

PITFALLS OF BIOMEDICAL RESEARCH

Presented at the Nigerian Orthopaedic Association Zone 1 Meeting at EKSUTH, April 4th, 2019

Prof Oluwadiya Kehinde

Department of Surgery

Ekiti State university

Ado-Ekiti

www.oluwadiya.com



IMPORTANT

Design, analysis and interpretation
of studies are closely interwoven!

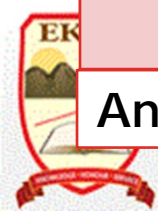


DELUGE...

There is an increasing number of publications on the flaws and errors in much of published medical literature:

- The scandal of poor medical research- **DG Altman, 1994**
- Statistical errors in medical research, a chronic disease?- **J Young, 2007**
- Improved reporting of statistical design and analysis: guidelines, education, and editorial policies. **Mazumdar M et al 2010**
- Bay Area Research symposium 2010- "Why most published research findings are false" - **John Ioannidis**

And many more...



MEDICAL LIES?



PLOS Medicine | www.plosmedicine.org

0696

August 2005 | Volume 2 | Issue 8 | e124

Open access, freely available online

Essay

Why Most Published Research Findings Are False

John P.A. Ioannidis

Summary

factors that influence this problem and is characteristic of the field and can
some populations the confounding may also demonstrate in which the



WHY MOST PUBLISHED RESEARCH FINDINGS ARE FALSE

Summary

There is increasing concern that most current published research findings are false. The probability that a research claim is true may depend on study power and bias, the number of other studies on the same question, and, importantly, the ratio of true to no relationships among the relationships probed in each scientific field. In this framework, a research finding is less likely to be true when the studies conducted in a field are smaller; when effect sizes are smaller; when there is a greater number and lesser preselection of tested relationships; where there is greater flexibility in designs, definitions, outcomes, and analytical modes; when there is greater financial and other interest and prejudice; and when more teams are involved in a scientific field in chase of statistical significance.

Simulations show that for most study designs and settings, it is more likely for a research claim to be false than true.

Moreover, for many current scientific fields, claimed research findings may often be simply accurate measures of the prevailing bias. In this essay, I discuss the implications of these problems for the conduct and interpretation of research.



MEDICAL LIES?



Hall of shame

- 80% of non-randomized studies were wrong
- 25% of supposedly gold-standard randomized trials
- 10% of large randomized trials



EXCERPTS FROM DG ALTMAN

- *"We need less research, better research, and research done for the right reasons"*
- *"Doctors need not be experts in statistics, but they should understand the principles of sound methods of research. If they can also analyze their own data, so much the better. **Amazingly, it is widely considered acceptable for medical researchers to be ignorant of statistics. Many are not ashamed (and some seem proud) to admit that they don't know anything about statistics"***

The scandal of poor medical research- DG Altman, 1994



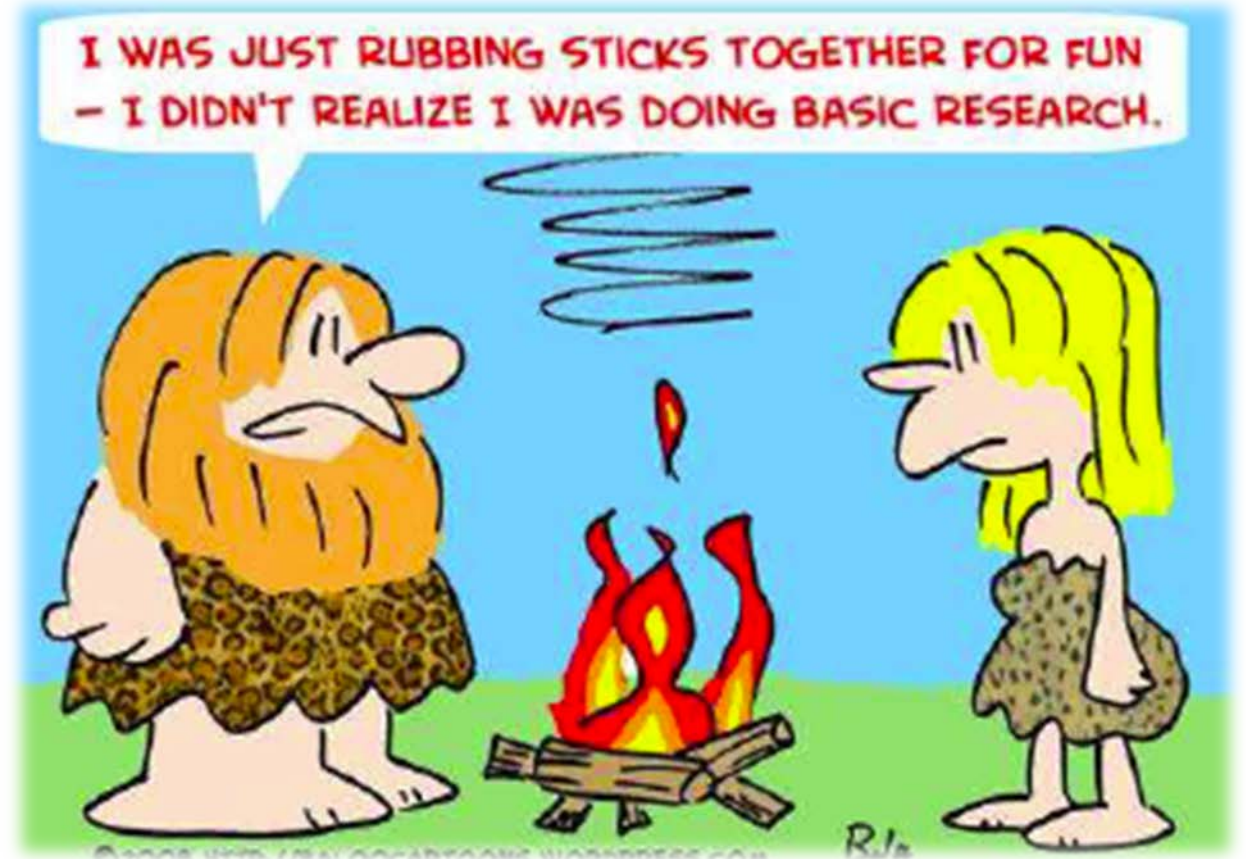
WHY THEN, ARE ERRORS SO COMMON IN MEDICAL RESEARCH?

