1 Title Page.

2	Avoiding Seven of the Most Common Statistical Errors in Scientific Literature		
3	Jordan M Potter, PharmD Candidate 2019		
4	University of Kentucky		
5	jmpo228@uky.edu		
6	270-217-4160		
7	Alexander H Flannery, PharmD, BCCCP, BCPS		
8	University of Kentucky		
9	Alex.Flannery@uky.edu		
10	502-609-5754		
11	Key Words: statistical errors, interpretation, analysis, biostatistics		
12	Pages: 15		
13	Tables: 2		
14	Figures: 1		
15	Appendices: 0		
16	Financial Disclosures: None reported		
17	Conflicts of Interest: None reported		
18			

## 1 ABSTRACT

Conducting statistical analysis is a major component of sharing results. Investigators collect enormous amounts of 2 information when testing a hypothesis and translate that data with statistical tests to present their findings. 3 Scholarship of teaching and learning (SoTL) research should adhere to the same rigorous standards required of 4 5 those dealing with other areas of scholarship and as such, this paper discusses seven common statistical errors likely to be encountered in the scientific literature. While not an all-inclusive list, the following represents areas with a 6 high potential for statistical error and misinterpretation in clinical research and SoTL. More specifically, this paper 7 demonstrates the consequences of assuming normal distributions, examines each component of sample size 8 calculation and its impact on a study, reviews the historical emphasis on p values, examines the components of 9 subgroup analyses, describes methodology to avoid false interpretations of presented results, and identifies 10 resources for ongoing professional development in biostatistics. 11

# 12 INTRODUCTION

In scientific research, conducting statistical analysis is a major component of sharing results. Investigators collect enormous amounts of information when testing a hypothesis and translate that data with statistical tests to present their findings. The accuracy and validity of the results are reviewed rigorously to ensure correctness before being added to the scientific literature. Healthcare providers and educators alike then utilize this literature to make informed decision based on the statistical inferences presented. Unfortunately, an alarming rate of error in statistical analysis exists in the literature resulting in misinformation and misinterpretation. Concern for the implications of these errors has been raised for more than 30 years.<sup>1,2</sup>

Pocock et al reviewed 45 articles published in the *British Medical Journal*, the *Lancet*, and the *New England Journal of Medicine* in 1987 and demonstrated the prevalence of statistical problems even in the most respected journals.<sup>2</sup> In an *Internal Emergency Medicine* article in 2013, Costantino et al countered, stating that of 125 articles reviewed none of the errors would have changed the results of the work. However, in their analysis Costantino et al reports excluding 'more significant items such as study design, outcome, and bias' and still found that 82% of articles analyzed contained errors.<sup>3</sup>

The Accreditation Council for Pharmacy Education has set forth standards for colleges of pharmacy for educating students about biostatistics and study interpretation.<sup>4</sup> Yet a steady rise in published errors exists with a concurrent decline in practitioner confidence.<sup>5-7</sup> In 2009, Bookstaver et al conducted a survey of postgraduate year

2

1 (PGY1) pharmacy residents consisting of 10 knowledge-based biostatistics and study design questions, as well as
2 Likert-type scales to assess confidence.<sup>7</sup> Of those who responded to at least one knowledge assessment (n=166),
3 the overall mean biostatistics knowledge score was 47.3% ±18.50%.<sup>7</sup> In nearly all fields of healthcare education,
4 low levels of confidence in regard to biostatistics are reported, leading students, residents, practitioners, and
5 educators to rely on authors' abstract conclusions to interpret studies.<sup>5-7</sup> The overall perception of confidence and
6 training in statistical analysis may impact junior and senior faculty alike, regardless of their level of training.

As a service to its authors, the *American Journal of Pharmaceutical Education* (AJPE), published a special article outlining several areas of improvement to increase likelihood of publication in response to commonly noticed flaws in submitted manuscripts.<sup>8</sup> In the 2016 article, Persky and Romanelli provided suggestions to potential authors and reviewers, including recommendations for statistics and p values, adding to previously published guidelines regarding the scholarship of teaching and learning (SoTL).<sup>8-10</sup> This article is a follow-up to those papers to further discuss common errors in statistical analyses of clinical research and SoTL.

