# P-values, confidence intervals, and the Bayesian alternative

Steve Simon

P.Mean Consulting

www.pmean.com

# 2. Abstract

- P-values and confidence intervals are the fundamental tools used in most inferential data analyses. They are possibly the most commonly reported statistics in the medical literature. Unfortunately, both p-values and confidence intervals are subject to frequent misinterpretations.

# 3. Abstract

- In this two hour webinar, you will learn the proper interpretation of p-values and confidence intervals, and the common abuses and misconceptions about these statistics. You will also see a simple application of Bayesian analysis which provides an alternative to p-values and confidence intervals.

# 4. Learning objectives

- In this seminar, you will learn how to:
  - distinguish between statistical significance and clinical significance;
  - define and interpret p-values;
  - explain the ethical issues associated with inadequate sample sizes.
  - explain the difference between informative and diffuse priors;
  - interpret statistics from a posterior distribution.

# 5. Outline

1. Icebreaker
2. Pop quiz
3. Definitions
4. What is a p-value?
5. Practice exercises
6. What is a confidence interval?
7. Practice exercises
8. A simple example of Bayesian data analysis.
9. Repeat of pop quiz

# 6. Icebreaker

- What is the most complicated statistic you have ever had to compute?
    1. percentage
    2. mean
    3. standard deviation
    4. t-test
    5. correlation coefficient
    6. linear regression model
    7. survival curve
    8. logistic regression model
    9. other

# 7. Pop quiz

A research paper computes a p-value of 0.45. How would you interpret this p-value?

1. Strong evidence for the null hypothesis
2. Strong evidence for the alternative hypothesis
3. Little or no evidence for the null hypothesis
4. Little or no evidence for the alternative hypothesis
5. More than one answer above is correct.
6. I do not know the answer.

# 8. Pop quiz

A research paper computes a confidence interval for a relative risk of 0.82 to 3.94. What does this confidence interval tell you.

1. The result is statistically significant and clinically important.
2. The result is not statistically significant, but is clinically important.
3. The result is statistically significant, but not clinically important.
4. The result is not statistically significant, and not clinically important.
5. The result is ambiguous.
6. I do not know the answer.

# 9. Pop quiz

A Bayesian data analysis can incorporate subjective opinions through the use of

1. Bayes rule.
2. data shrinkage.
3. a prior distribution.
4. a posterior distribution.
5. p-values.
6. I do not know the answer.

# 10. Definitions: Population

- A population is a collection of items of interest in research. The population represents a group that you wish to generalize your research to. Populations are often **defined in terms of demography, geography, occupation, time, care requirements, diagnosis, or some combination of the above**.

# 11. Definitions: Population

- A population is a collection of items of interest in research. The population represents a group that you wish to generalize your research to. Populations are often **defined in terms of demography, geography, occupation, time, care requirements, diagnosis, or some combination of the above**.

# 12. Definitions: Population

- An example of a population would be all infants born in the state of Missouri during the 1995 calendar year who have one or more visits to the Emergency room during their first year of life.

# 13. Definitions: Sample

- A sample is a subset of a population. A random sample is a subset where every item in the population has the same probability of being in the sample. **Usually, the size of the sample is much less than the size of the population**. The primary goal of much research is to use information collected from a sample to try to characterize a certain population.

# 14. Definitions: Type I Error

- In your research, you specify a null hypothesis (typically labeled H0) and an alternative hypothesis (typically labeled Ha, or sometimes H1). By tradition, the null hypothesis corresponds to no change. A Type I error is **rejecting the null hypothesis when the null hypothesis is true**.

# 15. Definitions: Type I Error

- **Example:** Consider a new drug that we will put on the market if we can show that it is better than a placebo. In this context, H0 would represent the hypothesis that the average improvement (or perhaps the probability of improvement) among all patients taking the new drug is equal to the average improvement (probability of improvement) among all patients taking the placebo. A Type I error would be **allowing an ineffective drug onto the market**.

# 16. Definitions: Type II Error

- A Type II error is **accepting the null hypothesis when the null hypothesis is false. Many studies have small sample sizes that make it difficult to reject the null hypothesis,** even when there is a big change in the data. In these situations, a Type II error might be a possible explanation for the negative study results.

# 17. Definitions: Type II Error

- **Example:** Consider a new drug that we will put on the market if we can show that it is better than a placebo. In this context, H0 would represent the hypothesis that the average improvement (or perhaps the probability of improvement) among all patients taking the new drug is equal to the average improvement (probability of improvement) among all patients taking the placebo. A Type II error would be **keeping an effective drug off the market**.