- Mistakes and errors
- Deliberate fraud
 - Pressure
 - Publish or perish
 - Grants
 - Recognition and fame



REMEMBER.....

- The primary purpose of research is to conduct a scientific, or, scholarly investigation into a phenomenon, or to answer a burning question.
- Research may be defined as a systematic approach to problem solving.



WHAT IS A PITFALL?

- A hidden danger or unsuspected difficulty (Webster)
- A pitfall is a conceptual error into which, ***because of its specious plausibility***, people frequently and easily fall.
- It is "***the taking of a false logical path***" that may lead the unwary to absurd conclusions, a hidden mistake capable of destroying the validity of an entire argument.



STEPS IN PERFORMING RESEARCH

- Research Problem → **What, When**
- Literature Review → **What, When, How, Why**
- Conceptual & Theoretical Frameworks → **What, Why**
- Variables & Hypotheses → **What, How**
- Research Design → **How**
- Population & sample → **Who, What**
- Data Collection → **How**
- Data Analysis → **Why**
- Results and findings

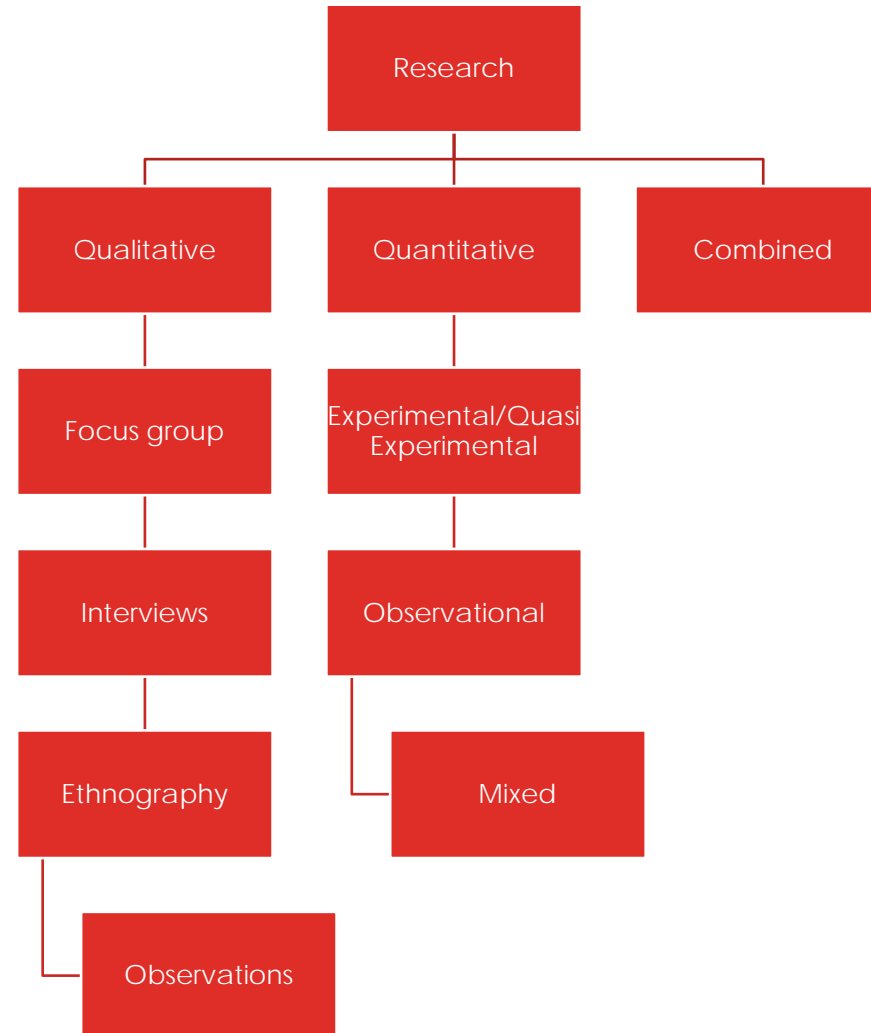


PITFALLS: RESEARCH DESIGN

- Not choosing the right study design
- Not seeking the advice of a statistician on study design
- Not specifying the priori hypotheses
- Investigator Loose Procedure Pitfall
- Not anticipating potential confounders
- Not specifying the randomization and blinding procedures



NOT CHOOSING THE RIGHT STUDY DESIGN



NOT SEEKING STATISTICIAN ADVICE ON STUDY DESIGN

- What measurement levels should be used for each variable?
- Sample size determination
- Need for pilot study?
- Priori hypothesis



NOT SPECIFYING THE PRIORI HYPOTHESES

- Can lead to data dredging or multiple testing with its attendant errors (More on this later)
- A serious potential pitfall is present when investigators collect a large amount of data and have not pre-planned how they are to analyze the data. If an investigator is blessed with a abundance of data ... he can select those data which confirm his hypothesis that a relationship exists.

