

Module 3F: Hypothesis Testing

Applied Epistemology: A Framework for how we know things scientifically

Agenda:

1. H_0/H_A : Our model of the test universe (the distribution of the variable)
2. **Test & assumptions:** are the assumptions met? Is the test valid?
3. **Quantitative evidence: p-value**, or critical value.
 - False positive = Type I (α), False Negative = Type II (β), Type III errors
 - Sensitivity, Specificity, Power \rightarrow confusion matrix, ROC/AUC curve
 - Positive Predictive Power, Negative Predictive Power
 - Confusion Matrix
 - **ROC/AUC curve**
4. **Conclusion & uncertainty/estimation**

Two independent studies are performed to test the same null hypothesis.

What is the probability that one or both of the studies obtains a significant result and rejects the null hypothesis ***even if the null hypothesis is true***? Assume that in each study there is a **0.05** probability of rejecting the null hypothesis.

$$\begin{aligned} P[\text{rejecting} | H_0 \text{ is true}] &= P[\text{study 1 reject} | H_0] + P[\text{study 2 reject} | H_0] - P[\text{study 1 reject} | H_0] * P[\text{study 2 reject} | H_0] \\ &= 0.05 + 0.05 - 0.025 = 0.0975 \end{aligned}$$

You can consider the previous question in one of two equally valid ways:

$$\begin{aligned} &P(\text{at least 1 study obtains significant results}) \\ &= 1 - P(\text{neither study obtains significant results}) \\ &= 1 - (1 - 0.05)^2 = 0.0975 \end{aligned}$$

$$\begin{aligned} &P[1^{\text{st}} \text{ study significant OR } 2^{\text{nd}} \text{ study is significant}] \\ &= (0.05) + (0.05) - (0.05)^2 = 0.0975 \end{aligned}$$

The experimenter thinks that they are using an $\alpha=0.05$, but they are actually using an $\alpha=0.0975$

223 admissions reasons, 223×12 hypothesis

Austin et al (2006): sifted through health care data for >10 million residents and **223** different reasons for admissions; **12 astrological signs**.

Conclusion: 72 conditions were significantly associated with a particular zodiac sign.

- **This is actually 223×12 hypothesis being tested (~2500 hypothesis)**
- You expect 134 statistically significant associations just due to chance so 72 is < 134 (calculated with $\alpha = 0.05$)

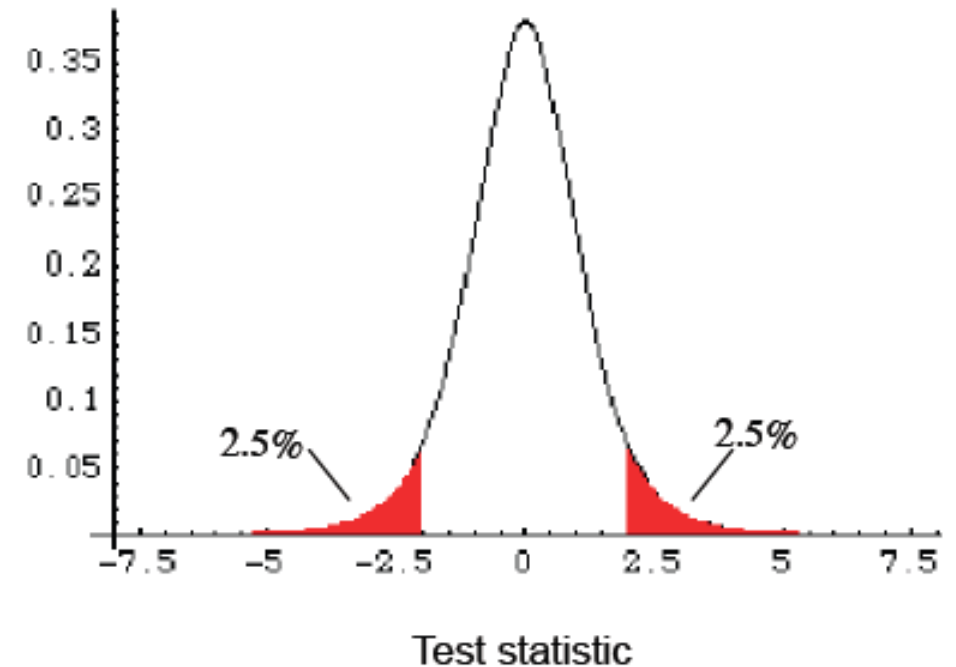
Bonferroni correction

$\alpha^* = \alpha / \text{num of hypothesis} = 0.05 / (223 \times 12) = 0.0000187$

- (GWAS tests hundreds of thousands if not millions at a time)

One tailed and two tailed tests:

- Most tests are two-tailed tests
- This means that a deviation in either direction would reject the null hypothesis
 - this means that α is divided into $\frac{\alpha}{2}$ on the one side and $\frac{\alpha}{2}$ on the other



One Tailed Tests:

Only used when the other tail is nonsensical

- o **Example:**

- o Comparing grades on a multiple-choice test to random guessing

- o **Example:**

- o Do daughters resemble their biological fathers?
 - o Experiment involves a subject who examines photo of one girl and two adult men and guesses the father
 - o If subjects pick father correctly > 0.5 then the hypothesis being tested would FTR
 - o Wouldn't make sense that daughters would, on average, resemble their biological fathers less than other men.

- Some parting words & popular misconceptions:
 - FTR does not mean ACCEPT
 - We ***never*** accept the null hypothesis
 - If FTR the null hypothesis, we can conclude that the data is compatible with the hypothesis
- If the result is statistically significant there is a temptation to believe that the effect is large. DO NOT GIVE IN THIS TO ERRONEOUS BELIEF.
 - Nor does it mean that the effect is interesting
 - If the sample size is large (and measurements have little variation) then even inconsequential differences will be significant
- P-values are calculated from the data itself. In contrast, the alpha value is set by the experimenter prior to conducting the experiment. P-values and alpha are related BUT THEY ARE NOT THE SAME!

- Why use hypothesis testing at all?
 - Why don't we skip hypothesis testing since confidence intervals give us similar information **plus** gives us information about the actual magnitude of the parameter?
 - *Main purpose of hypothesis testing is to determine if sufficient evidence has been presented to support a scientific claim*
- Deploy these tools wisely
 - Just because something is statistically significant does not mean it is biologically important or interesting
 - **Almost any null hypothesis can be rejected with a large enough sample**