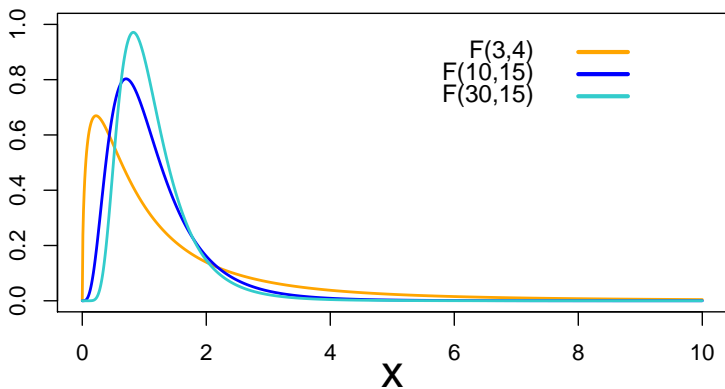


# Null Hypothesis Significance Testing

## Gallery of Tests

18.05 Spring 2022



# Announcements/Agenda

## Announcements

- None

## Agenda

- Discussion of Studio 7
- Frequentist methods don't know the probability a hypothesis is true.
- Other tests: chi-square for homogeneity, F-test (ANOVA) for equal means.
- Pay attention to how all the tests follow the same pattern. Differences are in details of how the test statistic is computed and, therefore, the type of null distribution.
- Multiple testing

## Discussion of Studio 7 and simulation

- What is a simulation?
  - Run an experiment with pseudo-random data instead of real-world real random data.
  - By doing this many times we can estimate the statistics for the experiment.
- Why do a simulation?
  - In the real world we are not omniscient.
  - In the real world we don't have infinite resources.
- What was the point of Studio 7?
  - To simulate some simple significance tests and compare various frequencies.
  - Simulated  $P(\text{reject}|H_0) \approx \text{significance } \alpha$
  - Simulated  $P(\text{reject}|H_A) \approx \text{power}$
  - $P(H_0|\text{reject})$  can be anything; depends on the distribution of hypotheses, which we almost never know.

## Studio 7 simulation results

	H0	HA	
reject	n_reject_and_H0	n_reject_and_HA	
	2. 468	2. 0	2. 468
	3. 0	3. 2101	3. 2101
	4. 344	4. 628	4. 972
non-reject	n_nonreject_and_H0	n_nonreject_and_HA	
	2. 9532	2. 0	2. 9532
	3. 0	3. 7899	3. 7899
	4. 6604	4. 2424	4. 9028
	2. 10000	2. 0	2. 10000
	3. 0	3. 10000	3. 10000
	4. 6948	4. 3052	4. 10000

	Problem 2	Problem 3	Problem 4
$P(\text{reject} H_0)$	$0.0468 \approx \alpha$	–	$\approx \alpha$
$P(\text{reject} H_A)$	–	$\approx \text{power}$	$\approx \text{power}$
$P(H_0 \text{reject})$	1	0	$\approx 0.32$
$P(H_A \text{reject})$	0	1	$\approx 0.65$

## Concept question: $t$ -test odds

We run a two-sample  $t$ -test for equal means, with  $\alpha = 0.05$ , and obtain a  $p$ -value of 0.04. What are the odds that the two samples are drawn from distributions with the same mean?

- (a) 19/1    (b) 1/19    (c) 1/20    (d) 1/24    (e) unknown

## General pattern of NHST

You want to decide whether to reject  $H_0$  in favor of  $H_A$ .

Design:

- Design experiment to collect data relevant to hypotheses.
- Choose test statistic  $x$  with known null distribution  $\phi(x | H_0)$ .
- Choose the significance level  $\alpha$  and find the rejection region.
- For a simple alternative  $H_A$ , use  $\phi(x | H_A)$  to compute the power.

Alternatively, you can choose both the significance level and the power, and then compute the necessary sample size.

Implementation:

- Run the experiment to collect data.
- Compute the statistic  $x$  and the corresponding  $p$ -value.
- If  $p < \alpha$ , reject  $H_0$ .

## Chi-square test for homogeneity

**Setting:** We have several data sets (for example results of applying several different treatments; or polling data from several different states).