SoTL should adhere to the same rigorous standards required of those dealing with other areas of scholarship 13 and as such, this paper discusses seven common statistical errors likely to be encountered in the scientific literature.<sup>9</sup> 14 While not an all-inclusive list, the following represents areas with a high potential for statistical error and 15 misinterpretation in clinical research and SoTL. More specifically, this paper demonstrates the consequences of 16 assuming a normal distribution of data, differentiates the components of a sample size calculation and how 17 components impact conclusions from a research study, challenges the historical emphasis on p values, examines 18 the components of a well-done subgroup analysis, describes methodology to avoid false interpretations of presented 19 results, and identifies resources for ongoing professional development in biostatistics. 20

# 21 COMMON STATISTICAL ERRORS IN THE LITERATURE

As articulated by April McGrath in her 2016 article in *Teaching and Learning Inquiry*, "Being aware of these factors...will put scholars in a better position to evaluate...their own research and that of others."<sup>11</sup>

Missing the distribution. When reporting continuous variables, it is common for researchers to summarize the data with a mathematical mean. Calculating the mean allows for quick comparisons between groups and gives readers a sense of what the 'average' value is for any given variable. The underlying assumption when calculating a mean is that the continuous data is normally distributed. By definition, about 67% of the values of a normal distribution are within  $\pm 1$  standard deviation of the mean, and about 95% are within  $\pm 2$  standard deviations.<sup>12</sup>

Unfortunately, most biological data is not normally distributed. Issues then arise when readers see a mean and 1 standard deviation reported and falsely assume that the values from the sample are normally distributed, when in 2 3 reality they may be skewed due to outliers. To demonstrate this, imagine a group of 5 students have taken an exam and receive scores of 96, 89, 94, 97, and 13. In assessing the exam for appropriateness, you review the average test 4 score (77.8%) and question the difficulty of the exam. This is a simple example, but it is clear to see that these 5 6 scores are not evenly distributed and suggested a falsely lower average score when in reality the average was much more likely to be in the mid-90s. In addition to misrepresenting descriptive statistics, falsely assuming a normal 7 distribution of continuous variables may yield inaccurate results from statistical analysis. To appropriately analyze 8 9 non-normally distributed data, researchers must use non-parametric tests, rather than the parametric tests used for normally distributed data.<sup>13</sup> 10

There are multiple ways to avoid this common error. What is likely the most simple solution to avoid this error is to create a histogram of the data which is being collected. This graphical representation will allow a visual check of normal distribution. If all data points are available, another simple check for normal distribution is a comparison of the mean, median, and mode. If the data is normally distributed, these three values will be equal. As a final possibility to test the normality of data, goodness-of-fit there statistical tests may be run including the Shapiro-Wilk and Kolmogorov-Smirnov tests.<sup>14</sup> Unfortunately, all of these assessments require authors to present all of the data and readers are often left to assume that sample distributions are normal.

Details of the sample size calculation are key to making accurate conclusions. Calculating an 18 appropriate sample size (n) is an integral step in the research process.<sup>15</sup> Deciding the number of people or other 19 experimental units to involve in a study is a seemingly simple task. However, sample size calculations must balance 20 the financial burden of increasing sample size while maintaining the ability to detect an important effect should one 21 exist. In other words, a sample size too large wastes money, time, and resources and a sample size too small may 22 lack the power to answer the study question.<sup>15</sup> There are four components of a sample size calculation: alpha (type 23 I error), beta (type II error), the minimal relevant difference, and the baseline event occurrence. Optimizing the 24 sample size is challenging and as such, the details of the sample size calculation are key to making accurate 25 conclusions as a reader. It is important to note here that a statistical test can either reject or fail to reject a null 26 hypothesis, but can never prove the hypothesis to be true.<sup>13,16</sup> Any claims by the author of 'proving the alternate 27 hypothesis' should raise concern. 28