# 18. What is a p-value?

- A p-value is a **measure of how much evidence we have against the null hypothesis**. The null hypothesis, traditionally represented by the symbol H0, represents the hypothesis of no change or no effect. The smaller the p-value, the more evidence we have against H0.

# 19. What is a p-value?

- The p-value is also a measure of how likely we are to get a certain sample result or a result "more extreme," assuming H0 is true. The type of hypothesis (right tailed, left tailed or two tailed) will determine what "more extreme" means.

# 20. What is a p-value?

- The p-value is also a measure of how likely we are to get a certain sample result or a result "more extreme," assuming H0 is true. The type of hypothesis (right tailed, left tailed or two tailed) will determine what "more extreme" means.

# 21. What is a p-value?

- It is easiest to understand the p-value in a data set that is already at an extreme. Suppose that a drug company alleges that **only 50% of all patients who take a certain drug will have an adverse event of some kind**. You believe that the adverse event rate is much higher. **In a sample of 12 patients, all twelve have an adverse event**.

# 22. What is a p-value?

- **The data supports your belief because it is inconsistent with the assumption of a 50% adverse event rate**. It would be like flipping a coin 12 times and getting heads each time.

# 23. What is a p-value?

- The p-value, the probability of getting a sample result of **12 adverse events in 12 patients** assuming that the adverse event rate is 50%, is a measure of this inconsistency. **The p-value, 0.000244, is small enough that we would reject the hypothesis that the adverse event rate was only 50%.**

# 24. What is a p-value?

**A large p-value should not automatically be construed as evidence in support of the null hypothesis**. Perhaps the failure to reject the null hypothesis was caused by an inadequate sample size. When you see a large p-value in a research study, you should also look for one of two things:

1. a **power calculation** that confirms that the sample size in that study was adequate for detecting a clinically relevant difference; and/or

2. a **confidence interval** that lies entirely within the range of clinical indifference.

# 25. What is a p-value?

You should also be cautious about a small p-value, but for different reasons. **In some situations, the sample size is so large that even differences that are trivial from a medical perspective can still achieve statistical significance.**

# 26. Practice exercise: interpret the p-values in the following abstract.

**1. The Outcome of Extubation Failure in a Community Hospital Intensive Care Unit: A Cohort Study**. Seymour CW, Martinez A, Christie JD, Fuchs BD. *Critical Care 2004*, 8:R322-R327 (20 July 2004) **Introduction:** Extubation failure has been associated with poor intensive care unit (ICU) and hospital outcomes in tertiary care medical centers. Given the large proportion of critical care delivered in the community setting, our purpose was to determine the impact of extubation failure on patient outcomes in a community hospital ICU. **Methods:** A retrospective cohort study was performed using data gathered in a 16-bed medical/surgical ICU in a community hospital. During 30 months, all patients with acute respiratory failure admitted to the ICU were included in the source population if they were mechanically ventilated by endotracheal tube for more than 12 hours. Extubation failure was defined as reinstitution of mechanical ventilation within 72 hours (n = 60), and the control cohort included patients who were successfully extubated at 72 hours (n = 93). **Results:** The primary outcome was total ICU length of stay after the initial extubation. Secondary outcomes were total hospital length of stay after the initial extubation, ICU mortality, hospital mortality, and total hospital cost. Patient groups were similar in terms of age, sex, and severity of illness, as assessed using admission Acute Physiology and Chronic Health Evaluation II score ($P > 0.05$). Both ICU (1.0 versus 10 days; $P < 0.01$) and hospital length of stay (6.0 versus 17 days; $P < 0.01$) after initial extubation were significantly longer in reintubated patients. ICU mortality was significantly higher in patients who failed extubation (odds ratio = 12.2, 95% confidence interval [CI] = 1.5–101; $P < 0.05$), but there was no significant difference in hospital mortality (odds ratio = 2.1, 95% CI = 0.8–5.4; $P < 0.15$). Total hospital costs (estimated from direct and indirect charges) were significantly increased by a mean of US$33,926 (95% CI = US$22,573–45,280; $P < 0.01$). **Conclusion:** Extubation failure in a community hospital is univariately associated with prolonged inpatient care and significantly increased cost. Corroborating data from tertiary care centers, these adverse outcomes highlight the importance of accurate predictors of extubation outcome.