**Homogeneity** (the null hypothesis) means that the data sets are all drawn from the same distribution, e.g. that all the treatments are equally effective, or that Connecticut voters have the same political opinions as those in New York.

**Example.** Three treatments for a disease are compared in a clinical trial, yielding the following data:

	Treatment 1	Treatment 2	Treatment 3
Cured	50	30	12
Not cured	100	80	18

Use a chi-square test to compare the cure rates for the three treatments, i.e. to test if all three cure rates are the same.

## Solution

$H_0$  = all three treatments have the same cure rate.

$H_A$  = the three treatments do not all have the same cure rate.

### Expected counts

- Under  $H_0$  the MLE for the cure rate is  
(total cured)/(total treated) =  $92/290 = 0.317$ .
- Assuming  $H_0$ , the expected number cured for each treatment is estimated by the number\_treated  $\times 0.317$ .
- This gives the following table of observed and expected counts (observed: on left, in black; expected: on right, in blue).
- We include the marginal totals (in orange). These are all used to compute the expected counts.

	Treatment 1	Treatment 2	Treatment 3	
Cured	50, 47.6	30, 34.9	12, 9.5	92
Not cured	100, 102.4	80, 75.1	18, 20.5	198
Total	150	110	30	290

*continued*



## Solution continued

Likelihood ratio statistic:  $G = 2 \sum O_i \ln(O_i/E_i) = 2.12$

Pearson's chi-square statistic:  $X^2 = \sum \frac{(O_i - E_i)^2}{E_i} = 2.13$

### Degrees of freedom

- **Using the formula:** Test for homogeneity  
 $df = (2 - 1)(3 - 1) = 2$ .
- **By counting:** The marginal totals are fixed because they are needed to compute the expected counts. So we can freely put values in 2 of the cells and then all the others are determined:  
degrees of freedom = 2.

**Null distribution:**  $\phi(X^2|H_0) \sim \chi^2(2)$  (Chi-square with 2 degrees of freedom.)

**p-value:**  $p = 1 - \text{pchisq}(2.12, 2) = 0.346$  (right-sided)

The data does not support rejecting  $H_0$ . We do not conclude that the treatments have differing efficacy.

## Board question: Khan's restaurant

Sal is thinking of buying a restaurant and asks about the distribution of lunch customers. The owner provides row one below. Sal records the data in row two himself one week.

	M	T	W	R	F	S
Owner's distribution	0.1	0.1	0.15	0.2	0.3	0.15
Observed # of cust.	30	14	34	45	57	20

Set the significance level ahead of time.

$H_0$ : the owner's distribution is correct.

$H_A$ : the owner's distribution is not correct.

Compute both  $G$  and  $X^2$ .

Run a chi-square [goodness-of-fit test](#) on the null hypotheses:

## Board question: genetic linkage

In 1905, William Bateson, Edith Saunders, and Reginald Punnett were examining flower color and pollen shape in sweet pea plants by performing crosses similar to those carried out by Gregor Mendel.

Purple flowers (P) is dominant over red flowers (p).

Long seeds (L) is dominant over round seeds (l).

F<sub>0</sub>: PPLL × ppll (initial cross)

F<sub>1</sub>: PpLl × PpLl (all second generation plants were PpLl)

F<sub>2</sub>: 2132 plants (third generation)

$H_0$  = independent assortment: color and shape are inherited independently.

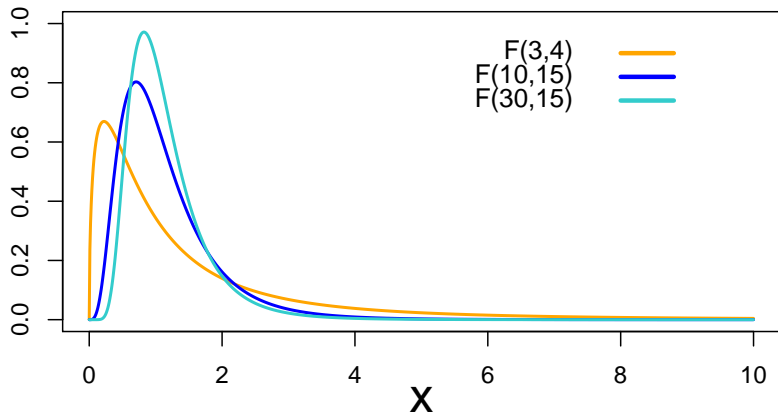
Here is the data from their experiment.

	purple, long	purple, round	red, long	red, round
Expected	?	?	?	?
Observed	1528	106	117	381

Determine the expected counts for  $F_2$  under  $H_0$  and find the  $p$ -value for a Pearson chi-square test. Explain your findings biologically.