The values of alpha and beta are traditionally chosen to be 0.05 and 0.20, respectively. These values, while 1 arbitrary, are conventional in the scientific literature.<sup>13</sup> Any other stated values should be investigated by the reader, 2 3 if no compelling explanation is given in the text. The value of alpha establishes the acceptable probability of a type I error for the given hypothesis test. Type I error is falsely rejecting the null hypothesis when it is true, therefore, 4 5 with a traditional alpha of 0.05, the researchers accept a <5% chance of finding a false-positive conclusion.<sup>15,17</sup> Beta, or type II error, is the probability of failing to reject the null hypothesis when it is false.<sup>17</sup> The conventional beta of 6 0.20 tells readers that the researchers accept a <20% chance of a false-negative conclusion.<sup>15</sup> Power is the 7 complement of beta (1-beta). It is equivalent for investigators to report a power of 0.80 or 80% versus stating a beta 8 of 0.20 or 20%. 9

The third component of a calculating a sample size is determining the smallest effect of interest. The concept 10 of estimating an effect size prior to collecting the data seems counterintuitive. However, the estimated effect size is 11 critical to the calculation of the number of patients, students, or data points needed for the study. Since sample size 12 is inversely proportional the square of the expected effect size, even small changes in the expected difference have 13 major implications on the sample size required.<sup>15</sup> To demonstrate this, a hypothetical example is provided. Say you 14 15 created a training program for pharmacy residents to increase publication rates of pharmacy resident projects. As a proud creator, you believe that this program can increase publications by 25%. To test your hypothesis, you 16 randomize residents to either receive the course or complete standard residency training. Completing a sample 17 calculation (Table 1) yields a required sample of 58 residents per group. While recruiting 216 residents may be 18 feasible, estimating that your training will increase publication rates from 50% to 75% is unrealistic and will likely 19 result in failure to reject the null hypothesis of no statistical difference between residents that completed the training 20 and those who did not. Estimating smaller effect sizes allow for greater sensitivity to detect a difference, but requires 21 exponentially more data points (Table 1). As a researcher and a reader, it is important to recognize this balance of 22 producing accurate results, while maintaining a reasonable sample size. 23

The final component of the sample size calculation is baseline event and variance for the population. Investigators primarily utilize previous studies or pilot studies and their own knowledge and opinions to estimate the difference of baseline event and variance.<sup>17</sup> This has inherent unreliability, however, it is the only way to preemptively estimate effects in new populations. Sample size calculation is integral in providing adequate power to detect differences between groups. Yet,
only 16% of RCTs in major journals had sufficient statistical power (80%) to detect 25% relative difference and
only 36% were powered to detect a 50% relative difference.<sup>18</sup> When reviewing positive claims, remember that the
details of the sample size calculation are key to making accurate conclusions.

5 Don't overly rely on p values: other measures/methods may tell you much more. As stated in Perskey and Romanelli's Insights, Pearls, and Guidance on Successfully Producing and Publishing Educational Research, 6 "...impactful published papers almost always have statistical analysis," and that means p values are plentiful.<sup>8</sup> P 7 values are considered to be the evidence in favor of random chance as an explanation for a result.<sup>17</sup> This 8 mathematical probability is then compared to the predetermined alpha, which as stated earlier is an arbitrary cutoff. 9 Then, an all-or-none decision is made as to whether the difference between groups was statistically significant or 10 not. P values are especially prone to misinterpretation and as academics it is imperative to recognize that even 11 12 statistically significant results tell only part of the story.

First, it is necessary to recognize that statistical significance does not measure effect size or relevance and 13 importance of a result. Statistical significance is highly dependent on sample size and even very small effects can 14 15 appear statistically significant when the sample size is large enough.<sup>17</sup> In the ALBIOS trial, a multicenter randomized controlled trial with 1818 ICU patients with septic shock or severe sepsis, the albumin group had a 16 significantly lower heart rate than those in the crystalloid group (p=.002).<sup>19</sup> Shown only in the appendix, however, 17 are the heart rates of the two groups; 89±20 versus 92±20 beats per minute.<sup>19</sup> As a reader, it is imperative to 18 differentiate statistical significance from a relevant result in the field of study. Similarly, a non-statistically 19 significant result cannot be equated to a statement of no effect. A p value of p=.049 (assuming an alpha of .05) 20 would be considered significant, whereas p=.051 would not be statistically significant. In the same manner that a 21 decrease of three beats per minute in heart rate was statistically significant, but not clinically relevant in the ALBIOS 22 trial, a meaningful difference could yield a non-significant result in a different study. A careful assessment of effect 23 size is always necessary and is most easily conducted with review of the confidence intervals. Another erroneous 24 practice is equating the relative size of a p value to being more or less significant or likely to be true. P values do 25 not measure the probability that the hypothesis is true, or the probability that the result was a random effect.<sup>20</sup> 26 Therefore, p=.000001 is no more significant than p=.01, nor is the result any more likely to be true. 27

