

① Two-Sample Inferences

In Chapters 6 & 7, you learnt how to test a particular value of a parameter (proportion p or mean μ) in one population.

Sometimes, you don't compare one population to a particular value of its parameter (p_0 or μ_0), but two populations together

- Ex:
- Average SAT score of men VS women in UCSD
 - Average height of San Diegans VS New Yorkers
 - Average age of husbands & wives
 - Average income of children compared to their parents

In all these examples, you have 2 samples coming from 2 populations

X_1, \dots, X_n : Population 1

Y_1, \dots, Y_m : Population 2

⚠ n and m have no reason to be equal in general.

(2)

9.0 Paired Samples

Sometimes though, the members of the 2 populations have some direct link where each member of one population is paired with a member of the other. In such a case, $m = n$.

- Ex:
- Weight Pre-diet / Weight Post-diet
 - Husband Age / Wife Age

We have $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu_x, \sigma_x^2)$ (all parameters unknown)
 $Y_1, \dots, Y_m \sim N(\mu_y, \sigma_y^2)$

If we want to test $H_0: \mu_x - \mu_y = 0$, ($\text{i.e. } \mu_x = \mu_y$)
 we just consider the differences

$$\boxed{\Delta_i = X_i - Y_i} \quad (1 \leq i \leq n)$$

and we perform a 1 sample t-test on $\Delta_1, \dots, \Delta_n$

(3)

Ex: You decide to research global warming in the U.S. You choose 62 random cities and look up the high temperature on Jan 1st 1970 and Jan 1st 2018.

Clearly, the data are paired: hot locations will have high readings in both 1970 and 2018. Similarly for cold locations.

You calculate $\Delta = \text{Temp}_{2018} - \text{Temp}_{1970}$ for each location and find the differences Δ_i ; have $\bar{\Delta}_{62} = 1.1^\circ\text{F}$ and $S_\Delta^2 = \frac{1}{62-1} \sum_{i=1}^{62} (\Delta_i - \bar{\Delta}_{62})^2 = 24.1$.

If we care about the true (unknown) average temperature rise $\mu_{2018} - \mu_{1970}$ in the U.S., our confidence interval would be

$$CI_{1-\alpha} = \left[\bar{\Delta}_{62} - t_{\frac{1-\alpha}{2}, 62-1} \frac{S_\Delta}{\sqrt{64}}, \bar{\Delta}_{62} + t_{\frac{1-\alpha}{2}, 62-1} \frac{S_\Delta}{\sqrt{64}} \right]$$

$$\underset{\alpha=5\%}{=} [-0.164^\circ\text{F}, 2.346^\circ\text{F}]$$

Rk: Nothing new statistically speaking, just the "difference trick".

(4) In what follows, we'll only focus on unpaired data, meaning that the X_i 's and Y_j 's are independent.

9.2 Testing $\mu_x = \mu_y$

We focus here on tests with $H_0: \mu_x - \mu_y = 0$. The alternative would be $>$, $<$ or \neq depending on the situation but the method is the same.

The variances σ_x^2 and σ_y^2 are not known, but the method we'll use will differ when $\sigma_x^2 = \sigma_y^2$ and when $\sigma_x^2 \neq \sigma_y^2$. We'll learn later (Chapter 9.3) how to actually test $H_0: \sigma_x^2 = \sigma_y^2$.