# 27. Practice exercise: interpret the p-values in the following abstract.

**2. Elevated White Cell Count in Acute Coronary Syndromes: Relationship to Variants in Inflammatory and Thrombotic Genes**. Byrne CE, Fitzgerald A, Cannon CP, Fitzgerald DJ, Shields DC. *BMC Medical Genetics 2004*, 5:13 (1 June 2004)
**Background:** Elevated white blood cell counts (WBC) in acute coronary syndromes (ACS) increase the risk of recurrent events, but it is not known if this is exacerbated by pro-inflammatory factors. We sought to identify whether pro-inflammatory genetic variants contributed to alterations in WBC and C-reactive protein (CRP) in an ACS population. **Methods:** WBC and genotype of interleukin 6 (IL-6 G-174C) and of interleukin-1 receptor antagonist (IL1RN intronic repeat polymorphism) were investigated in 732 Caucasian patients with ACS in the OPUS-TIMI-16 trial. Samples for measurement of WBC and inflammatory factors were taken at baseline, i.e. Within 72 hours of an acute myocardial infarction or an unstable angina event.
**Results:** An increased white blood cell count (WBC) was associated with an increased C-reactive protein (r = 0.23, $p < 0.001$) and there was also a positive correlation between levels of β-fibrinogen and C-reactive protein (r = 0.42, $p < 0.0001$). IL1RN and IL6 genotypes had no significant impact upon WBC. The difference in median WBC between the two homozygote IL6 genotypes was 0.21/mm3 (95% CI = -0.41, 0.77), and -0.03/mm3 (95% CI = -0.55, 0.86) for IL1RN. Moreover, the composite endpoint was not significantly affected by an interaction between WBC and the IL1 ($p = 0.61$) or IL6 ($p = 0.48$) genotype. **Conclusions:** Cytokine pro-inflammatory genetic variants do not influence the increased inflammatory profile of ACS patients.

# 28. Practice exercise: interpret the p-values in the following abstract.

**3. Is There a Clinically Significant Gender Bias in Post-Myocardial Infarction Pharmacological Management in the Older (>60) Population of a Primary Care Practice?** Di Cecco R, Patel U, Upshur REG. *BMC Family Practice 2002*, 3:8 (3 May 2002) **Background:** Differences in the management of coronary artery disease between men and women have been reported in the literature. There are few studies of potential inequalities of treatment that arise from a primary care context. This study investigated the existence of such inequalities in the medical management of post myocardial infarction in older patients. **Methods:** A comprehensive chart audit was conducted of 142 men and 81 women in an academic primary care practice. Variables were extracted on demographic variables, cardiovascular risk factors, medical and non-medical management of myocardial infarction. **Results:** Women were older than men. The groups were comparable in terms of cardiac risk factors. A statistically significant difference (14.6%: 95% CI 0.048–28.7 $p = 0.047$) was found between men and women for the prescription of lipid lowering medications. 25.3% ($p = 0.0005$, CI 11.45, 39.65) more men than women had undergone angiography, and 14.4 % ($p = 0.029$, CI 2.2, 26.6) more men than women had undergone coronary artery bypass graft surgery. **Conclusion:** Women are less likely than men to receive lipid-lowering medication which may indicate less aggressive secondary prevention in the primary care setting.

# 29. **What is a confidence interval?**

- We statisticians have a habit of **hedging our bets**. We always insert qualifiers into our reports, warn about all sorts of assumptions, and never admit to anything more extreme than probable. There's a famous saying: **"Statistics means never having to say you're certain."**

# 30. **What is a confidence interval?**

- We qualify our statements, of course, because we are always **dealing with imperfect information**. In particular, we are often asked to make statements about a population (a large group of subjects) using information from a sample (a small, but carefully selected subset of this population). No matter how carefully this sample is selected to be a fair and unbiased representation of the population, **relying on information from a sample will always lead to some level of uncertainty**.

# 31. What is a confidence interval?

- **A confidence interval is a range of values that tries to quantify uncertainty associated with the sampling process.** Consider it as a **range of plausible values**.

# 32. What is a confidence interval?

- A wide confidence interval implies poor precision; we can only specify plausible values to a broad and uninformative range. A narrow confidence interval implies good precision; we can place sharp limits on the range of plausible values.