(Lipset, Trow, & Coleman, 1970; Selvin, 1970)



INVESTIGATOR LOOSE PROCEDURE PITFALL

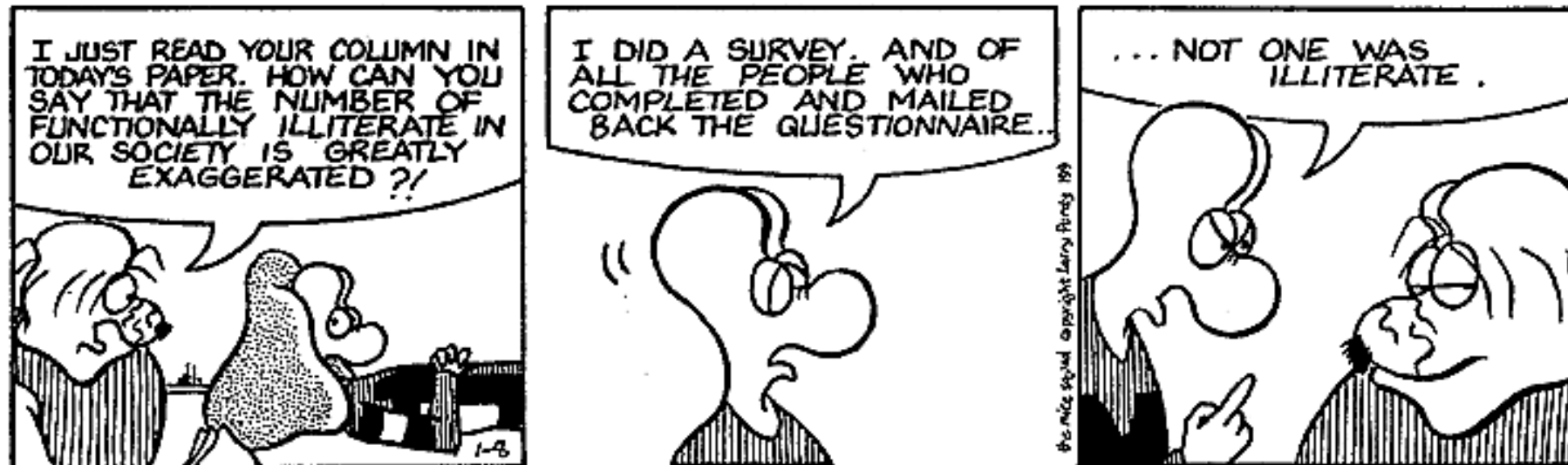
- Not specifying the outcome measures
- Not anticipating potential confounders
- Not specifying the randomization and blinding procedures



SAMPLING PROCEDURE PITFALLS

- Representative sampling is one of the most fundamental tenets of inferential statistics: the observed sample must be representative of the target population in order for inferences to be valid

THE MICE SQUAD



SAMPLING METHODS

Probability Sampling

- Simple random sampling
- Stratified random sampling
- Systematic sampling
- Cluster (area) sampling
- Multistage sampling

Non-Probability Sampling

- Deliberate (quota) sampling
- Convenience (Availability) sampling
- Purposive sampling
- Snowball sampling

SAMPLING PROCEDURE PITFALLS

- Not calculating the correct sample size
 - Small (inadequate) samples
 - Overlarge samples
 - (More on these later)



SAMPLING PROCEDURE PITFALLS

- Using the wrong population
 - Hospital data
 - Employee data
 - Ignoring potential cofounders in the population
 - Wrong sampling method