Equal Variances

Thm: let $\begin{cases} X_1, \dots, X_m \stackrel{iid}{\sim} N(\mu_x, \sigma_x^2) \\ Y_1, \dots, Y_n \stackrel{iid}{\sim} N(\mu_y, \sigma_y^2) \end{cases}$ be independent with $\sigma_x^2 = \sigma_y^2 = \sigma^2$

$$\text{then } T = \frac{(\bar{X}_m - \bar{Y}_m) - (\mu_x - \mu_y)}{\sqrt{\frac{1}{m} + \frac{1}{n}} S_p} \sim t_{m+n-2} \text{ - distribution}$$

$$\text{where pooled } S_p^2 = \frac{(m-1)S_x^2 + (n-1)S_y^2}{m+n-2} = \frac{\sum_{i=1}^m (X_i - \bar{X}_m)^2 + \sum_{j=1}^n (Y_j - \bar{Y}_m)^2}{m+n-2}$$

(5) Proof: One easily checks that $Q = \frac{\bar{X}_m - \bar{Y}_m - (\mu_x - \mu_y)}{\sqrt{\frac{1}{n} + \frac{1}{m}}}$ has distribution $N(0, \sigma^2)$.

Furthermore, if $W = (m-1)S_x^2 + (m-1)S_y^2$, we notice that:

$$\left\{ \begin{array}{l} \cdot \frac{(m-1)S_x^2}{\sigma^2} \sim \chi_{m-1}^2 \\ \cdot \frac{(m-1)S_y^2}{\sigma^2} \sim \chi_{m-1}^2 \\ \cdot S_x^2 \perp S_y^2 \text{ since the } X_i \text{'s and } Y_j \text{'s are independent} \end{array} \right.$$

(See Review on Gaussian vectors)

Hence, $\frac{W}{\sigma^2} \sim \chi_{m+m-2}^2$

Finally, $W \perp Q$ (Same), so that $T = \sqrt{m+m-2} \frac{Q}{W} \sim F_{m+m-2, m+m-2}$

Rk: • The assumption that $\sigma_x = \sigma_y = \sigma$ is key in the proof and cannot be released. It allows to compute a pooled variance S_p^2 , and hence, improves the power of the tests associated to this result.

• If $m+m-2 \geq 100$, you may use the quantiles of the η -table, since they get super close to those of the t -tables with $df \geq 100$.

⑥

Unequal variances

Thm: (Welch) let $\{X_1, \dots, X_m\} \stackrel{iid}{\sim} N(\mu_x, \sigma_x^2)$

$\{Y_1, \dots, Y_n\} \stackrel{iid}{\sim} N(\mu_y, \sigma_y^2)$

be independent

Then

$$T = \frac{(\bar{X}_m - \bar{Y}_n) - (\mu_x - \mu_y)}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}} \stackrel{\text{Approximately has distribution}}{\sim} t_{\nu},$$

$$\text{where } \nu = \frac{\left(\frac{s_x^2}{m} + \frac{s_y^2}{n} \right)^2}{\frac{s_x^4}{m^2(m-1)} + \frac{s_y^4}{n^2(n-1)}}$$

Proof: (rough)

$$T = \frac{\left(\frac{\bar{X}_m - \bar{Y}_n - (\mu_x - \mu_y)}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}} \right)}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}} \stackrel{\text{?}}{\sim} \frac{N(0, 1)}{\chi_{df}}$$

$$\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}} = \sqrt{\left(\frac{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}} \right)^2}$$

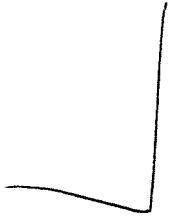
Need to find the best $\nu = df$ for this approximation.

(7) To find the best approximation (in distribution) of

$$\frac{\frac{S_x^2}{m} + \frac{S_y^2}{m}}{\frac{S_x^2}{m} + \frac{S_y^2}{m}} \text{ by } \frac{V}{J} \text{ where } V \sim \chi^2_1,$$

We equate the first two moments of both sides, and get

$$J = \frac{\left(\frac{S_x^2}{m} + \frac{S_y^2}{m} \right)^2}{\frac{S_x^4}{m^2(n-1)} + \frac{S_y^4}{m^2(n-1)}}.$$



- Rk:
- This is not a rigorous proof, only a heuristic. What we've shown is that the best value for J is the announced one, not that it actually does approximate well the distribution.
 - In practice, it does provide a good test statistic.
 - See "Satterthwaite approximation" for a more general statement of such an approximation by χ^2 distributions.