# 33. What is a confidence interval?

- Consider a recent study of **homoeopathic treatment of pain and swelling after oral surgery** (Lokken 1995). When examining swelling 3 days after the operation, they showed that **homoeopathy led to 1 mm less swelling on average**. The **95% confidence interval, however, ranged from -5.5 to 7.5 mm**. This interval implies that **neither a large improvement due to homoeopathy nor a large decrement could be ruled out**.

# 34. **What is a confidence interval?**

- When you see a confidence interval in a published medical report, you should look for two things. First, **does the interval contain a value that implies no change or no effect**? For example, with a confidence interval for a difference look to see whether that interval includes zero. With a confidence interval for a ratio, look to see whether that interval contains one.

# 35. **What is a confidence interval?**

- Here's an example of a confidence interval that contains the null value. **The interval shown below implies no statistically significant change**.

Null
value

# 36. What is a confidence interval?

- Here's an example of a confidence interval that excludes the null value. If we assume that larger implies better, then **the interval shown below would imply a statistically significant improvement**.

Null
value
|                    ←——————→

# 37. **What is a confidence interval?**

- Here's a different example of a confidence interval that excludes the null value. **The interval shown below implies a statistically significant decline**.

Null
value

# 38. What is a confidence interval?

- You should also see **whether the confidence interval lies partly or entirely within a range of clinical indifference**. Clinical indifference represents values of such a trivial size that you would not want to change your current practice.

# 39. What is a confidence interval?

- **If a confidence interval is contained entirely within your range of clinical indifference**, then you have clear and convincing evidence to **keep doing things the same way**.

# 40. What is a confidence interval?

- One the other hand, **if part of the confidence interval lies outside the range of clinical indifference**, then you should consider the possibility that **the sample size is too small**.

# 41. **What is a confidence interval?**

- If your **confidence interval excludes the null value but still lies entirely within the range of clinical indifference**, then you have a result with **statistical significance, but no practical significance**.

# 42. **What is a confidence interval?**

- Finally, if your **confidence interval excludes the null value and lies outside the range of clinical indifference**, then you have **both statistical and practical significance**.

Null
value

Range of clinical indifference

# 43. **Practice exercises**

- Read the abstracts presented above. Interpret the confidence intervals presented in these abstracts.

# 44. **Bayesian data analysis**

- There's a wonderful example of Bayesian data analysis at work that is simple and fun. It's taken directly from an article by Jim Albert in the Journal of Statistics Education (1995, vol. 3 no. 3) which is available on the web at

  – [www.amstat.org/publications/jse/v3n3/albert.html](www.amstat.org/publications/jse/v3n3/albert.html).

# 45. **Bayesian data analysis**

- I want to use his second example, involving a comparison of ECMO to conventional therapy in the treatment of babies with severe respiratory failure. In this study, 28 of 29 babies assigned to ECMO survived and 6 of 10 babies assigned to conventional therapy survived. Refer to the Albert article for the source of the original data.

# 46. **Bayesian data analysis**

- Wikipedia gives a nice general introduction to the concept of Bayesian data analysis with the following formula:
  - $P(H|E) = P(E|H) \, P(H) \, / \, P(E)$
    - where H represents a particular hypothesis, and E represents evidence (data). P, of course, stands for probability.

# 47. **Bayesian data analysis**

- The first step is to specify P(H), which is called the prior probability. Specifying the prior probability is probably the one aspect of Bayesian data analysis that causes the most controversy. The prior probability represents the degree of belief that you have in a particular hypothesis prior to collection of your data.

# 48. **Bayesian data analysis**

- The prior distribution can incorporate data from previous related studies or it can incorporate subjective impressions of the researcher. What!?! you're saying right now. Aren't statistics supposed to remove the need for subjective opinions?

  Actually, a bit of subjectivity is a good thing.

# 49. **Bayesian data analysis**

- First, it is impossible to totally remove subjective opinion from a data analysis. You can't do research without adopting some informal rules. These rules may be reasonable, they may be supported to some extent by empirical data, but they are still applied in a largely subjective fashion.

# 50. **Bayesian data analysis**

Here are some of the subjective beliefs that I use in my work:

1. you should always prefer a simple model to a complex model if both predict the data with the same level of precision.

2. you should be cautious about any subgroup finding that was not pre-specified in the research protocol.

3. if you can find a plausible biological mechanism, that adds credibility to your results.