Lastly, p values are only credible if all the assumptions from the statistical test were met. Since p values 1 are a direct output from statistical analysis, the underlying assumptions for the respective test must be met for the p 2 value to be accurate. Examples may include using a parametric test for non-normally distributed data or applying a 3 linear regression without first verifying that the relationship between variables is indeed linear.<sup>20</sup> P values alone 4 don't provide evidence with regard to a hypothesis, yet unfortunately, many articles only report p values.<sup>21</sup> Be sure 5 to examine the whole picture when making conclusions from p values. Look for effect sizes and confidence intervals 6 to aid in your assessment. These convey two critical factors that a p value lacks, the magnitude of an effect and the 7 relative importance of an effect. 8

Multiple hypothesis testing is not without risks. Testing multiple hypotheses may seem benign. Costs 9 associated with conducting research to answer study questions are high and increase constantly. Logically, 10 researchers want to maximize the return on investment and answer as many questions as possible with the data they 11 collect. Technology continues to advance and the wide availability of user-friendly statistical software have made 12 it possible for investigators to run hundreds of tests with ease. Using software in this manner is inappropriate and 13 introduces significant risk of error into the analysis. This is demonstrated in Figure 1, where type I error is plotted 14 as a function of the number of subgroup analyses run. The risk of error in running 100 subgroup analyses may seem 15 obvious, but it is not uncommon for studies to run 20 subgroup analyses. Conducting 20 hypothesis tests increases 16 the probability of type I error from .05 to .64 (assuming  $\alpha$ =.05). When reviewing work with multiple subgroup 17 analysis, be aware that the likelihood that at least one of the analyses will be 'statistically significant' by chance 18 alone also increases. Redundant or repetitive statistical analysis of research data conducted in an attempt to find 19 significance is often referred to as 'data dredging.'22 20

Outcomes and subgroups should be determined *a priori*, or before any data collection and analysis has 21 occurred. This is not always possible and often times analyses are performed *post-hoc*, or after data has been 22 collected. In this case, the number of analyses, the selection process for these analyses, and the significant and non-23 significant results of these analyses should all be disclosed by the author.<sup>16</sup> Authors may also use statistical 24 correction factors, such as the Bonferroni or similar procedures, to adjust the levels of significance accordingly 25 when conducting multiple statistical tests or comparisons.<sup>22</sup> Corrections should be stated by the author so that 26 readers and reviewers can be assured that significance levels have been adjusted appropriately to be more 27 conservative.23 28

Be appropriately skeptical of subgroup analyses. Subgroups analysis can be defined as analysis of 1 intervention effects within subgroups of the sample. Subgroups are beneficial because they can expand the number 2 3 of questions that may be asked from a limited data set, such as, 'Did the sicker patients benefit from the drug?' or 'Do certain educational strategies benefit certain subgroups of students?' However, the analysis of subgroups should 4 5 be approached cautiously. The Journal of the American Medical Association published a 'users guide' to subgroup analysis and provide a series of questions that should be asked when assessing the validity of subgroup analyses: 6 Can chance explain the apparent subgroup effect?; Is the effect consistent across studies?; Was the subgroup 7 hypothesis one of a small number of hypotheses developed a priori with direction specified?; and Is there strong 8 preexisting support?<sup>24</sup> 9