PITFALLS IN DATA COLLECTION

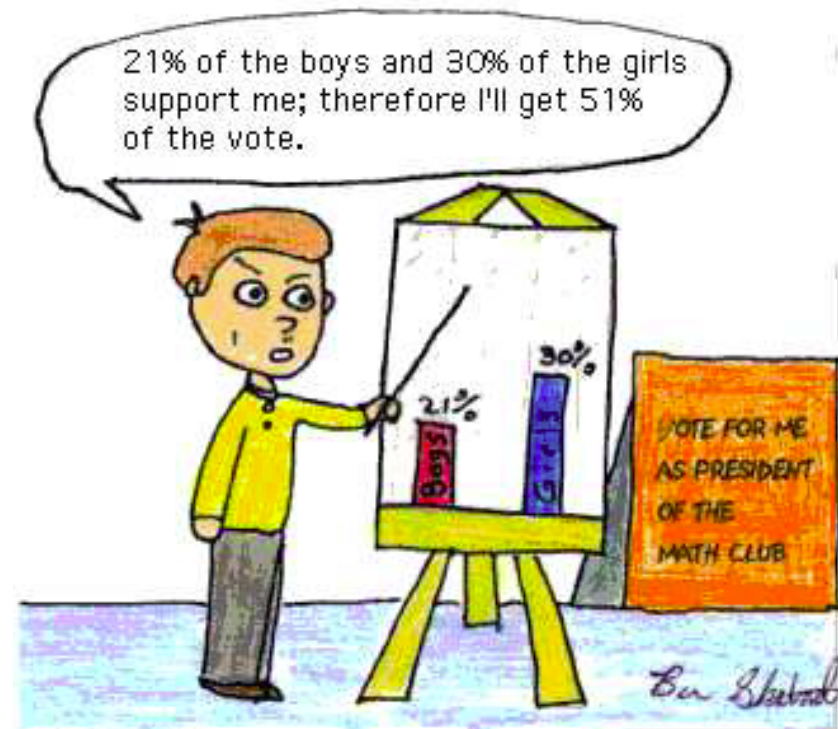
- Failure to follow the procedure as laid down in the methodology (*Experimenter Failure to Follow the Procedure Effect*)
- Poorly trained research personnel
- Using poorly developed questionnaires
- Poor supervisory procedure (*Investigator Loose Procedure Effect*)
- Poor supervision
- Outright fraud



PITFALLS IN DATA ANALYSIS

The Problem of Statistical Errors in research

- Widespread
- Long-standing and becoming more common!
- Potentially serious
- Largely unknown
- Usually concerns basic, not advanced, statistics



PITFALLS IN DATA ANALYSIS

- Investigators at times fail to report that the data did not support their original hypothesis.
- Instead, after they have studied the data, they derive a new hypothesis that is supported by the data and then "verify" the new hypothesis by performing a statistical test on the same data from which it was derived
- *Although investigators may derive a new hypothesis from a completed study, the new hypothesis needs to be tested and verified in a subsequent study.*

(Lipset, Trow, & Coleman, 1970; Selvin, 1970)



PITFALLS IN DATA ANALYSIS

- Investigators at times collect incidental data that are not directly related to the hypotheses they are testing.
- If they fail to confirm their original hypotheses, they then perform a large number of statistical tests on the remaining data and report whatever significant results are obtained as "findings."
- *This can easily lead to misleading conclusions*



PITFALLS IN DATA ANALYSIS

Failure to report negative results.

- Investigators may discard all data of an experiment as bad data if not in agreement with theory, and start over
- The problem here is that if the investigator obtains positive results in a later study and publishes them without mentioning his earlier negative results, the reader is likely to conclude wrongly that the positive results are more stable, more easily replicable, or more valid than is actually the case



PITFALLS IN DATA ANALYSIS

- When an investigator obtains negative results that fail to confirm his hypothesis he is likely to check for computational errors in the data analysis or to run another data analysis
- However, when the original analysis confirms the investigators' hypothesis, it is unlikely that he will check for computational errors or run another analysis

(Friedlander, 1964).



PITFALLS IN DATA ANALYSIS

Using descriptive statistics incorrectly

- Use the mean and standard deviation ONLY to report normally distributed data: "Mean (SD) height was 72 cm (4.3 cm)."
- Use the median and interquartile range to report non normally distributed data: "Median (IQR) length was 9 cm (6 to 25 cm)."



PITFALLS IN DATA ANALYSIS

- **Using descriptive statistics incorrectly**
- The shape of the distribution (normal or skewed) may determine the class of statistical test used to analyze the data ("parametric" or "nonparametric," respectively).
- Most biological data are not normally distributed; the median and IQR should be used in such situations



PITFALLS IN DATA ANALYSIS

Over-emphasis on p-values

- Statistical significance does not guarantee clinical significance.
- **Example:**
 - a study of about 60,000 heart attack patients found that those admitted to the hospital on weekdays had a significantly longer hospital stay than those admitted to the hospital on weekends ($p < .03$), but the magnitude of the difference was too small to be important: 7.4 days (weekday admits) vs. 7.2 days (weekend admits).



PITFALLS IN DATA ANALYSIS

Over-emphasis on p-values

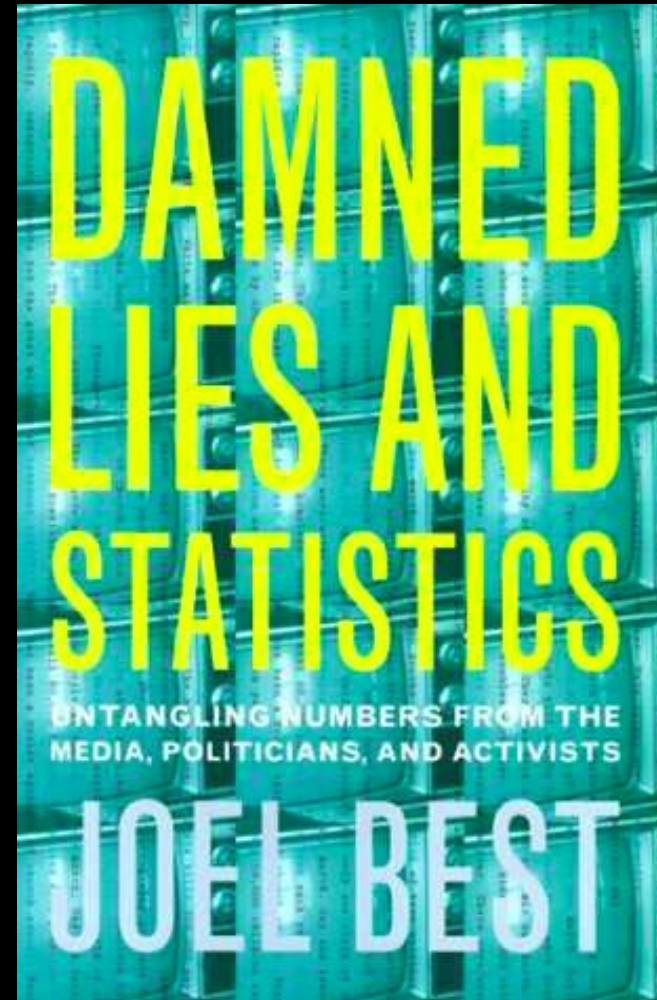
- Statistical significance does not guarantee clinical significance.
 - Clinically unimportant effects may be statistically significant if a study is large (and therefore, has a small standard error and extreme precision).
- Pay attention to effect sizes and confidence intervals