(8)

9.3 Testing $\sigma_x = \sigma_y$ (The F-test)

As highlighted by the previous Section 9.2, equality of variances helps designing more powerful tests and inference procedures. Also, variances provide an informative parameter that may describe populations. In what follows, we describe a test statistic to decide whether $\sigma_x = \sigma_y$ or not, from data.

First, let's describe a family of distributions that'll help in the process.

Def: (Fisher distribution) —

Given $d_1, d_2 \in \mathbb{N}$, we denote by F_{d_1, d_2} the distribution of the random variable

$$\frac{U_1/d_1}{U_2/d_2}, \text{ where } \begin{cases} \cdot U_1 \sim \chi^2_{d_1} \\ \cdot U_2 \sim \chi^2_{d_2} \\ \cdot U_1 \perp U_2 \end{cases}$$

F_{d_1, d_2} is called the Fisher distribution with (d_1, d_2) d.f.

An immediate consequence of Cochran yields the following

⑨

Thm: let $\begin{cases} X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu_x, \sigma_x^2) \\ Y_1, \dots, Y_m \stackrel{iid}{\sim} N(\mu_y, \sigma_y^2) \end{cases}$ be independent.

Then

$$\frac{\sigma_x^2}{\sigma_y^2} \times \frac{S_y^2}{S_x^2} \sim F_{m-1, n-1}$$

Proof: From Cochran, $\begin{cases} \cdot (n-1) S_x^2 \sim \chi^2_{n-1} \\ \cdot (m-1) S_y^2 \sim \chi^2_{m-1} \end{cases}$, and since the X_i 's and

Y_j 's are independent, so are $(n-1) S_x^2$ and $(m-1) S_y^2$, which proves the result.]

Rk: This allows to test $H_0: \sigma_x = \sigma_y$ without even knowing (nor having inferred) μ_x and μ_y .

Depending on the alternative hypothesis, the rejection region has 1 or 2 pieces.

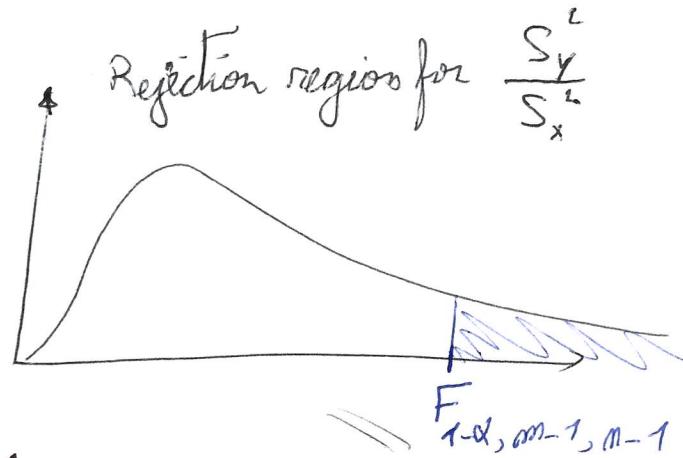
In all the cases the test statistic to use is

