## Comparing Two Independent Samples

The sampling distribution for ...

- $\overline{X}_1$  is approximately Normal with *Mean*= $\mu_1$ ,
- $\overline{X}_2$  is approximately Normal with Mean= $\mu_2$ ,  $SD = \frac{\sigma_2}{\sqrt{n}}$
- $\overline{X}_1 \overline{X}_2$  is approx Normal with  $Mean = \mu_1 \mu_2$ ,  $SD = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

 $SD = \frac{o_1}{\sqrt{n_1}}$ 

$$z_{s} = \frac{\left(\overline{X}_{1} - \overline{X}_{2}\right) - \left(\mu_{1} - \mu_{2}\right)}{SD_{\overline{X}_{1} - \overline{X}_{2}}} \qquad \Rightarrow \qquad t_{s} = \frac{\left(\overline{X}_{1} - \overline{X}_{2}\right) - \left(\mu_{1} - \mu_{2}\right)}{SE_{\overline{X}_{1} - \overline{X}_{2}}}$$

$$H_o: \mu_1 - \mu_2 = D_0$$

$$t_s = \frac{(\overline{x}_1 - \overline{x}_2) - D_0}{SE_{\overline{x}_1 - \overline{x}_2}}$$

$$\sigma_1 = \sigma_2 \quad \Rightarrow \quad t_s = \frac{(\overline{x}_1 - \overline{x}_2) - D_0}{\sqrt{s_\rho^2 \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \text{ with } df = n_1 + n_2 - 2 \quad (\text{Equal Variance T-test})$$

Rejection Region at the  $\alpha$  level of significance:

$$\begin{array}{lll} \mbox{Reject Ho in favor of} & H_a: \mu_1 - \mu_2 > D_0 & if \quad t_s \geq t_\alpha \\ & H_a: \mu_1 - \mu_2 < D_0 & if \quad t_s \leq -t_\alpha \\ & H_a: \mu_1 - \mu_2 \neq D_0 & if \quad |t_s| \geq t_{a/2} \end{array}$$

$$\sigma_1 \neq \sigma_2 \quad \Rightarrow \quad t_s = \frac{(\overline{x}_1 - \overline{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

(Unpooled or Unequal Variance T-test)

Satterthwaite's approximation for the degrees of freedom:





Reduction in cholesterol in mg/dl is measured on day 22

5100		110 011	00001
Bean	15	26.46	5.90
Oat	15	32.23	9.56

Example: 30 mice given one of two diets for 21 days.





Non-parametric Tests for 2 Independent Samples

- Test  $H_o: \tilde{\mu}_1 = \tilde{\mu}_2$  as opposed to  $H_o: \mu_1 = \mu_2$
- Use recoded data: the ranks of the observations as opposed to the original values.
- Do not rely on the assumption of normally distributed populations  $\rightarrow$  Nonparametric

## A T-test on the Ranks

- Rank the combined data from both samples (ties get assigned the average rank).
- Do a 2-sample t-test on the ranks, instead of the observed values.
  - o This test is roughly equivalent to the Wilcoxon Rank Sum Test.

## Wilcoxon Rank Sum Test (a.k.a. Mann-Whitney Test)

- Rank the combined data from both samples (ties get assigned the average rank).
- T =Sum of the ranks in sample 1 (called *W* or *S* in some software).

When H<sub>0</sub> is true, the sampling distribution for T is approx. normal (if  $n_1 > 10 \& n_2 > 10$ ) with

• Mean: 
$$\mu_T = \frac{n_1(n_1 + n_2 + 1)}{2}$$

• Variance: 
$$\sigma_T^2 = \frac{n_1 n_2}{12} (n_1 + n_2 + 1)$$

• Test Statistic: 
$$z_s = \frac{T - \mu_T}{\sigma_T}$$

Rejection Region at the  $\alpha$  level of significance:  $H_a: \tilde{\mu}_1 > \tilde{\mu}_2$  if  $z_s \ge z_{\alpha}$ 

$$H_a: \mu_1 < \mu_2 \quad \text{if} \quad z_s \le -z_\alpha$$

 $H_a: \tilde{\mu}_1 \neq \tilde{\mu}_2 \quad if \quad \mid z_s \mid \geq z_{\alpha/2}$ 

When there are ties in the ranks there is a correction to the variance:

$$\sigma_T^2 = \frac{n_1 n_2}{12} (n_1 + n_2 + 1) - \left(\frac{n_1 n_2}{12} \frac{\sum_{j=1}^k t_j (t_j^2 - 1)}{(n_1 + n_2)(n_1 + n_2 - 1)}\right)$$

JMP:

- (Analyze: Fit Y by X: <red ▼> Nonparametric: Wilcoxon Test)
- When there are more than 2 samples, this is known as the Kruskal-Wallis Test