

To determine if effects seen in the sample appear in the larger population,

- can use classical hypothesis testing
- there is also,
 - * Fisher null hypothesis testing
 - * Neyman-Pearson decision theory
 - * Bayesian Inference

Classical Hypothesis Testing achieves,

given a sample and an opponent effect, what is the probability of seeing such an effect by chance.

- ① Quantify the size of the opponent effect by choosing a test statistic (e.g. a diff in means)
- ② Define a Null Hypothesis: a model of the real system based on the effect not being real (e.g. the distributions between both groups is identical)
- ③ Compute the p-value - the probability of seeing the opponent effect if the null hypothesis is true.
- ④ Interpret result. If p-value is low, the result is statistically significant - the effect is more likely to appear in the larger population.

⇒ Similar to proof by contradiction.

Null Hypothesis Test. / of Sample Data

Can Toss Example

Toss a coin 1000 $H=400, T=600$ (you will have to make inference)

Create Model $50\% H, 50\% T$ (is the effect really in a sample?)

Score → toss a coin 250x (1000/4) (applying to a larger population?)

↳ Count the number of times where $H \neq T$ occurs as many or more times than the data's statistic.

↳ Use that count and the number of experiments to determine the probability of the results happening given the model - based on the future not having a real effect.

↳ This is p-value

- p-value ≤ 0.05 is by convention stat. sig.
- in practice p-value threshold depends on the test statistic and the model of the null hypothesis

↳ maybe less if is sig.
1% - 10% is borderline.
50%+ is insignificant.





Null Hypothesis: Pregnancy Length



One sided versus two-sided - considering one or both sides of the distribution.

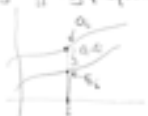
Chi-Squared Test - better for testing proportional differences - rather than using the total deviation.

① Testing Assumptions

unbiased estimate of the population

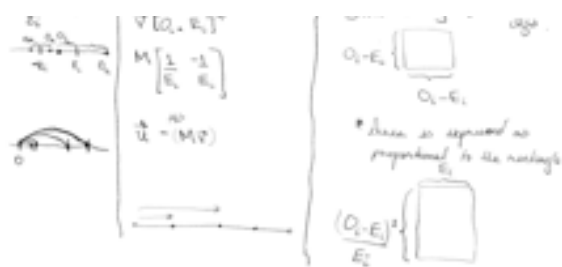
- p-value depends on the choice of test statistic. as the test statistic influences the perceived probability of certain experiments.
- The model of the null hypothesis also influences the perceived probability of experiments by making them more or less likely to happen.
- squaring gives weight to largely differing proportions

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

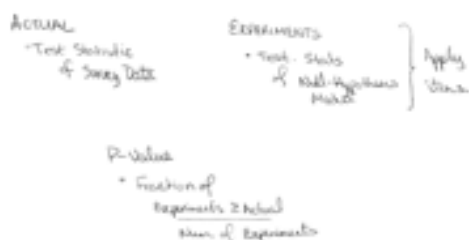
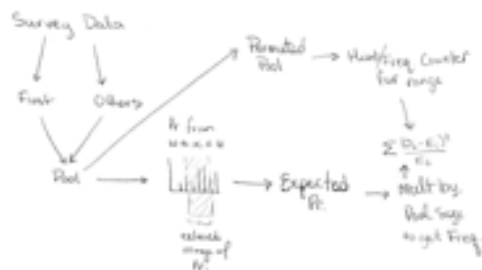


Mathematical Intuition

Two-sided test: Transformation of $(\mu_1 - \mu_0)^2$ | Distributions for μ_1 and μ_0 Standard



Diagram



Errors Surrounding the P-Value.

False-Positive Rate - probability of wrongly considering something significant.

False-Negative Rate - probability that the hypothesis test will fail when the effect is real!

False Positives occur with an opportunity of $1/20 = 5\%$.



Cumulative distribution of Null Hypothesis

Cumulative P-val for random alt in distrib. $P\text{-val} = 1 - \text{CDF}(x) = \text{CCDF}(x)$

If $\text{CCDF}(x) \leq 1/20 \rightarrow \text{stat. is significant}$

Since only 5th Percentile of the distribution would qualify on statistical anomalies, that means that if there was real data which would pass the hypothesis test, $1/20$ chances of the data being distributed in way that would then fall in the 5th (1%) Percentile.

False Negatives

Depends on the amount of data which represents the effect. - the effect size.

Option 1 - compute false negative rate from $1 - 1/20 = 0.95$



This yields the **significance test**.

The correct positive rate is called the **power** of the test, or the sensitivity.

The complement of the false negative rate - is the correct positive rate.

	negative	positive
false	x	y
correct	a	w

A test power of 80% is considered acceptable - for detecting difference. Below this threshold the test is underpowered.

A negative hypothesis test doesn't imply the absence of a relation, merely the absence of evidence in the sample data to prove that relationship.

Since we know the data correctly not a positive correlation, the only other way to classify the result is incorrectly - that is, likely as a negative. The data couldn't be classified correctly as negative nor wrongfully/falsely as positive.

Analogy: like in life when you hear/see the right things, and in your situation with the world, usually produces results of success with many of those events observed (even if not perfect).

Accuracy of Inference - Replicating Results.

Statistical tools - like hypothesis tests, are inaccurate, and they can not be utilized without the introduction of inaccuracy or loss of precision.

Therefore, exploration and then analysis/descriptive analysis of the same dataset are inherent bias: notably your findings might be used to construct your hypotheses and your findings are subsequently found within your sample.

ALTERNATIVES - also known as Cohen & Tukey's Procedures.

- Split the exploration & testing data portion.
 - Adjust the α level - we want a α level (e.g. 0.05) to compensate (i.e. necessarily reflect the inaccuracy that comes from the strength of your evidence) via Holm-Bonferroni Method, Bonferroni or Sidak.
 - result replication
 - first paper is considered exploratory
 - the second is confirmatory
 - Control False Discovery Rate (FDR) allows type I error to be higher if you have a high number of hypotheses, but control proportion of false discoveries (i.e. α' is more accurate).
- FWER $\leq \alpha$ $\alpha' = 1 - (1 - \alpha)^{1/n}$

- Let H_1, \dots, H_m be a family of null hypotheses
 $\& R_1, \dots, R_m$ the corresponding p -vals
- Order p -vals (smallest to largest) $R_{(1)}, \dots, R_{(m)}$ according to p -val
- Given significance level α , let k be the min index
such that $R_{(k)} > \frac{\alpha}{m+1-k}$
- Reject hypotheses $H_{(1)}, \dots, H_{(k-1)}$ d.k.
- If $k=1$ keep all
- If no k , reject all

- Most researchers accept an alpha of 0.05.
- Type 1 Error: erroneously reject null hypothesis, when true for the population the data does
- Type 2 Error: erroneously accept null hypothesis when false for the population.
- Correction should be applied to p -values when two or more statistical analyses have been performed on the same sample data.
This is due to the increase in familywise type 1 error rate.
- Family-wise error rate (FWER) (inflation):
- the probability of making at least 1 type 1 error
for a family of tests - essentially a series of tests on data.

$$FWER \leq 1 - (1 - \alpha_{adj})^k$$

$$\alpha_{adj} = \text{alpha level for}$$