# 51. **Bayesian data analysis**

- Second, the use of a range of prior distributions can help resolve controversies involving conflicting beliefs. For example, an important research question is whether a research finding should "close the book" to further research. If data indicates a negative result, and this result is negative even using an optimistic prior probability, then all researchers, even those with the most optimistic hopes for the therapy, should move on.

# 52. **Bayesian data analysis**

- Third, while Bayesian data analysis allows you to incorporate subjective opinions into your prior probability, it does not require you to incorporate subjectivity. Many Bayesian data analyses use what it called a diffuse or non-informative prior distribution. This is a prior distribution that is neither optimistic nor pessimistic, but spreads the probability more or less evenly across all hypotheses.

# 53. **Bayesian data analysis**

- Here's a simple example of a diffuse prior that Dr. Albert used for the ECMO versus conventional therapy example. Let's assume that the true survival rate could be either 0, 10%, 20%, ..., 100% in the ECMO group and similarly for the conventional therapy group. This is not an optimal assumption, but it isn't terrible either, and it allows us to see some of the calculations in action.

# 54. **Bayesian data analysis**

- With 11 probabilities for ECMO and 11 probabilities for conventional therapy, we have 121 possible combinations. How should we arrange those probabilities? One possibility is to assign half of the total probability to combinations where the probabilities are the same for ECMO and conventional therapy and the remaining half to combinations where the probabilities are different. Split each of these probabilities evenly over all the combinations.

# 55. **Bayesian data analysis**

- To simplify the display, we multiplied each probability by 1000 and rounded.

| | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 45 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 0.1 | 5 | 45 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 0.2 | 5 | 5 | 45 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 0.3 | 5 | 5 | 5 | 45 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 0.4 | 5 | 5 | 5 | 5 | 45 | 5 | 5 | 5 | 5 | 5 | 5 |
| 0.5 | 5 | 5 | 5 | 5 | 5 | 45 | 5 | 5 | 5 | 5 | 5 |
| 0.6 | 5 | 5 | 5 | 5 | 5 | 5 | 45 | 5 | 5 | 5 | 5 |
| 0.7 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 45 | 5 | 5 | 5 |
| 0.8 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 45 | 5 | 5 |
| 0.9 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 45 | 5 |
| 1 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 45 |

# 56. **Bayesian data analysis**

- The second step in a Bayesian data analysis is to calculate P(E | H).

| | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1 |
|-----|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0.2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 0.3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 0 |
| 0.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 17 | 0 |
| 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 31 | 0 |
| 0.6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 38 | 0 |
| 0.7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 30 | 0 |
| 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 13 | 0 |
| 0.9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# 57. **Bayesian data analysis**

- P(E | H) is the probability of the observed data under each hypothesis. There are formulas for this, using the binomial distribution. I used Excel to calculate these probabilities for me. Here's an example of the formula I used.

```
=binomdist(28,29,0.9,FALSE)*binomdist(6,10,0.6,FALSE)
```

# 58. **Bayesian data analysis**

- Now multiply the prior probability of each hypothesis by the likelihood of the data

|     | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1 |
|-----|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| 0   | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 |
| 0.1 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 |
| 0.2 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 4   | 0 |
| 0.3 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 2   | 25  | 0 |
| 0.4 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 6   | 77  | 0 |
| 0.5 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 10  | 141 | 0 |
| 0.6 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 1   | 13  | 173 | 0 |
| 0.7 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 4   | 10  | 138 | 0 |
| 0.8 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 45  | 67  | 0 |
| 0.9 | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 1   | 77  | 0 |
| 1   | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 |

# 59. **Bayesian data analysis**

- These numbers do not add up to 1, so we need to rescale them.

|  | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1 |
|-----|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0.1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0.2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 |
| 0.3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 32 | 0 |
| 0.4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 97 | 0 |
| 0.5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 13 | 178 | 0 |
| 0.6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 16 | 218 | 0 |
| 0.7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 13 | 174 | 0 |
| 0.8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 57 | 84 | 0 |
| 0.9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 97 | 0 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

# 60. **Bayesian data analysis**

- How likely are we to believe the hypothesis that ECMO and conventional therapy have the same survival rates? Just add the cells along the diagonal (0+0+...+5+57+97+0) to get 159 out of a thousand. Prior to collecting the data, we placed the probability that the two rates were equal at 500 out of a thousand, so the data has greatly (but not completely) dissuaded us from this belief.