PITFALLS IN DATA ANALYSIS

Over-emphasis on p-values

- Statistical significance does not imply a cause-effect relationship.
- Always interpret results in the context of the study design.



PITFALLS IN DATA ANALYSIS

Myths about significant values

- **Myth 1:** "If a result is not significant, it proves there is no effect."
- **Myth 2:** "The obtained significance level indicates the reliability of the research finding."
- **Myth 3:** "The significance level tells you how big or important an effect is."
- **Myth 4:** "If an effect is statistically significant, it must be clinically significant."



PITFALLS IN DATA ANALYSIS

Problems with relying on p -values

- When sample size is low, p is usually not insignificant
 - ☒ if effect size is small, try bigger sample
- When sample size is very big, p can easily be very small even for tiny effects
 - ☒ e.g., mean IQ of men is 0.8pts higher than IQ of women, in a sample of 10,000: statistically significant, but is it clinically significant?
- When many tests are run, one of them is bound to turn up significant by random chance
 - ☒ multiple comparisons → inflated Type-I errors



PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging

- Performance of 2 or more related hypothesis tests using the same data set
- **Example 2:**
 - Suppose we consider the safety of a drug in terms of the occurrences of different types of side effects. As more types of side effects are considered, it becomes more likely that the new drug will appear to be less safe than existing drugs in terms of at least one side effect.



PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging

- For a single test, with significant level at 0.05 means that there is only a 5 percent chance that it is a spurious finding resulting solely from chance variations.
- For two tests: the probability that at least one such analysis will yield a spurious, significant finding is greater than 5 percent.
- To determine the new probability level:

The probability that a significant result would be obtained in either of the two tests = $.95 \times .95 = 0.902$

Subtract this from 1.

$1 - .902 = .098$.

That is 9.8%



PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging

- For 10 tests:
 - The probability that a significant result will be obtained in any of the ten tests is $(0.95)^{10}$
=0.59
The new probability $1-0.59$
=0.41 i.e. 41%

Formula is:

$$1-(1-n)^x$$

n=Significant level

x=number of independent tests



PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging: Real life example

- In 1980, researchers at Duke randomized 1073 heart disease patients into two groups, **but treated the groups equally**.
- Not surprisingly, there was no difference in survival.
- Then they divided the patients into 18 subgroups based on prognostic factors.
- In a subgroup of 397 patients (with three-vessel disease and an abnormal left ventricular contraction) survival of those in "group 1" was significantly different from survival of those in "group 2" ($p < .025$).
- ***How could this be since the "treatment" was equal for all the groups?***

(Lee et al. "Clinical judgment and statistics: lessons from a simulated randomized trial in coronary artery disease," *Circulation*, 61: 508-515, 1980.)



PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging: Real life example

- If we compare survival of “treatment” and “control” within each of 18 subgroups, that’s 18 comparisons.
- If these comparisons were independent, the chance of at least one false positive would be...

Formula is:

