

Summary of the hypothesis testing

All tests are of the following form:

- Population parameter: θ
- Null hypothesis $H_0: \theta = \text{some fixed number}$
- Estimate of θ based on samples from the population: $\hat{\theta}$
- Approximate $Z = \frac{\hat{\theta} - E(\hat{\theta})}{\text{sd}(\hat{\theta})}$ assuming H_0 is true.
- Convert the (approximate) Z -score to p -value by studying the sampling variability of Z (or the approximation of Z).
- if the p -value is really small (less than some cut off, say),
i.e. if $|Z|$ is large, then this is evidence against H_0 .
If the p -value is not small (greater than some cut off, say)
then the data is consistent with H_0 (but H_0 could still be false).

List of approximate z-scores & their sampling variability for the
Non-exact "quick guess method".

(2)

<u>Population Parameter</u>	<u>estimate</u>	<u>approximate z-score</u>	<u>approximate sampling variability</u>
(1) Population average μ (population s.d. σ known)	$\hat{\mu} = \frac{x_1 + \dots + x_n}{n}$	$Z = \frac{\hat{\mu} - \mu}{(\sigma/\sqrt{n})}$	$Z \approx N(0,1)$.
(2) Population average μ (population s.d. σ unknown)	$\left\{ \begin{array}{l} \hat{\mu} = \frac{x_1 + \dots + x_n}{n} \\ \hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2} \end{array} \right.$	$Z = \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{n}}$	$Z \approx N(0,1)$
(3) Population proportion p	$\hat{p} = \frac{x_1 + \dots + x_n}{n}$ where $x_i = \begin{cases} 0 \\ 1 \end{cases}$	$Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}}$	$Z \approx N(0,1)$
(4) The difference b/w two population averages $\mu_x - \mu_y$	$\left\{ \begin{array}{l} \hat{\mu}_x - \hat{\mu}_y \text{ when} \\ \hat{\mu}_x = \frac{x_1 + \dots + x_n}{n} \\ \hat{\mu}_y = \frac{y_1 + \dots + y_m}{m} \\ \hat{\sigma}_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_x)^2} \\ \hat{\sigma}_y = \sqrt{\frac{1}{m-1} \sum_{i=1}^m (y_i - \hat{\mu}_y)^2} \end{array} \right.$	$Z = \frac{(\hat{\mu}_x - \hat{\mu}_y) - (\mu_x - \mu_y)}{\sqrt{\frac{\hat{\sigma}_x^2}{n} + \frac{\hat{\sigma}_y^2}{m}}}$	$Z \approx N(0,1)$.

These are all of the form $\frac{\hat{\theta} - E(\hat{\theta})}{sd(\hat{\theta})}$ where MF_1 & MF_2 are used to compute $E(\hat{\theta})$ & $sd(\hat{\theta})$.

Modifications and assumptions
to make the "quick guess
method" official

(3)

For (1): You need n large enough, usually $n > 30$, for the CLT to kick in so $Z \approx N(0,1)$ becomes a good approximation.

For (2): If you assume the population to be normally distributed. Then Z is exactly a draw from $t^{(n-1)}$, i.e. a t -distribution with $n-1$ degrees of freedom.
Note: if $n > 30$ then $t^{(n-1)}$ is approx normal so $Z \approx N(0,1)$ works well.

For (3): You need n large (typically larger than what is needed for (1)) for the CLT to imply $Z \approx N(0,1)$

For (4): To get an exact pdf for Z , you need to assume both populations are Normal and $\sigma_x = \sigma_y$. Now replace

$$\sqrt{\frac{\hat{\sigma}_x^2}{n} + \frac{\hat{\sigma}_y^2}{m}} \text{ with } \hat{\sigma}_p \sqrt{\frac{1}{n} + \frac{1}{m}}$$

where $\hat{\sigma}_p$ is the "pooled est of s.d. from both samples $X_1, \dots, X_n, Y_1, \dots, Y_m$ " defined as

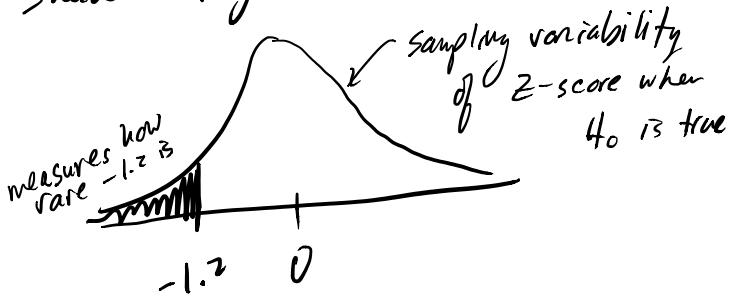
$$\hat{\sigma}_p := \sqrt{\frac{n-1}{n+m-2} \hat{\sigma}_x^2 + \frac{m-1}{n+m-2} \hat{\sigma}_y^2}$$

Now the "new" $Z = \frac{(\hat{\mu}_x - \hat{\mu}_y) - (\mu_x - \mu_y)}{\hat{\sigma}_p \sqrt{\frac{1}{n} + \frac{1}{m}}}$ is exactly a draw from $t^{(n+m-2)}$.

All the p-values done in this class are "one-sided"

To convert an approximate Z -score to a p-value I will always use only one side of the pdf for Z :

e.g. if the approximate Z -score = -1.2 then the p-value is given by the shaded region below



e.g. if the approximate Z -score = 2.5 then the p-value is given by the shaded region below