$$F = \frac{S_y^2}{S_x^2}$$

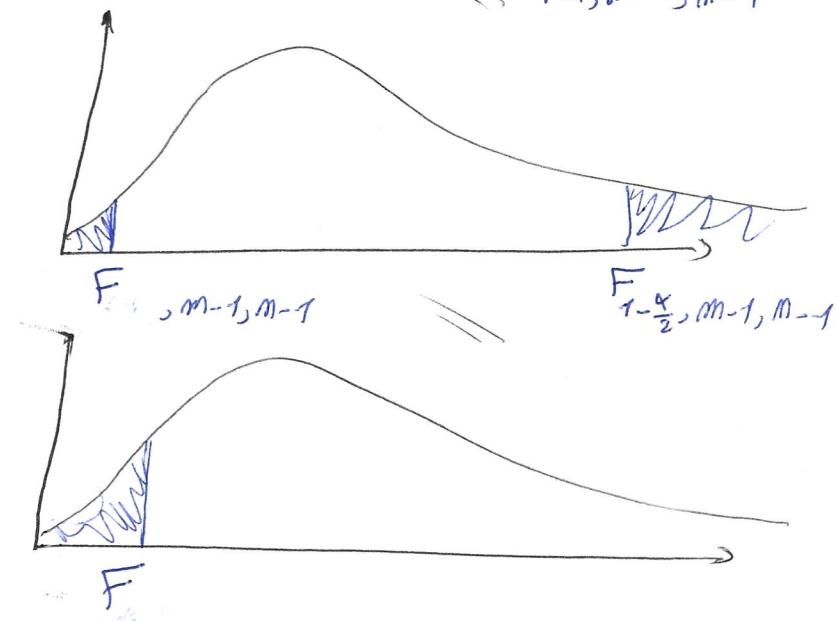
(or equivalently $\frac{S_x^2}{S_y^2}$)

(10)

$$H_1: \sigma_x < \sigma_y$$



$$H_1: \sigma_x \neq \sigma_y$$



$$H_1: \sigma_x > \sigma_y$$

- For testing means, the F-test is often used to determine if variances are equal.

1) State hypotheses $H_0: \mu_x = \mu_y$, $H_1: \mu_x \left\{ \begin{matrix} < \\ \neq \\ > \end{matrix} \right\} \mu_y$

2) Check the assumptions

- Paired data?
- Normality?
- Independence?
- Equal or not-equal variances?

→ Perform a F-test

3) Perform the right t-test.

9.4 Binomial Data: Testing $H_0: p_x = p_y$

Back to categorical data, we now want to infer equality of the proportions p_x and p_y of a trait in two populations.

Data then take the form of two sequences of 0-1 variables

$$\begin{cases} X_1, \dots, X_n \stackrel{i.i.d}{\sim} \text{Bernoulli}(p_x) = B(p_x) \\ Y_1, \dots, Y_m \stackrel{i.i.d}{\sim} \text{Bernoulli}(p_y) = B(p_y) \end{cases}$$

Thm: If n and m are large enough so that $\begin{cases} np_x, n(1-p_x) \geq 10 \\ mp_y, m(1-p_y) \geq 10 \end{cases}$

then

$$\frac{\left(\hat{p}_x - \hat{p}_y\right) - (p_x - p_y)}{\sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{n} + \frac{\hat{p}_y(1-\hat{p}_y)}{m}}} \sim N(0, 1)$$

Where $\hat{p}_x = \frac{X_1 + \dots + X_n}{n}$ and $\hat{p}_y = \frac{Y_1 + \dots + Y_m}{m}$

$$\begin{aligned} \hat{p}_x &= \bar{X}_n \\ \hat{p}_y &= \bar{Y}_m \end{aligned}$$

Proof:

From the C.L.T. $\begin{cases} \sqrt{n}(\hat{p}_x - p_x) \xrightarrow[n \rightarrow \infty]{D} N(0, p_x(1-p_x)) \\ \sqrt{m}(\hat{p}_y - p_y) \xrightarrow[m \rightarrow \infty]{D} N(0, p_y(1-p_y)) \end{cases}$

(12) Hence, since \hat{P}_x and \hat{P}_y are independent,

$$(\hat{P}_x - \hat{P}_y) - (P_x - P_y) \sim N(0, \frac{P_x(1-P_x)}{n} + \frac{P_y(1-P_y)}{m})$$

Hence,

$$\frac{(\hat{P}_x - \hat{P}_y) - (P_x - P_y)}{\sqrt{\frac{P_x(1-P_x)}{n} + \frac{P_y(1-P_y)}{m}}} \xrightarrow[m, m \rightarrow \infty]{D} N(0, 1)$$