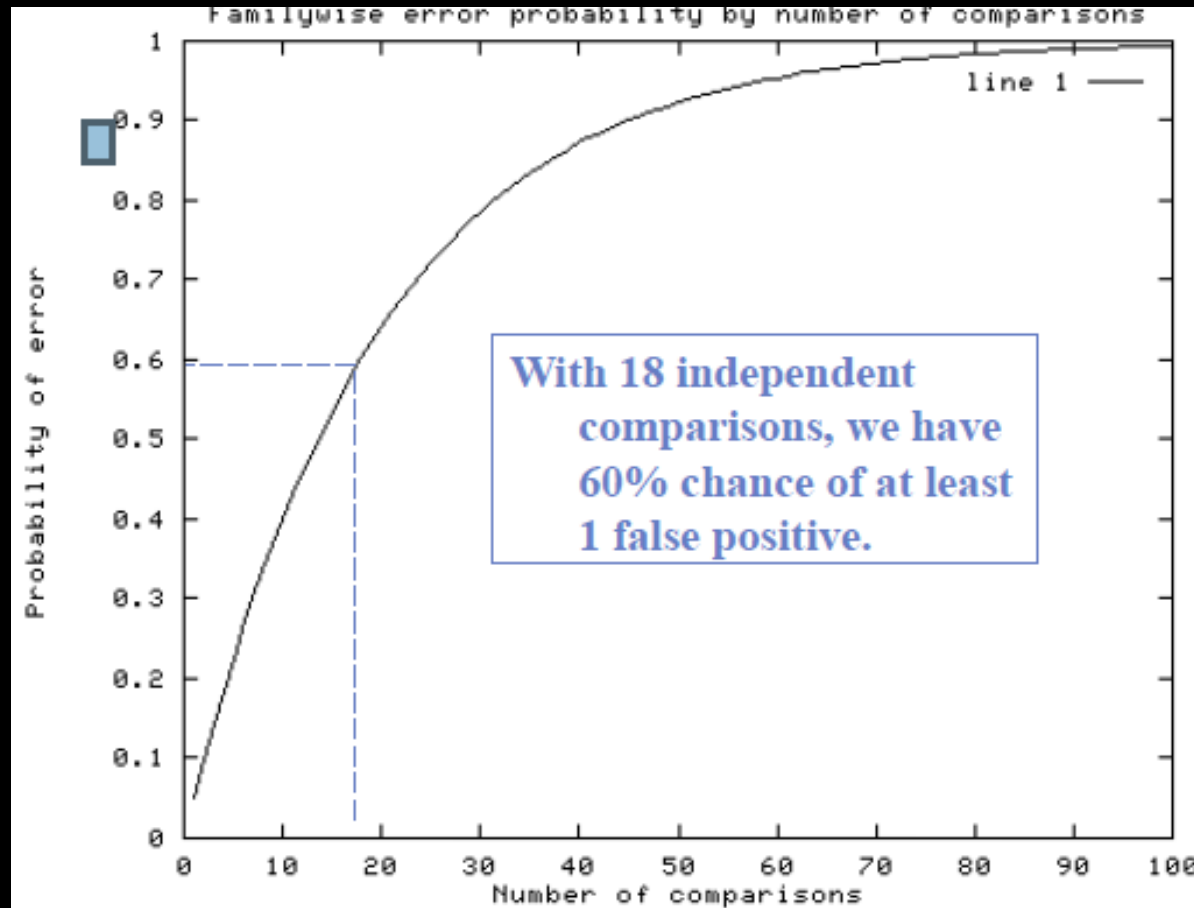
$$\begin{aligned} &1 - (.95)^{18} \\ &= 0.60 \\ &60\% \end{aligned}$$

(Lee et al. “Clinical judgment and statistics: lessons from a simulated randomized trial in coronary artery disease,” *Circulation*, 61: 508-515, 1980.)



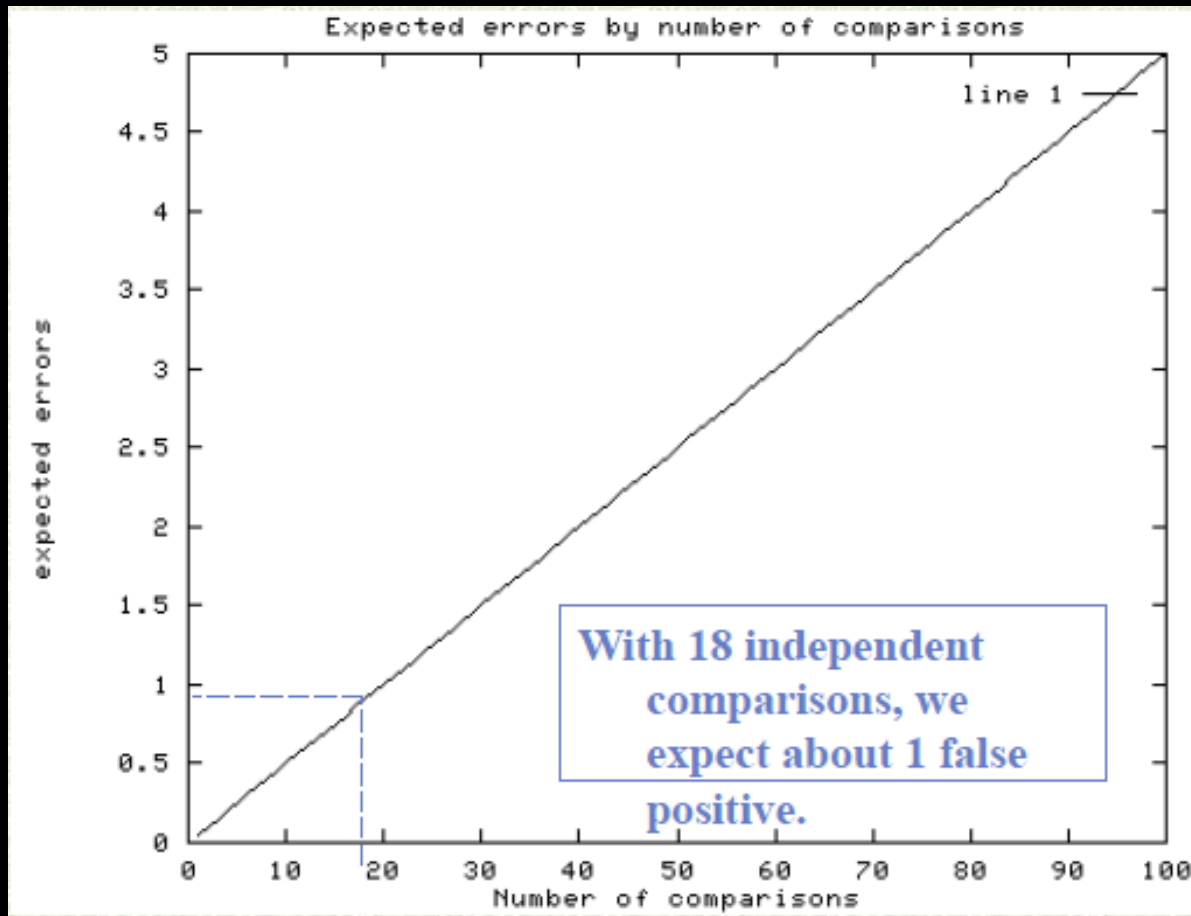
PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging: Real life example



PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging: Real life example



$.05 \cdot k$ significant p-values ($p < .05$) are expected to arise just by chance, where k is the number of tests run.



PITFALLS IN DATA ANALYSIS

Multiple Testing or Data Dredging: EXAMPLES

<u>Source</u>	<u>Example</u>
Multiple outcomes	a cohort study looking at the incidence of breast cancer, colon cancer, and lung cancer
Multiple predictors	an observational study with 40 dietary predictors or a trial with 4 randomization groups
Subgroup analyses	a randomized trial that tests the efficacy of an intervention in 20 subgroups based on prognostic factors
Multiple definitions for the exposures and outcomes	an observational study where the data analyst tests multiple different definitions for “moderate drinking” (e.g., 5 drinks per week, 1 drink per day, 1-2 drinks per day, etc.)
Multiple time points for the outcome (repeated measures)	a study where a walking test is administered at 1 months, 3 months, 6 months, and 1 year
Multiple looks at the data during sequential interim monitoring	a 2-year randomized trial where the efficacy of the treatment is evaluated by a Data Safety and Monitoring Board at 6 months, 1 year, and 18 months



PITFALLS IN DATA ANALYSIS

Hypothetical example:

- Researchers wanted to compare nutrient intakes between women who had fractures and women who had no fractures.
- They used a food-frequency questionnaire and a food diary to capture food intake.
- From these two instruments, they calculated daily intakes of all the vitamins, minerals, macronutrients, antioxidants, etc.
- Then they compared fracturers to non-fracturers on all nutrients from both questionnaires.
- They found a statistically significant difference in vitamin K between the two groups ($p < .05$).
- They had a lovely explanation of the role of vitamin K in injury repair, bone, clotting, etc.



PITFALLS IN DATA ANALYSIS

Hypothetical example:

- What's going on?
- Findings are almost certainly artifactual (false positive!).....



PITFALLS IN DATA ANALYSIS

Factors indicative of chance findings:

1. Analysis are explanatory	The authors have mined the data for associations rather than testing a limited number of a priori hypotheses.
2. Many tests have been performed, but only a few p-values are “significant”.	If there are no associations present, $.05 \cdot k$ significant p-values ($p < .05$) are expected to arise just by chance, where k is the number of tests run.
3. The “significant” p-values are modest in size.	The closer a p-value is to .05, the more likely it is a chance finding. According to one estimate*, about 1 in 2 p-values $< .05$ is a false positive, 1 in 6 p-values $< .01$ is a false positive, and 1 in 56 p-values $< .0001$ is a false positive.
4. The p-values are not adjusted for multiple comparisons	Adjustment for multiple comparisons can help control the study-wide false positive rate.



ERRORS

- **Type I error**

- Claiming a difference between two samples when in fact there is none.
- Also called an α error.
- Typically 0.05 is used.
- Seen typically in multiple testing



ERRORS

- **Type II error**

- Claiming there is no difference between two samples when in fact there is.
- Also called β error.
- The probability of not making a Type II error is $1 - \beta$, which is also called the *power* of the test.
- It usually cannot be detected without a proper power analysis



ERRORS

		Test Result	
		No difference H_0	Shows Difference H_1
Truth	No difference H_0	No Error	Type I α
	Shows Difference H_1	Type II β	No Error



HOW TO INCREASE THE POWER

i.e. reduce type II error

- Increase the sample size
- Reduce variation between measurements
- The effect of intervention should be stronger



SAMPLE SIZE CALCULATION

- **Also called “*power analysis*”.**
 - When designing a study, one needs to determine how large a sample is needed.
 - Power is the ability of a study to avoid a Type II error.
 - Sample size calculation yields the number of study subjects needed, given a certain desired power to detect a difference and a certain level of P value that will be considered significant.
 - Many studies are completed without proper estimate of appropriate study size.
 - This may lead to an erroneous “negative” study outcome.



SAMPLE SIZE CALCULATION

Depends on:

- Level of Type I error: 0.05 typical
- Level of Type II error: 0.20 typical
- One sided vs. two sided: nearly always two-sided
- Inherent variability of population
 - Usually estimated from preliminary data or previous study
- The difference that would be meaningful between the two assessment arms.



Download from:
goo.gl/zgC2pA

MODULE 2

SAMPLE SIZE CONSIDERATIONS

Prof Oluwadiya K S

www.oluwadiya.com

College of Medicine,

Ekiti State University, Ado-Ekiti, Nigeria

A sample is a percentage of the total population in statistics. And the aim is to use the data from a sample to make inferences about the population as a whole. Determining sample size is one of the most ticklish problems in research. What should be the size "n" of the sample. If the sample size is too small, it may not achieve the objectives of the study, and if it is too large, we may incur huge costs and waste resources. Ethics, economics, time and other constraints dictates that we keep our sample size optimal. In essence, the sample must be just the right size: neither too large or too small.

