## "Old Fashioned" Descriptive Statistics – You should come into the course already knowing these

| Statistic | Parameter | Point Estimate | Formula | Interpretation | Notes / Discussion |
|---|---|---|---|---|---|
| Sum of squares | $\sigma^2 \times df$ | SS | $SS = \sum_{i=1}^{n}(x_i - \bar{x})^2$ | No easy interpretation. | • Accompany descriptive statistics w/ EDA when possible. |
| Mean | $\mu$ | $\bar{x}$ | $\bar{x} = \frac{1}{n}\sum_{i=1}^{n} x_i$ | A measure of central location; "expectation" | • Mean and standard deviation are best when distribution is bell-shaped or at least symmetrical. |
| Variance | $\sigma^2$ | $s^2$ | $s^2 = \frac{SS}{n-1}$ | A measure of "spread"; expressed in units squared | • Assuming normal distribution, 68% of data points lie within $\pm1$ standard deviation from the mean, 95% within $\pm2$. For other data use Chebychev's rule (at least 75% of data lie within $\pm2$ standard deviations from the mean). |
| Standard Deviation | $\sigma$ | $s$ | $s = \sqrt{s^2}$ or $\sqrt{\frac{SS}{n-1}}$ | A measure of variability, expressed in data units. More appropriate for descriptive purposes. | |

## 5-point summaries and boxplots - You should come into the course already knowing these

| Statistic | Formula | Interpretation | 5-point Summary | Notes / Discussion |
|---|---|---|---|---|
| Median | Median has depth of $\frac{n+1}{2}$ | A measure of central location | Q0 – Minimum<br>Q1 – First Quartile<br>Q2 – Median<br>Q3 – Third quartile<br>Q4 – Maximum | • When comparing side-by-side boxplots, discuss their locations (look at box, whiskers, and medians), discuss spread (look at IQR's), relative to each other. Also discuss ranges (whisker-spread), overlaps, and outside values. |
| Interquartile Range $(IQR)$ | $IQR = Q3 - Q1$ | A measure of spread, aka "hinge-spread" | | • The box contains 50% of data, drawn from Q1 to Q3, and the height of this box is the IQR (hinge-spread). The line inside the box is Q2 (median). |
| Lower Fence $(F_l)$ | $F_l = Q1 - 1.5(IQR)$ | Use to determine:<br>Lower inside value<br>Lower outside value(s) | | • The lower whisker is drawn from Q1 to the lower inside value. The upper whisker is drawn from Q3 to the upper inside value. Fences are <u>never</u> drawn.<br>• Remember to plot dots for outside values. |
| Upper Fence $(F_u)$ | $F_u = Q3 + 1.5(IQR)$ | Used to determine:<br>Upper inside value<br>Upper outside value(s) | | • Good method to compare groups, esp. when distribution is asymmetrical. |

**Testing MEANS -- always interpret in context of descriptives and EDA; remember fallacies of sig. testing; remember validity assumptions trump distributional conditions**