In addition to the inherent risk of error associated with testing multiple hypotheses previously described, 10 readers must be appropriately skeptical of an author's claims based on subgroup analysis. One limitation of 11 subgroup analysis is the possibility that the differences or associations found could be spuriously positive.<sup>12</sup> For 12 example, in 1988 the Second International Study of Infarct Survival (ISIS-2) investigators reported an apparent 13 subgroup effect in the randomized trial comparing streptokinase and aspirin in suspected cases of acute myocardial 14 infarction.<sup>25</sup> "Patients ... born under the zodiac signs of Gemini or Libra did not experience the same reduction in 15 vascular mortality attributable to aspirin that patients with other zodiac signs had."<sup>25</sup> Obviously, there is no 16 biological explanation for these statistically significant findings and the investigators made this claim to 17 demonstrate the increased risk of type I error when conducting subgroup analysis.<sup>24,25</sup> Another major disadvantage 18 of subgroup analysis is the risk of making a type II error due to smaller sample size in the subgroup. Readers may 19 interpret the lack of a difference between groups to mean that there is no difference when in reality a difference 20 may exist and the group lacks the statistical power to detect it. Subgroups that are prespecified a priori, or at 21 baseline, increase the validity of the analysis. In these instances, stratification and regression techniques may be 22 used to adjust the overall comparison for subgroups.<sup>24</sup> 23

As a reader, a certain level of skepticism is necessary when reviewing the results of any subgroup analysis. Subgroup analyses are highly susceptible to false positives from multiple comparisons, false negatives due to inadequate power, and often lack transparency in the validity of the conduction.<sup>24</sup> Due to these limitations the information garnered from the analysis of subgroups, even under the best circumstances, should be reserved for generating hypothesis for further study.

Details of modeling are just as important as the rest of the methods section. Even the simplest models, 1 examining the relationship between only two variables, are prone to error.<sup>17</sup> One of the first steps in evaluating a 2 model is assessing its plausibility. Correlation does not confer causation. This adage may seem obvious for variables 3 that appear correlated without any reason for connection. However, when looking at multiple factors, correlations 4 can be deceptive. For instance variable X may demonstrate a correlation statistical significance with variable Y, 5 solely because X and Y are both dependent on a third variable Z.<sup>26</sup> When presenting results, the investigator should 6 provide details on the modeling procedure. These explanations should be thorough and include the methods used 7 to create the model, the assumptions of the data that were made, the limitations of applicability of the model, 8 potential sources of bias, and the methods of validation that were used in the creation of the model.<sup>26</sup> In addition to 9 these criteria, the authors should present an assessment of the fit of the model. Use of a linear regression model is 10 predicated on the fact that the relationship between variables is in fact linear. A variable may yield a statistically 11 significant linear regression despite the underlying relationship not being linear.<sup>17</sup> This should be assessed with an 12 analysis of 'residuals'. Further testing of model fit should be described by goodness-of-fit tests, such as the Hosmer-13 Lemeshow test.<sup>27</sup> This test measures the statistical significance of any differences between the observed and 14 15 predicted outcomes over the risk groups. A model that is properly fitted with yield a non-significant difference and suggest that the model is appropriate.<sup>27</sup> This test can be easily misinterpreted. For instance, a retrospective 16 observational study reviewing two different forms of anticoagulation in patients with a traumatic brain injury 17 claimed "an unexpected finding ... was the superiority of LMWH over UH with regards to mortality" despite a 18 resulting Hosmer-Lemeshow goodness-of-fit test of p=.066.<sup>28</sup> Readers should also recognize collinearity, two 19 variables put in the model that measure the same thing and predict each other. 20

Testing too many covariates to fit a model can lead to 'over-fitting' of the model. As a general rule of 21 thumb, for every covariate to be tested, there should be somewhere between ten and twenty outcomes of interest. 22 Similar to 'data dredging' of multiple hypothesis testing, the number of covariates you include in the model increase 23 risk of spurious findings. This is not uncommon to see in the literature, but a striking example is the retrospective 24 analysis of postoperative urinary tract infections by Sultan et al.<sup>29</sup> This study had a large sample size (n=891) of 25 patients who had a Foley catheter placed by either a surgical resident, operating room nurse, or medical resident.<sup>29</sup> 26 Sultan et al claim that medical students need more monitoring because "patients with Foleys placed by medical 27 students were at over 4-fold increased risk of [catheter-associated urinary tract infection]" (OR 4.09, p=.02) 28

compared to nurses.<sup>29</sup> The issue in making this claim is that only 2.4% of sample experienced the outcome of interest 1 and the authors put 20 factors into a logistic regression model where there were only 22 events. In this instance, 2 parameters are biased toward the extremes when the number of variables approaches the number of events and it 3 becomes nearly impossible to meaningfully conduct a regression. Correcting for these errors would only widen 4 confidence intervals for the studies claim that medical students need more supervision. Taking all of this into 5 account, it is evident that the results of the author's multivariate analysis (OR 4.09, CI 1.22-13.7; OR 2.16, CI 0.75-6 6.2) simply does not support the conclusion.<sup>29</sup> Like many of the previously listed, the statistically significant 7 findings generated from modeling should be further researched, rather than a basis for conclusive findings. 8

