# Computational Public Health Statistics (HS 267) : Formulas (Part 1)
(Initially compiled by Jane Pham; modified by B. Gerstman; Version of March 6, 2005)

## "Old Fashioned" Descriptive Statistics

| 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. | • When comparing two or more groups, address central locations, variability and sample size(s). |
| 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 symmetrical. |
| Variance | $\sigma^2$ | $s^2$ | $s^2 = \frac{SS}{n-1}$ | A measure of "spread"; expressed in units <u>squared</u> | • Assuming normal distribution, 68% of data points lie within $\pm 1$ standard deviation from the mean, 95% within $\pm 2$, and nearly all data within $\pm 3$ 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. | • For all data, Chebychev's rule, at least 75% of data lie within $\pm 2$ standard deviations from the mean. |

## 5-point summaries and boxplots

| 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 their spread (look at IQR's), relative to each other. Also discuss their 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. NOTE: fences are <u>never</u> formally drawn. |
| Upper Fence $(F_u)$ | $F_u = Q3 + 1.5(IQR)$ | Used to determine:<br>Upper inside value<br>Upper outside value(s) | | • Remember to plot dots for outside values.<br>• Good method to compare group locations and spreads, esp. when distribution is assymetrical. |

# Testing MEANS

| Test | K groups / Assumption | Hypothesis | Formulas | Notes / Discussion |
|------|----------------------|------------|----------|-------------------|
| Independent $t$ –test (Student) | $k = 2$ groups<br><br>Equal variances assumed | $H_o$ assumes there's no difference in means<br><br>$H_o : \mu_1 = \mu_2$<br>$H_1 : \mu_1 \neq \mu_2$ | Pooled estimate of variance<br>$$s_p^2 = \frac{(df_1)(s_1^2) + (df_2)(s_2^2)}{df}$$<br>Standard error of the mean difference<br>$$se_{\bar{x}_1 - \bar{x}_2} = \sqrt{s_p^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$<br>$\underline{t}$ statistic , degrees of freedom<br>$$t_{stat} = \frac{\bar{x}_1 - \bar{x}_2}{se_{\bar{x}_1 - \bar{x}_2}} , \quad df = df_1 + df_2$$ | • A common error is forgetting to square the standard deviations in the equation for pooled estimate of variance. Check to see whether variances or standard deviations are given.<br>• When we pool the variances, we suppress non-uniformity of $s_1$ and $s_2$. |
| Independent t –test (Behrens-Fisher) | $K = 2$ groups<br><br>Equal variances <u>NOT</u> assumed | $H_o$ assumes there's no difference in means<br><br>$H_o : \mu_1 = \mu_2$<br>$H_1 : \mu_1 \neq \mu_2$ | Standard error of the mean difference<br>$$se_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$<br>$\underline{t}$ statistic<br>$$t_{stat} = \frac{\bar{x}_1 - \bar{x}_2}{se_{\bar{x}_1 - \bar{x}_2}}$$<br>degrees of freedom, the integer below $df'$<br>$$df' = \frac{\left( s_1^2/n_1 + s_2^2/n_2 \right)^2}{\frac{\left( s_1^2/n_1 \right)^2}{n_1 - 1} + \frac{\left( s_2^2/n_2 \right)^2}{n_2 - 1}}$$ | • This t-test can also 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. |
| ANOVA | $k = 2$ or more groups<br><br>Equal variances assumed | $H_o$ assumes there's no difference in means<br><br>$H_o : \mu_1 = \mu_2 = \mu_3 \dots$<br>$H_1 : H_o$ is false<br>or<br>(at least one pop. mean differs) | Variance Between groups<br>$$SS_B = \sum_{i=1}^{k} n_i \left( \bar{x}_i - \bar{x} \right)^2 , \quad s_B^2 = \frac{SS_B}{df_B} , \quad df_B = k - 1$$<br>Variance Within groups<br>$$SS_W = \sum_{i=1}^{k} (n_1 - 1)s_i^2 , \quad s_W^2 = \frac{SS_W}{df_W} , \quad df_W = N - k$$<br>F statistic → $\quad F_{stat} = \frac{s_B^2}{s_W^2}$ | • Keep in mind that ANOVA can only reveal whether at least one mean differs or not. It doesn't tell you which means differ, that's why we might need to do post-hoc comparisons. |

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

| Test | K groups / Assumption | Hypothesis | Formulas | Notes / Discussion |
|---|---|---|---|---|
| Least Significance Difference (LSD) | $K = 2$ groups for each comparison<br><br>Number of comparisons:<br>$m = {_k}C_2$ | Ex: 3 groups, 3 comparisons<br>Test 1:<br>$H_o : \mu_1 = \mu_2$ vs. $H_1 : \mu_1 \neq \mu_2$<br>Test 2:<br>$H_o : \mu_1 = \mu_3$ vs. $H_1 : \mu_1 \neq \mu_3$<br>Test 3:<br>$H_o : \mu_2 = \mu_3$ vs. $H_1 : \mu_2 \neq \mu_3$ | Standard error of the mean difference<br>$$se_{\bar{x}_i - \bar{x}_j} = \sqrt{s_w^2 \left( \frac{1}{n_i} + \frac{1}{n_j} \right)},\ \text{where}\ s_w^2\ \text{is}$$<br>obtained from ANOVA<br>$$t_{stat} = \frac{\bar{x}_i - \bar{x}_j}{se_{\bar{x}_i - \bar{x}_j}}, \quad df = N - k$$ | • LSD performs multiple t-tests between two groups each time to infer where the mean difference is significant.<br>• Another alternative (to the Bonf. Method) to adjust for the problem of multiplicity is to apply a more stringent error or $\alpha$ level.<br>• Consider looking at descriptive/summary statistics or using EDA to explain where the difference exists, p values are not always necessary. |
| Bonferroni's Method | Same as LSD Method | Same as LSD Method | For each p-value obtained using LSD or independent t-test, that p-value is adjusted by using the following formula:<br>$$p_{Bonf} = p \times m$$ | |