## $F$ -distribution

- Notation:  $F_{a,b}$ ,  $a$  and  $b$  degrees of freedom
- Derived from normal data
- Range:  $[0, \infty)$
- mean  $\approx 1$



## $F$ -test = one-way ANOVA

Like  $t$ -test but for  $n$  groups of data with  $m$  data points each.

$y_{i,j} \sim N(\mu_i, \sigma^2)$ ,  $y_{i,j}$  =  $j^{\text{th}}$  point in  $i^{\text{th}}$  group. All  $y_{i,j}$  independent.

$H_0$ : means are all equal, i.e.  $\mu_1 = \dots = \mu_n$ .

Assumptions: data points independent, **group variances  $\sigma^2$  equal**.

$MS_B$  = between group variance =  $\frac{m}{n-1} \sum_{i=1}^n (\bar{y}_i - \bar{y})^2$

$MS_W$  = within group variance = mean of  $s_1^2, \dots, s_n^2 \approx \sigma^2$ .

Idea: If  $\mu_i$  are equal,  $MS_B \approx \sigma^2$ , so this ratio should be near 1.

test statistic:  $f = \frac{MS_B}{MS_W}$ . Under  $H_0$ :  $f \sim F_{n-1, n(m-1)}$ .

Note: Formulas easily generalize to unequal group sizes:

<https://en.wikipedia.org/wiki/F-test>

## Board question: Recovery

The table shows recovery time in days for three medical treatments.

**(a)** Set up and run an F-test testing if the average recovery time is the same for all three treatments. Use significance level 0.05.

**(b)** Based on the test, what might you conclude about the treatments?

$T_1$	$T_2$	$T_3$
6	8	13
8	12	9
4	9	11
5	11	8
3	6	7
4	8	12

**Note:** For  $\alpha = 0.05$ , the critical value of  $F_{2,15}$  is 3.68.

## Concept question: multiple-testing

**(a)** Suppose we have 6 treatments and want to know if the average recovery time is the same for all of them. If we compare two at a time, how many two-sample  $t$ -tests do we need to run?

- (i)** 1      **(ii)** 2      **(iii)** 6      **(iv)** 15      **(v)** 30

## Concept question: multiple-testing

**(a)** Suppose we have 6 treatments and want to know if the average recovery time is the same for all of them. If we compare two at a time, how many two-sample  $t$ -tests do we need to run?

- (i)** 1      **(ii)** 2      **(iii)** 6      **(iv)** 15      **(v)** 30

**(b)** Suppose we use the significance level 0.05 for each of the 15 tests. Assuming the null hypothesis, what is the best estimate of the probability that we reject at least one of the 15 null hypotheses?

- (i)**  $< 0.05$       **(ii)** 0.05      **(iii)** 0.10      **(iv)**  $> 0.25$



## Discussion

Recall that there is an  $F$ -test that tests if all the means are the same. What are the trade-offs of using the  $F$ -test rather than many two-sample  $t$ -tests?

## Discussion

Recall that there is an  $F$ -test that tests if all the means are the same. What are the trade-offs of using the  $F$ -test rather than many two-sample  $t$ -tests?

The advantage of the  $F$ -test is that it is a single test. So, there is no need to adjust the significance of multiple tests to achieve the correct overall significance. The disadvantage is that it only tests if all the means are the same. It does not say which pairs are different.

## Board question: chi-square for independence

(From Rice, *Mathematical Statistics and Data Analysis*, 2nd ed. p.489)

Consider the following contingency table of counts

Education	Married once	Married multiple times	Total
College	550	61	611
No college	681	144	825
Total	1231	205	1436

Use a chi-square test with significance level 0.01 to test the hypothesis that the number of marriages and education level are independent.

MIT OpenCourseWare

<https://ocw.mit.edu>

18.05 Introduction to Probability and Statistics

Spring 2022

For information about citing these materials or our Terms of Use, visit:

<https://ocw.mit.edu/terms>.