It's relatively risky...but it doesn't have to be. An important distinction to make when reviewing and 9 interpreting estimated effect size between two groups is the distinction between the odds ratio, relative risk, and 10 absolute risk for a dichotomous outcome.<sup>30</sup> These summary statistics are not interchangeable and result in different 11 values due to the different calculations required. To start, odds are the probability of occurrence of an event or 12 outcome compared to the probability of the event or outcome not occurring. The odds ratio (OR) is then the ratio 13 of the odds of an event or outcome occurring in one group to the odds of the same event or outcome occurring in 14 the other group. Odds ratios are the estimate of effect size for retrospective data.<sup>31</sup> Alternatively, for prospective 15 data results can be expressed in terms of risk, the probability of occurrence of an event or outcome. Typically risk 16 is presented as risk reduction and can be expressed in relative terms as relative risk reduction (RRR) or absolute 17 terms as absolute risk (ARR) and number needed to treat (NNT).<sup>31</sup> RRR estimates the percentage of baseline risk 18 that is reduced as a result of an intervention.<sup>31</sup> The estimate, however, only applies to the study population and has 19 a tendency to over-estimate the beneficial effects in the intervention group.<sup>30</sup> 20

Despite recommendation from the CONSORT guidelines for randomized controlled trials, authors do not 21 always report both the relative and absolute effect sizes.<sup>32</sup> Calculating risk in absolute terms provides context for 22 the reduction of risk of an outcome by accounting for likelihood that the event will occur in the population at large. 23 Since authors often do not report ARR, it is necessary to calculate it to accurately assess the results of the study. To 24 demonstrate ARR calculation from RRR, consider this hypothetical example provided by Streiner and Norman in 25 their *Chest* commentary looking at the relative risk of death between two groups.<sup>30</sup> In the hypothetical study, 20 26 individuals died in the intervention group (n=100) and 40 individuals died in the control group (n=100). The relative 27 risk of death in the intervention group is half that of the control group. Then, instead of a sample of 200 individuals, 28

the hypothetical study was conducted in 20,000 with the same absolute number of deaths, respectively, yielding RR=0.5 again. This clearly demonstrates how only reporting the relative risk can be very misleading. Calculating the absolute risk reduction incorporates the baseline risk associated with the outcome (see citation for calculation).<sup>30</sup> In the first example, ARR=0.2, or a 20% absolute risk reduction in death in the intervention arm. In the second example, with the same relative risk, the absolute risk reduction is 0.002, or a 0.2% absolute reduction of risk of death. Converting these values into NNT, the number of participants required to receive the intervention before preventing one outcome, further demonstrates this difference (NNT=5 vs. NNT=500).

8 Speaking to results in relative terms is risky. The example in *Chest* clearly demonstrates how relative risk 9 reduction can be a massive overestimate of effect size.<sup>30</sup> Calculating risk and risk reduction in absolute terms, if not 10 presented in the manuscript, is an easy calculation that should always be conducted as a reader to draw appropriate 11 conclusions about the study.

#### 12 CONCLUSION

The misinterpretation and abuse of statistical analyses has been established many years, yet their presence is nearly ubiquitous. Only about 20% of published manuscripts are free from statistical error.<sup>33</sup> Despite an increasing number of guidance documents/calls for standardization, the key problem remains that there are no interpretations of these concepts that are simultaneously simple, intuitive, correct, and foolproof.<sup>21</sup> As previously mentioned, this article is not an all-inclusive list of statistical errors in the literature. There are, however, many additional resources available to foster knowledge in this area, some of which are given in Table 2. Quality instruction and teaching are the most important goals for higher education, but research efforts help shape best educational practices.<sup>10</sup>