But from the law of large number, $\begin{cases} \hat{P}_x \xrightarrow[n \rightarrow \infty]{a.s.} P_x \\ \hat{P}_y \xrightarrow[m \rightarrow \infty]{a.s.} P_y \end{cases}$, so that

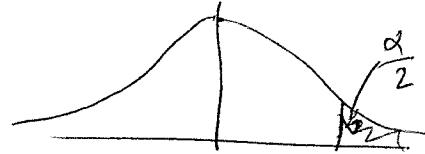
from Slutsky's theorem,

$$\frac{(\hat{P}_x - \hat{P}_y) - (P_x - P_y)}{\sqrt{\frac{\hat{P}_x(1-\hat{P}_x)}{n} + \frac{\hat{P}_y(1-\hat{P}_y)}{m}}} \xrightarrow[m, m \rightarrow \infty]{D} N(0, 1)$$

Rk: The conditions $nP_x, nP_y, n(1-P_x), n(1-P_y) > 10$ are here to ensure the convergence result holds. Indeed, if one of these conditions don't hold, the associated Bernoulli is skewed and the rate in the C.L.T. get really slow (See Berry-Esseen Thm).

(13) This result provides a confidence interval for $P_x - P_y$:

$$\left[(P_x - P_y) \pm z^* \sqrt{\frac{\hat{P}_x(1-\hat{P}_x)}{m} + \frac{\hat{P}_y(1-\hat{P}_y)}{m}} \right], \text{ where } z^* = N(0, 1)_{1-\frac{\alpha}{2}}$$



For testing $P_x = P_y$, one should use another statistic than the one above. Indeed, under the null, $P_x = P_y = p$, so that

$$\sqrt{\frac{P_x(1-P_x)}{m} + \frac{P_y(1-P_y)}{m}} = \sqrt{\left(\frac{1}{m} + \frac{1}{m}\right)p(1-p)}$$

As a consequence, we can use another estimator of $P_x = P_y = p$ to estimate the variance term. This would result in a better approximation in Slutsky, and hence a more powerful test.

We let $\hat{P}_e = \frac{X_1 + \dots + X_m + Y_1 + \dots + Y_m}{m+m} = \frac{m\hat{P}_x + m\hat{P}_y}{m+m}$

(14)

Prop: If $\hat{p}_x = \hat{p}_y = \hat{p}$, and that $(m+n)p, (m+n)(1-p) \geq 10$, then

$$\frac{\hat{p}_x - \hat{p}_y}{\sqrt{\left(\frac{1}{m} + \frac{1}{n}\right)\hat{p}_e(1-\hat{p}_e)}} \sim N(0, 1)$$

Proof: Same as before, except you use that $\hat{p}_e \rightarrow p$ is Slutsky -

Rk: You see that the condition $(m+n)p, (m+n)(1-p) \geq 10$ is weaker than previously. Actually, that'll provide a better power to the test that use this statistic

• When to pool / not to pool ?

* Confidence Interval : Do not pool

* Test : Pool

The reason for this is that for a test, you do your reasoning under $H_0: p_x = p_y = p$ but not for C. I.

15

- When performing a test for proportions, keep in mind to check that you actually are working with a Binomial model:

- counts derived from II Bernoulli trials
- Number of trials fixed in advance
- Probability of success constant across trials

9.5 Confidence intervals for two-sample problems

As for one-sample problems, the knowledge of the distributions of the above statistics (and others) provides confidence intervals.

Make sure you can write C.I.'s of level α in all the cases below:

Parameter of interest	Case	
$\mu_x - \mu_y$	• σ_x and σ_y known	
	• σ_x and σ_y unknown	• $\sigma_x = \sigma_y$ • $\sigma_x \neq \sigma_y$
$\frac{\sigma_x^2}{\sigma_y^2}$	• μ_x and μ_y known	
	• either μ_x or μ_y known	
	• neither μ_x nor μ_y known	

Can you do it for $\frac{\sigma_x}{\sigma_y}$ too?