| *Test* | *Hypothesis* | *Formulas* | *Notes / Discussion* |
|---|---|---|---|
| Equal variance *t* test<br><br>"Student's *t*"<br>"Pooled t"<br>(not recommended in practice) | $H_o : \mu_1 = \mu_2$<br>$H_1 : \mu_1 \neq \mu_2$ | $s_p^2 = \dfrac{(df_1)(s_1^2) + (df_2)(s_2^2)}{df}$<br><br>$se_{\bar{x}_1 - \bar{x}_2} = \sqrt{s_p^2 \left(\dfrac{1}{n_1} + \dfrac{1}{n_2}\right)}$<br><br>$t_{stat} = \dfrac{\bar{x}_1 - \bar{x}_2}{se_{\bar{x}_1 - \bar{x}_2}}$ , $df = df_1 + df_2$ | • A common error is forgetting to square the standard deviations before pooling. Check to see whether variances or standard deviations are given.<br>• When we pool the variances, we suppress non-uniformity of $s_1$ and $s_2$.<br>• Conditions: independent samples, Normality, equal variance |
| Unequal variance *t* test<br><br>("Behrens-Fisher problem;" Welch's solution) | $H_o : \mu_1 = \mu_2$<br>$H_1 : \mu_1 \neq \mu_2$ | Standard error of the mean difference<br><br>$se_{\bar{x}_1 - \bar{x}_2} = \sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}$<br><br>t statistic<br><br>$t_{stat} = \dfrac{\bar{x}_1 - \bar{x}_2}{se_{\bar{x}_1 - \bar{x}_2}}$<br><br>$df' = \dfrac{\left(s_1^2/n_1 + s_2^2/n_2\right)^2}{\dfrac{\left(s_1^2/n_1\right)^2}{n_1 - 1} + \dfrac{\left(s_2^2/n_2\right)^2}{n_2 - 1}}$ | • This *t* test can be used when equal variances are assumed or not assumed.<br>• Consider using this test if you performed a significance test for variances and concluded heteroscedasticity.<br>• Conditions: independent samples, Normality |
| ANOVA | $H_o : \mu_1 = \mu_2 = \mu_3 .....$<br>$H_1 : H_o$ is false<br>(at least one pop. mean differs) | Variance Between groups<br><br>$SS_B = \sum_{i=1}^{k} n_i \left(\bar{x}_i - \bar{x}\right)^2$ , $s_B^2 = \dfrac{SS_B}{df_B}$ , $df_B = k - 1$<br><br>Variance Within groups<br><br>$SS_W = \sum_{i=1}^{k} (n_1 - 1)s_i^2$ , $s_W^2 = \dfrac{SS_W}{df_W}$ , $df_W = N - k$<br><br>F statistic →    $F_{stat} = \dfrac{s_B^2}{s_W^2}$ | • ANOVA doesn't tell you which means differs, so you might need to do post-hoc comparisons.<br>• ANOVA has all the limitations of significance testing.<br>• Conditions: independent samples, Normality, equal variance |

## Post-hoc Comparisons: These tests are performed following ANOVA

| *Test* | *K groups / Assumption* | *Hypothesis* | *Formulas* | *Notes / Discussion* |
|---|---|---|---|---|
| Least Significance Difference (LSD) | 2 groups for each comparison<br><br>Number of comparisons:<br>$m = {_k}C_2$ | Example: for $k = 3$<br>$H_o : \mu_1 = \mu_2$<br>$H_o : \mu_1 = \mu_3$<br>$H_o : \mu_2 = \mu_3$ | $se_{\bar{x}_i - \bar{x}_j} = \sqrt{s_w^2 \left( \dfrac{1}{n_i} + \dfrac{1}{n_j} \right)}$, where $s_w^2$ is obtained from ANOVA<br><br>$t_{stat} = \dfrac{\bar{x}_i - \bar{x}_j}{se_{\bar{x}_i - \bar{x}_j}}$ , $df = N - k$ | • Multiple *t* tests to infer which mean differences are significant.<br>• Bonf. adjusts for the Problem of Multiple Comparisons<br>• Look at summary statistics and EDA to put comparisons in context |
| Bonferroni's Method | Same as LSD Method | Same as LSD Method | For each *P* value obtained using LSD, $p_{Bonf} = p \times m$ | |

## Testing VARIANCES

| *Test* | *K groups / Assumption* | *Hypothesis* | *Formulas* | *Notes / Discussion* |
|---|---|---|---|---|
| F-ratio test | $k = 2$ groups | $H_o$ assumes there's no difference in variances (homocedasticity)<br><br>$H_o : \sigma_1^2 = \sigma_2^2$<br>$H_1 : \sigma_1^2 \neq \sigma_2^2$ | $F_{stat} = \dfrac{s_1^2}{s_2^2}$ or $\dfrac{s_2^2}{s_1^2}$, whichever is larger<br><br>Degrees of freedom<br>$df_1 = n_1 - 1$ , $df_2 = n_2 - 1$ | • When equal variances rejected, consider the following methods to compare means: EDA, descriptive/summary statistics, or unequal variance t-test.<br>• Keep the numerator and the denominator separate. The larger variance goes in the numerator and uses the df for this particular group. |
| Levene's test | $k = 2$ or more groups | $H_o : \sigma_1^2 = \sigma_2^2 = \sigma_3^2 ...$<br>$H_1 : H_o$ is false | Compute with SPSS, no hand calculations ☺ | |