Statistical analyses are a key part of the evidence for making conclusions that impact patient care, 20 educational methods, and allocation of resources. While cumbersome and often complex, formal training in 21 statistics is not a requirement for due diligence as a scholar. This article serves as a general primer so that commonly 22 encountered mistakes and misinterpretations of academic research will not continue to plague authors and readers 23 alike. Recent increases in the demand for evidence based decision making and quality improvement have placed a 24 greater emphasis on producing and evaluating the scientific literature.<sup>8,10</sup> The goal of this work is to provide authors, 25 instructors, and readers of the scientific literature a resource to avoid common pitfalls of statistical analysis in the 26 ever evolving field of health profession education. As academics, a focus on the interpretation and execution is an 27 integral part of practice and should be part of continued professional development. 28

## **1 ACKNOWLEDGMENTS**

The authors have no conflicts of interest and no financial disclosures to report.

#### **3 REFERENCES**

2

Glantz SA. Biostatistics: how to detect, correct, and prevent errors in the medical literature. *Circulation*.
 1980;61(1):1-7.

Pocock SJ, Hughes MD, Lee RJ. Statistical problems in the reporting of clinical trials. A survey of three
medical journals. *N Engl J Med.* 1987;317(7):426-432.

8 3. Costantino G, Casazza G, Cernuschi G, et al. Errors in medical literature: not a question of impact. *Intern*9 *Emerg Med.* 2013;8(2):157-160.

10 4. Accreditation Council for Pharmacy Education Standards 2016. Chicago, IL2015.

5. Windish DM, Huot SJ, Green ML. Medicine residents' understanding of the biostatistics and results in the
 medical literature. *Jama*. 2007;298(9):1010-1022.

Retsas A. Barriers to using research evidence in nursing practice. *Journal of advanced nursing*.
 2000;31(3):599-606.

Bookstaver PB, Miller AD, Felder TM, Tice DL, Norris LB, Sutton SS. Assessing pharmacy residents'
 knowledge of biostatistics and research study design. *Ann Pharmacother*. 2012;46(7-8):991-999.

Persky AM, Romanelli F. Insights, Pearls, and Guidance on Successfully Producing and Publishing
 Educational Research. *American Journal of Pharmaceutical Education*. 2016;80(5).

19 9. Poirier T, Crouch M, MacKinnon G, Mehvar R, Monk-Tutor M. Updated Guidelines for Manuscripts

Describing Instructional Design and Assessment: The IDEAS Format. *American Journal of Pharmaceutical Education*. 2009;73(3).

McLaughlin JE, Dean MJ, Mumper RJ, Blouin RA, Roth MT. A Roadmap for Educational Research in
 Pharmacy. *American Journal of Pharmaceutical Education*. 2013;77(10):218.

McGrath A. Searching for Significance in the Scholarship of Teaching and Learning and Finding None:
Understanding Non-Significant Results. *Teaching & Learning Inquiry: The ISSOTL Journal*. 2016;4(2).

Lang T. Twenty statistical errors even you can find in biomedical research articles. *Croatian medical journal*. 2004;45(4):361-370.

12

- Bajwa SJ. Basics, common errors and essentials of statistical tools and techniques in anesthesiology
   research. *J Anaesthesiol Clin Pharmacol.* 2015;31(4):547-553.
- 3 14. Sürücü B, Koç E. Assessing the validity of a statistical distribution: some illustrative examples from
  4 dermatological research. *Clinical and Experimental Dermatology*. 2008;33(3):239-242.
- 5 15. Noordzij M, Tripepi G, Dekker FW, Zoccali C, Tanck MW, Jager KJ. Sample size calculations: basic
  6 principles and common pitfalls. *Nephrol Dial Transplant*. 2010;25(5):1388-1393.
- 7 16. Thiese MS, Arnold ZC, Walker SD. The misuse and abuse of statistics in biomedical research. *Biochem*8 *Med (Zagreb)*. 2015;25(1):5-11.

9 17. Good PI. Common Errors in Statistics (and How to Avoid Them). Hoboken, New Jersey2003.

18. Moher D, Dulberg CS, Wells GA. Statistical power, sample size, and their reporting in randomized
controlled trials. *Jama*. 1994;272(2):122-124.

