# 1.017/1.010 Class 17
# Testing Hypotheses about Two Populations

## Tests of differences between two populations

To test if two populations $x$ and $y$ are different we can compare specified distributional properties $a_x$ and $a_y$ (means, variances, 90 percentiles, etc.).

Null hypothesis:

$$H0: \quad a_x = a_y = a_0 \quad \text{or} \quad a_x - a_y = 0$$

The hypothesis test may be based on "natural" (unbiased and consistent) estimators of $a_x$ and $a_y$, derived from the **independent** random samples $x_1$, $x_2,..., x_{Nx}$ and $y_1, y_2,...,y_{Ny}$ :

$$\hat{a}_x = \hat{a}_x(x_1, x_2,..., x_{Nx})$$

$$\hat{a}_y = \hat{a}_y(x_1, x_2,..., x_{Ny})$$

We can derive a two-sided rejection region $R_{a0}$ written in terms of a **standardized** statistic $z$ , following the same basic procedure as in the single population case (see Class 15):

$$z(\hat{a}_x,\hat{a}_y,a_x,a_y) = \frac{(\hat{a}_x - \hat{a}_y) - (a_x - a_y)}{SD[\hat{a}_x - \hat{a}_y]}$$

$$z(\hat{a}_x,\hat{a}_y,a_0,a_0) = \frac{(\hat{a}_x - \hat{a}_y)}{SD[\hat{a}_x - \hat{a}_y]}$$

$$R_{z0}: \quad z(\hat{a}_x,\hat{a}_y,a_0,a_0) \le z_L = F_z^{-1}(\frac{\alpha}{2})$$

$$z(\hat{a}_x,\hat{a}_y,a_0,a_0) \ge z_U = F_z^{-1}(1-\frac{\alpha}{2})$$

For **large samples** $z(\hat{a}_x,\hat{a}_y,a_0,a_0)$ has a unit normal distribution if $H0$ is true ($a_x = a_y = a_0$). Use `norminv` to compute $z_L$ and $z_U$ from $\alpha$.

We can also define a rejection region $R_{a0}$ written in terms of the **nonstandardized** estimates:

$$R_{a0} : \hat{a}_x - \hat{a}_y \le \Delta a_L = F_z^{-1}(\frac{\alpha}{2})SD[\hat{a}_x - \hat{a}_y]$$

$$\hat{a}_x - \hat{a}_y \ge \Delta a_U = F_z^{-1}(1 - \frac{\alpha}{2})SD[\hat{a}_x - \hat{a}_y]$$

The two-sided $p$-value is obtained from:

$$1 - p/2 = F_z[z] = F_z\left[\frac{(\hat{a}_x - \hat{a}_y)}{SD(\hat{a})}\right] \quad \hat{a}_x - \hat{a}_y \ge 0$$

$$p/2 = F_z[z] = F_z\left[\frac{(\hat{a}_x - \hat{a}_y)}{SD(\hat{a})}\right] \quad \hat{a}_x - \hat{a}_y \le 0$$

For large samples use `normcdf` to compute $p$ from $\dfrac{\hat{a}_x - \hat{a}_y}{SD[\hat{a}]}$ .

## Special Case: Large sample test of the difference between two means

If the property of interest is the mean then:

$$H0: \quad a_x = E[x] = a_y = E[y] \quad \text{or} \quad E[x] - E[y] = 0$$

Natural estimator of $E[x] - E[y]$ is $m_x - m_y$.

In large sample case $m_x - m_y$ is normal with mean and variance:

$$E[m_x - m_y] = E[x] - E[y] \quad \quad \text{(unbiased)}$$

$$Var[(m_x - m_y) = Var[m_x] + Var[m_y] = \frac{\sigma_x^2}{N_x} + \frac{\sigma_y^2}{N_y} \quad \text{(consistent)}$$

Construct a large sample test statistic $z \sim N(0,1)$ :

$$z = \frac{m_x - m_y}{\sqrt{\dfrac{\sigma_x^2}{N_x} + \dfrac{\sigma_y^2}{N_y}}} \approx \frac{m_x - m_y}{\sqrt{\dfrac{s_x^2}{N_x} + \dfrac{s_y^2}{N_y}}}$$

Two-sided rejection region written in terms of $m_x$ and $m_y$:

$$R_{a0} : m_x - m_y \leq \Delta a_L = F_z^{-1}(\frac{\alpha}{2})SD[\boldsymbol{m}_x - \boldsymbol{m}_y]$$

$$m_x - m_y \geq \Delta a_U = F_z^{-1}(1 - \frac{\alpha}{2})SD[\boldsymbol{m}_x - \boldsymbol{m}_y]$$

The two-sided p-value is obtained from:

$$1 - p/2 = F_z(z) = F_z\left[ (m_x - m_y)\left(\frac{s_x^2}{N_x} + \frac{s_y^2}{N_y}\right)^{-1/2} \right] \quad m_x \geq m_y$$

$$p/2 = F_z(z) = F_z\left[ (m_x - m_y)\left(\frac{s_x^2}{N_x} + \frac{s_y^2}{N_y}\right)^{-1/2} \right] \quad m_x \leq m_y$$

Example:  Comparing crop yields with and without fertilizer application

Consider two agricultural fields, one that is fertilized and one that is not.  Yield samples (kg/ha) from the two fields are as follows:

Fertilized ($x$):     66   41   77   80   52   98   99   74   81   78

Not fertilized ($y$): 65   88   55   124   66   72   96   71

Test the hypothesis $H0$:  Mean yields are the same with and without fertilizer

$m_x =$         $s_x =$        $N_x =$
$m_y =$         $s_y =$        $N_y =$

$z =$            $p =$