What are the factors that may affect the size of a sample?

The homogeneousness of the population: This is the uniformity of the population in relation to the variable of interest. If the population is homogeneous, a small sample will serve the purpose. But if heterogenous, then a large sample would be required.

- i. **The type of study:** If a survey is to be carried out for descriptive purposes (e.g. prevalence of a disease or characteristics), the sample size is based on the required precision of the prevalence estimate and margin of error. On the other hand, if it is comparative study, the sample size is based on how detailed a comparison is desired and power calculations.
- ii. **Grouping:** For adequate statistical analysis, each group in the study must have sufficient participants. Therefore, the larger the number of groups in a study, the larger the overall sample size of the study.
- iii. **Level of precision:** If a high level of precision is required, then a relatively larger sample would be required.
- iv. **The effect size:** Detecting very small differences requires larger samples than detecting large differences.
- v. **Analytical technique:** More sophisticated statistic techniques require larger sample sizes than simple techniques. That is why if you are using many statistic methods for your study, you must calculate the sample size for each technique. You will use the largest calculated sample size for the study. For example, if you propose to find difference between group means, chi square as well as Multiple regression. You will calculate the sample size requirements for all three statistics and use the largest as the sample size for your study.



FINAL THOUGHTS.....

- When reading a journal article.....



KEEP IN MIND THAT....

- No study is perfect
- All data is dirty is some way or another; research is what you do with that dirty data
- Measurement involves making choices



BE CRITICAL ABOUT NUMBERS

- Every statistic is a way of summarizing complex information into relatively simple numbers.
- How did the researchers arrive at these numbers?
- Who produced the numbers and what is their bias?
- How were key terms be defined & in how many different ways?



BE CRITICAL ABOUT NUMBERS

- How was the choice for the measurement made?
- What type of sample was gathered & how does that affect result?
- Is the statistical result interpreted correctly?
- If comparisons are made, were they appropriate?
- Are there competing statistics?

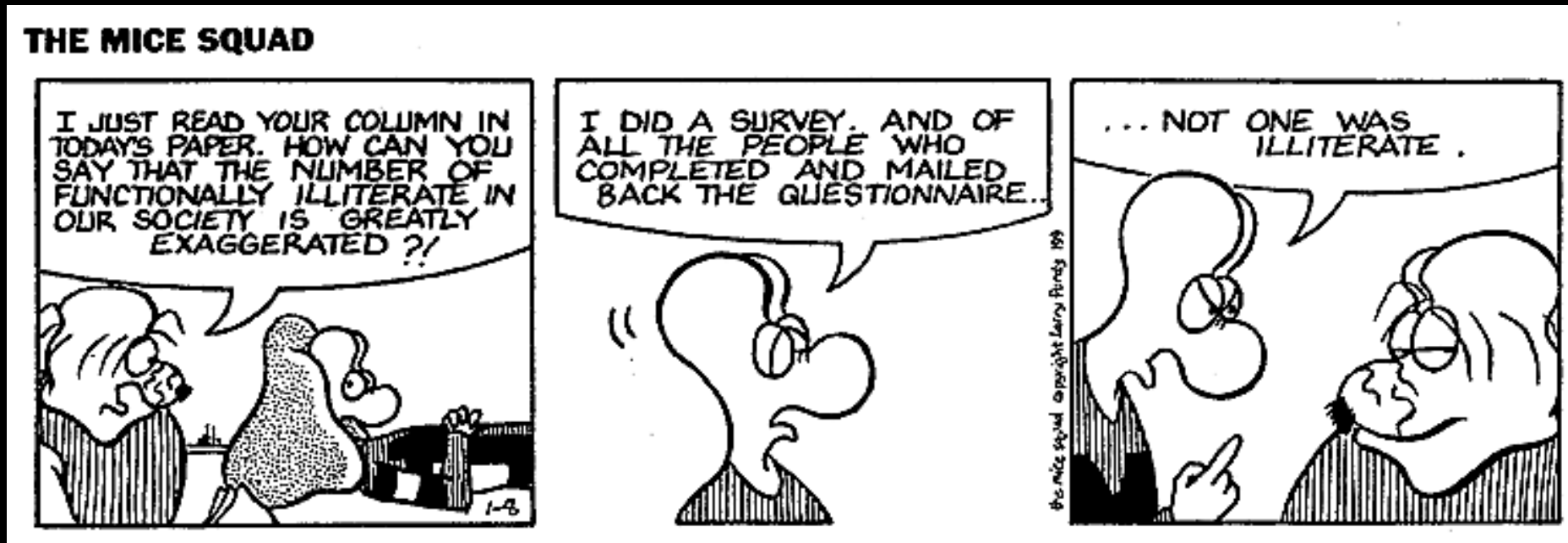


BE CRITICAL ABOUT NUMBERS

With one foot in a bucket of ice water, and one foot in a bucket of boiling water, you are, on the average, comfortable.



BE CRITICAL ABOUT NUMBERS: BIAS AND ERROR



Thanks for your attention
To ask questions, please
join the forum at
www.oluwadiya.com