- Caironi P, Tognoni G, Masson S, et al. Albumin replacement in patients with severe sepsis or septic shock.
   *N Engl J Med.* 2014;370(15):1412-1421.
- Wasserstein RL, Lazar NA. The ASA's Statement onp-Values: Context, Process, and Purpose. *The American Statistician*. 2016;70(2):129-133.
- 16 21. Abbott EF, Serrano VP, Rethlefsen ML, et al. Trends in P Value, Confidence Interval, and Power Analysis
- 17 Reporting in Health Professions Education Research Reports: A Systematic Appraisal. Acad Med. 2018;93(2):314-
- 18 323.

19 22. McGaghie WC, Crandall S. Data Analysis and Statistics. *Academic Medicine*. 2001;76(9):936-938.

20 23. Regehr G. Reporting of Statistical Analyses. *Academic Medicine*. 2001;76(9):938-939.

24. Sun X, Ioannidis JP, Agoritsas T, Alba AC, Guyatt G. How to use a subgroup analysis: users' guide to the
medical literature. *JAMA*. 2014;311(4):405-411.

23 25. Randomized trial of intravenous streptokinase, oral aspirin, both, or neither among 17,187 cases of
 24 suspected acute myocardial infarction: ISIS-2.ISIS-2 (Second International Study of Infarct Survival) Collaborative

25 Group. *J Am Coll Cardiol*. 1988;12(6 Suppl A):3A-13A.

26 26. Eddy DM, Hollingworth W, Caro JJ, et al. Model transparency and validation: a report of the ISPOR-

- 27 SMDM Modeling Good Research Practices Task Force--7. *Value Health*. 2012;15(6):843-850.
- 28 27. Hosmer DW Jr LS, Sturdivant RX. Applied Logistic Regression, 3rd Edition. New York, NY: Wiley; 2013.

- 1 28. Benjamin E, Recinos G, Aiolfi A, Inaba K, Demetriades D. Pharmacological Thromboembolic Prophylaxis
- in Traumatic Brain Injuries: Low Molecular Weight Heparin Is Superior to Unfractionated Heparin. *Ann Surg.*2017;266(3):463-469.
- Sultan I, Kilic A, Arnaoutakis G, Kilic A. Impact of Foley Catheter Placement by Medical Students on
  Rates of Postoperative Urinary Tract Infection. *J Am Coll Surg.* 2018.
- 6 30. Streiner DL, Norman GR. Mine is bigger than yours: measures of effect size in research. *Chest.*7 2012;141(3):595-598.
- 8 31. Ranganathan P, Pramesh CS, Aggarwal R. Common pitfalls in statistical analysis: Absolute risk reduction,
- 9 relative risk reduction, and number needed to treat. *Perspect Clin Res.* 2016;7(1):51-53.
- 10 32. Schulz KF, Altman DG, Moher D, Group C. CONSORT 2010 statement: updated guidelines for reporting
- 11 parallel group randomised trials. *BMJ*. 2010;340:c332.
- Wu S, Jin Z, Wei X, et al. Misuse of statistical methods in 10 leading Chinese medical journals in 1998 and *2008. ScientificWorldJournal.* 2011;11:2106-2114.

14

# 1 Tables

2

Table 1. Pearson Chi-square Test for Proportion Difference<sup>†</sup>

	Difference of 25.0%	Difference of 12.5%	Difference of 6.25%	
Group 1 proportion	0.5	0.25	0.25	
Group 2 proportion	0.75	0.375	0.3125	
Number of sides	2	2	2	
Null proportion difference	0	0	0	
Computed N per group	58	215	812	
Actual power	0.802	0.801	0.800	
<sup>†</sup> Assuming asymptotic normal distribution, normal approximation method for calculation, $\alpha$ =0.05, and $\beta$ =0.20				
3				
Table 2. Resources for Reviewing Biostatistics and Study Design   4				
Formal Resources				
Statistical software resources				
ASHP Foundation	6			
ACCP Programs	0			
Several books				
Tools for Developing the Thought Process				
Editorials				
Letters to the Editor				
JAMA (Guide to Statistics and Methods)				
BMJ (Endgames)				
Peer review				
Practice and listen				
Ask questions				
<sup>†</sup> See supplemental content for links				

See supplemental content for links