## Scatterplots

| Action | Look for | Additional notes |
|---|---|---|
| Determine X and Y | X (explanatory; independent)<br>Y (response; dependent) | • Label axes with variable names and units<br>• Watch for relations that are curved, U shaped, or otherwise non-linear.<br>• It is difficult to assess strength of association visually (too dependent on the scale of the plot)<br>• Outliers may be important or influential |
| Judge Linearity | Look at scatter plot. Can trend be described with a straight line? | |
| Assess Direction | Look at the direction of the slope (positive, negative, flat) | |
| Look for Outliers | Points outside scatter cloud (i.e., with large residuals) | |

## Correlation

| Statistic | Parameter | Point Estimate | Formula | Notes / Discussion |
|---|---|---|---|---|
| Sum of squares | | SS | $$SS_{xx} = \sum_{i=1}^{n}(x-\bar{x})^2$$ $$SS_{yy} = \sum_{i=1}^{n}(y-\bar{y})^2$$ $$SS_{xy} = \sum_{i=1}^{n}\left[(x-\bar{x})(y-\bar{y})\right]$$ | • $SS_{xx}$ quantifies the spread of variable X<br>  $SS_{yy}$ quantifies the spread of variable Y<br>  $SS_{xy}$ quantifies the extent in which two variables "go together"<br>• The strength of the correlation is in the absolute value of $r$, and <u>depends on the application</u>. General guidelines: weak if $|r|<0.3$, moderate if $0.3<|r|<0.7$, strong if $|r|>0.7$<br>• The coefficient of determination ($r^2$) is the proportion of Y's variability that can be "explained" by changes in X.<br>• Null hypothesis assumes no correlation. If we reject, then we are saying that $r$ is significant, and X and Y are correlated.<br>• Interpretation works best if there is a linear relationship, and inference/testing on $\rho$ assumes that X and Y are bivariate normal |
| Correlation Coefficient | $\rho$ | $r$ | $$r = \frac{SS_{xy}}{\sqrt{SS_{xx}SS_{yy}}}$$ | |
| Testing | $H_0 : \rho = 0$ vs.<br>$H_1 : \rho \neq 0$ | | $se_r = \sqrt{\dfrac{1-r^2}{n-2}}$ , $t_{stat} = r/se_r$<br>$df = n-2$ | |

## Linear Regression

| Statistic | Parameter | Point Estimate | Formula | Notes |
|---|---|---|---|---|
| Slope Coefficient: | β | $b$ | $b = \dfrac{SS_{XY}}{SS_{XX}} = r\dfrac{s_Y}{s_X}$ | • The slope indicates predicted change in Y per unit change in X (key statistic) |
| Intercept Coefficient | α | $a$ | $a = \bar{y} - b\bar{x}$ | • Y-intercept is the value on the Y-axis when X is equal to 0 (needed to anchor the line; not often interpreted. |
| Regression Model | | | $\hat{y} = a + bx$ where $\hat{y}$ is the predicted Y at level $x$, $a$ is the intercept estimate and $b$ is the slope estimate | • Distributional conditions: **L**inearity between X and Y **I**ndependence of each bivariate observation **N**ormality of the residuals **E**qual variance of the residual |
| Confidence Interval for β | | | $b \pm \left(t_{n-2,\,1-\alpha/2}\right)\left(se_b\right)$<br><br>where $se_b = \dfrac{se_{Y\|x}}{\sqrt{SS_{xx}}}$ is the standard error of the slope<br><br>$se_{Y\|x}$ = standard error of the regression = Residual Mean Square<br><br>Compute $se_b$ and/or $se_{Y\|x}$ with SPSS (too tedious to do by hand) | ▪ Validity conditions (good info, good sample, no confounding) trump distributional conditions. |
| Significance Test | | | $H_0$: β = 0 vs $H_0$: β ≠ 0<br><br>Null hypothesis claims population slope = 0 (no association )  Use $t$ test or ANOVA, as described in Lab Workbook | |

A few additional general comments:

- Interpretation starts with understanding what you hope to accomplish. Make your statistics valid and meaningful!
- Confidence intervals are usually *more uses* than significance test because of their ability to estimate effect size.
- Understand each step, not just the conclusion.  Step are: 1. Hypotheses statements, 2.Test statistics, 3. *P* value 4.Conclusion. The *P* value is a measure of evidence against $H_0$ (typical thresholds: 0.10, 0.05, and 0.01) but provides no information about the strength, direction, or importance of the relationship.
- Use computational tools when available. Don't feel compelled to calculate everything by hand.