# 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 statistic<br>$$F_{stat} = \frac{s_1^2}{s_2^2}\ \text{ or }\ \frac{s_2^2}{s_1^2},\ \text{whichever is larger}$$<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 | $H_1 : H_o$ is false | Compute with SPSS, no hand calculations ☺ | |

## Scatterplots

| Action | Examples | | Notes / Discussion |
|---|---|---|---|
| Determine X and Y | X (independent) | Y (dependent) | • X (independent) is often called the <u>predictor</u> variable. We hope to use X to predict changes in Y (regression). <br>• Y (dependent) is often called the <u>response</u> variable, because we are asking the question: will Y's value respond to changes in X? <br>• Label X and Y with variable names and units <br>• Watch out for data that appears to be curved, U shaped, or otherwise non-linear. <br>• Should be able to tell the <u>direction</u> of correlation <br>Positive correlation: high X corresponds to high Y <br>Negative correlation: high X corresponds to low Y <br>No correlation: changes in X don't affect Y (horizontal line or vertical cluster) <br>• <u>AVOID</u> trying to assess <u>strength</u> of correlation (weak/strong) as it is too dependent on the scale of the plot. <br>• Outliers may be important, or they may just show an error in data |
| | Sodium Intake <br>Cigarette Consumption <br>Incoming ACT Scores <br>Age <br>% Births Attended by Pro <br>Fluoride Concentration <br>% Reduced Fee Lunches <br>Weight | Systolic Blood Pressure <br>Lung Cancer Mortality <br>% Graduating <br>Height <br>Maternal Mortality <br>Cavities <br>% Bicycle Helmet Use <br>Eggs Produced | |
| Judge Linearity | See if the scattered points can be described by a straight line | | |
| Assess Correlation (Direction) | Look at the direction of the slope | | |
| Look for Outliers | Points that are outside scatter cloud | | |

## 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>$H_1: \rho \neq 0$ | | $$se_r = \sqrt{\frac{1-r^2}{n-2}}, \quad t_{stat} = r/se_r$$ $$df = n-2$$ | |

# Linear Regression

| Statistic | Parameter | Point Estimate | Formula | Notes / Discussion |
|---|---|---|---|---|
| Slope Coefficient: | $\beta$ | $b$ | $b = \dfrac{SS_{xy}}{SS_{xx}}$ | • The slope indicates a change in Y per unit change in X.<br>• Y-intercept is the value on the Y-axis when X is equal to 0. |
| Y-Intercept Coefficient | | $a$ | $a = \bar{y} - b\bar{x}$ | • Null hypothesis assumes a slope of 0 (no correlation). If we reject, then we are saying that X and Y are linearly related, with a change in Y per unit change in X. |
| *Linear Regression Model for Expected Y* | $\hat{y} = a + bx$ where $\hat{y}$ is the predicted Y at level of $x$, $a$ is the intercept and $b$ is the slope. | | | • Linear regression assumes… |
| *Confidence Interval for* $\beta$ | $b \pm \left(t_{n-2,\,1-\alpha/2}\right)\!\left(se_b\right)$ | | | **L**inearity between X and the expected value of Y<br>**I**ndependence of each observation (each X,Y pair is a single observation) |
| *Hypothesis Testing* | $H_0 : \beta = 0$ vs.<br>$H_1 : \beta \neq 0$ | | $se_b = \dfrac{se_Y}{\sqrt{SS_{xx}}}$ , $t_{stat} = (b-0)/se_b$<br>$df = n - 2$ | **N**ormality of the residuals<br>**E**qual variance of the residual normal distributions (centered at a point on the line for each value of X) |

## Few last comments about hypothesis testing:

- Listing of steps is recommended for beginners:
  1. State hypotheses $\left(H_o \text{ and } H_1\right)$

  2. Error threshold ($\alpha$ level) is used *only* for fixed-level testing, but is discouraged at this level and *not* used in practice.
  3. Calculate test statistic and draw the curve to help determine the $p$ value
  4. Conclusion: report precise $p$ value if available; if not available, report closest approximation (e.g., $p < .001$ or what-ever)
  How do we talk about the results? At the *simplest* level, $p \leq .05$ are considered significant (by convention) in practice; in practice the smaller the $p$ value, the better the evidence against $H_o$ and the more confidence we can put in rejecting the null.

- After performing the test, take a step back and view $p$ in relation to summary statistics and EDA. Make your conclusions meaningful!