# 61. **Bayesian data analysis**

You can calculate the probability that

- ECMO is exactly 10% better than conventional therapy

  - 0+0+...+1+13+84+0 = 98 /1000,

- ECMO is exactly 20% better

  - 0+0+...+13+218+0 = 231 / 1000,

- exactly 30% better

  - 0+0+...+7+178+0 = 185 /1000,

- and so forth.

# 62. Bayesian data analysis

- Here's something fun that Dr. Albert didn't show. You could take each of the cells in the table, compute a ratio of survival rates and then calculate the median of these ratios as 1.5.

| Relative risk | Prob. | Cumul. | |
|---|---|---|---|
| 0.8 / 0.9 = 0.89 | 1 | 1 | |
| 0.9 / 0.9 = 1.00 | 97 | 98 | |
| 0.8 / 0.8 = 1.00 | 57 | 155 | |
| 0.7 / 0.7 = 1.00 | 5 | 160 | |
| 0.9 / 0.8 = 1.12 | 84 | 244 | |
| 0.8 / 0.7 = 1.14 | 13 | 257 | |
| 0.7 / 0.6 = 1.17 | 1 | 258 | |
| 0.9 / 0.7 = 1.29 | 174 | 432 | |
| 0.8 / 0.6 = 1.33 | 16 | 448 | |
| 0.9 / 0.6 = 1.50 | 218 | 666 | << Median |
| 0.8 / 0.5 = 1.60 | 13 | 679 | |
| 0.9 / 0.5 = 1.80 | 178 | 857 | |
| 0.8 / 0.4 = 2.00 | 7 | 864 | |
| 0.9 / 0.4 = 2.25 | 97 | 961 | |
| 0.8 / 0.3 = 2.67 | 2 | 963 | |
| 0.9 / 0.3 = 3.00 | 32 | 995 | |
| 0.9 / 0.2 = 4.50 | 5 | 1000 | |

# 63. **Bayesian data analysis**

- Dr. Albert goes on to show an informative prior distribution. There is a fair amount of data to indicate that the survival rate for the conventional therapy is somewhere between 10% and 30%, but little or no data about the survival rates under ECMO.

# 64. **Bayesian data analysis**

- Here's a prior distribution that utilizes this historic information.

|      | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1 | Total |
|------|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|---|-------|
| 0    | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 | 0     |
| 0.1  | 0 | 50  | 50  | 48  | 39  | 34  | 25  | 17  | 11  | 6   | 0 | 280   |
| 0.2  | 0 | 54  | 54  | 51  | 42  | 36  | 27  | 18  | 12  | 6   | 0 | 300   |
| 0.3  | 0 | 38  | 38  | 36  | 29  | 25  | 19  | 13  | 8   | 4   | 0 | 210   |
| 0.4  | 0 | 22  | 22  | 20  | 17  | 14  | 11  | 7   | 5   | 2   | 0 | 120   |
| 0.5  | 0 | 11  | 11  | 10  | 8   | 7   | 5   | 4   | 2   | 1   | 0 | 60    |
| 0.6  | 0 | 5   | 5   | 5   | 4   | 4   | 3   | 2   | 1   | 1   | 0 | 30    |
| 0.7  | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 | 0     |
| 0.8  | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 | 0     |
| 0.9  | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 | 0     |
| 1    | 0 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0 | 0     |
| Total| 0 | 180 | 180 | 170 | 140 | 120 | 90  | 60  | 40  | 20  | 0 |       |

# 65. Repeat of pop quiz

A research paper computes a p-value of 0.45. How would you interpret this p-value?

1. Strong evidence for the null hypothesis
2. Strong evidence for the alternative hypothesis
3. Little or no evidence for the null hypothesis
4. Little or no evidence for the alternative hypothesis
5. More than one answer above is correct.
6. I do not know the answer.

# 66. Repeat of pop quiz

A research paper computes a confidence interval for a relative risk of 0.82 to 3.94. What does this confidence interval tell you.

1. The result is statistically significant and clinically important.
2. The result is not statistically significant, but is clinically important.
3. The result is statistically significant, but not clinically important.
4. The result is not statistically significant, and not clinically important.
5. The result is ambiguous.
6. I do not know the answer.

# 67. Repeat of pop quiz

A Bayesian data analysis can incorporate subjective opinions through the use of

1. Bayes rule.
2. data shrinkage.
3. a prior distribution.
4. a posterior distribution.
5. p-values.
6. I do not know the answer.