 3.50.	Inferential Stats
This means that the odds for the participants to continue to use the formatting tool in the intervention group were 3.50 times higher than for participants to continue to use the password manager in the baseline group.	Effect size explanation



Misconceptions about OR and Cohen's d





Behavioural Security and Privacy Group



Understanding of Effect Size Measures





Effect size type • Cohen's d • Odds ratio







Explain effect size measures.









Report standardized **and** non-standardized effect sizes.





Influencing Factors on Interpretation of Effect Size



- Size
- Context
- Point of view
- Other numerical values in the vignette



Influencing Factors on Interpretation of Effect Size



- Size
- Context
- Point of view
- Other numerical values in the vignette



Judgment Based on Size







Judgment Based on Context





"[Importance] depends

on the topic and the population. Using a password manager is not the same as curing cancer." (from survey)

On effect **size**:

"Comparing an elephant to the earth, it's small, but compared to a mouse it is big" (from interview)







Interpret and Discuss Effect Sizes.







Recommendations for Reporting



Reporting Statistical Results



- Specify if you are using confirmatory or exploratory analysis methods
- Report all effect sizes even non-significant ones
- Report standardized and non-standardized ES
- Interpret and discuss ES
- Explain ES measures



Make Reporting Usable for Future Work



- Report complete and appropriate descriptive statistics
- Make fully anonymized data sets available when needed
- Find a way to standardize statistical research output and make it machine-readable


Recommendations up for Discussion



- Move test statistics to supplementary material
- Consider not reporting p-values and focus on confidence intervals instead



Community Recommendations



Develop / Use Reporting Guidelines.

