

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, or is it just a fluke?)

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 fluke not having a real effect

↳ this is p-value

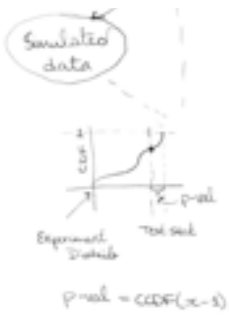
- p-value  $\leq 0.05$  is by convention stat. sig.
- in practice p-value threshold depends on the test statistic and the model of the null hypothesis

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

Get Data → apply → Test Stat

build your model → do experiments

iterations



Null Hypothesis: Pregnancy Length

**Survey Data**

First Born  $\mu$   $\mu_{1st}$

Other Born  $\mu$   $\mu_0$

**Testing Assumptions**

- 1. A common effect is that is the differences in means between two groups.
- 2. Assume the distrib is the same
- 3. Create null hypothesis model by preserving the detailed class
- 4. Measure probability of the statistic

repeating based on the data and the null hypothesis model given gives how no effect.

**Probability Mass Function**

diff in weeks

**Shuffle / Permute / Draw / Sample**

Arrange all datapoints

First  $\mu$   $\mu_{1st}$

Other  $\mu$   $\mu_0$

**repeated N-times**

no experiment data

**(3. p-val)**

diff in weeks

**Testing Assumptions**

- 1. Can also do correlation
- 2. Pearson's
- 3. Spearman's

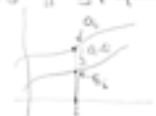
diff in weeks

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.   
 unbiased estimate of variance

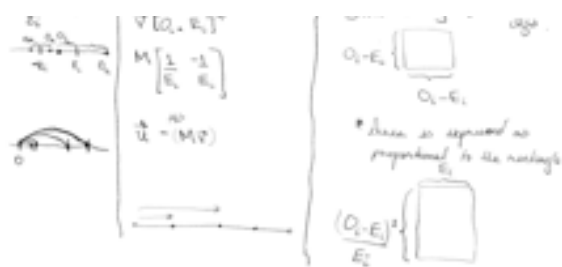
- p-value depends on the choice of test statistic. as the test statistic influences the perceived probability of certain experiments.
- The model of the null hypothesis also influences the perceived probability of experiments by making them more or less likely to happen.
- Squaring gives weight to largely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

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

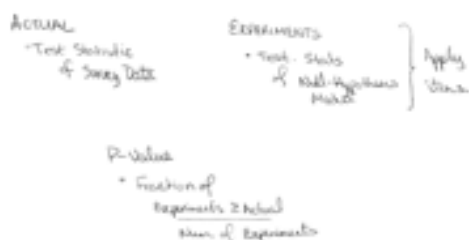
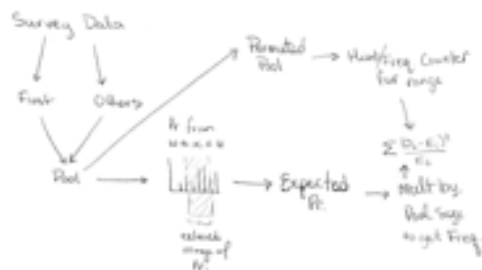


**Mathematical Notation**

Two-sided test: Transformation of  $(\chi^2 - \chi^2_{1-\alpha/2})$  Distribution has  $E_i$  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 upperbound 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 5<sup>th</sup> Percentile of the distribution would qualify on statistical anomalies, that means that if there was real data which would pass the hypothesis test,  $\approx 1/20$  chances of the data being distributed in way that would then fall in the 5<sup>th</sup> (1%) Percentile.

False Negatives

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

Option 1 - compute false negative rate from  $\frac{1}{N} \sum_{i=1}^N \text{P}(x_i < \text{threshold})$



This yields the rejection rate.

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 in the sampled 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 heard that the night changes, and in your situation with the world currently produces signs of success with many of those events observed (even negative)

### 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 French Connection to Tukey's Procedure.

- Split the exploration & testing data portion.
  - adjust the pre-set threshold to compensate (control the rate - we want a - wrong - false alarm)  
(i.e. necessarily reflect the inaccuracy that comes from the strength of your evidence)  
via Holm-Bonferroni Method, Benjamini or Sidak.
  - result replication  
→ first paper is considered exploratory  
→ the second is confirmatory
  - Control False Discovery Rate (FDR) allows type 1 error to regulation of the hypothesis but control proportion of false discoveries (i.e. more accurate)
- Holm-Bonferroni Method (Sidak is more accurate)
- FWER  $\leq \alpha$   $\alpha' = 1 - (1 - \alpha)^{1/n}$

- Let  $H_1, \dots, H_m$  be a family of null hypotheses
- &  $P_1, \dots, P_m$  be corresponding p-values
- Order p-values (smallest to largest)  $P_{(1)}, \dots, P_{(m)}$  according to hypothesis
- Given significance level  $\alpha$ , let  $k$  be the min index such that  $P_{(k)} > \frac{\alpha}{m+1-k}$
- 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.
$$FWER \leq 1 - (1 - \alpha_{adj})^k$$

$\alpha_{adj}$  - alpha level for