

MEANS $t_{crit} = t_{LS/\# \text{ of tails}}(DF)$

Nearly Normal Assumption
 -Data comes from dist. that is uni-modal & symmetric
 -With large enough samples (>40) we are safe to use normal model

One-Sample t-Interval
 When the assumptions and conditions are met, we're ready to find the **confidence interval for the population mean, μ** . The confidence interval is

$$\bar{y} \pm t_{n-1}^* \times SE(\bar{y})$$

One-Sample t-Test for the Mean
 The conditions for the one-sample t-test for the mean are the same as for the one-sample t-interval. We test the hypothesis $H_0: \mu = \mu_0$ using the statistic

$$t_{n-1} = \frac{\bar{y} - \mu_0}{SE(\bar{y})}$$

where the standard error of \bar{y} is: $SE(\bar{y}) = \frac{s}{\sqrt{n}}$

When the conditions are met and the null hypothesis is true, this statistic follows a Student's t-model with $n-1$ degrees of freedom. We use that model to obtain a P-value.

Confidence Interval for the Difference Between Two Means
 When the conditions are met, we're ready to find a **two-sample t-interval** for the difference between means of two independent groups, $\mu_1 - \mu_2$. The confidence interval is

$$(\bar{y}_1 - \bar{y}_2) \pm t_{df}^* \times SE(\bar{y}_1 - \bar{y}_2)$$

where the standard error of the difference of the means is

$$SE(\bar{y}_1 - \bar{y}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

The critical value t_{df}^* depends on the particular confidence level and on the number of degrees of freedom.

Two-Sample t-Test
 When the appropriate assumptions and conditions are met, we test the hypothesis $H_0: \mu_1 - \mu_2 = \Delta_0$

where the hypothesized difference Δ_0 is almost always 0. We use the statistic

$$t = \frac{(\bar{y}_1 - \bar{y}_2) - \Delta_0}{SE(\bar{y}_1 - \bar{y}_2)}$$

where the standard error of the difference of the means is

$$SE(\bar{y}_1 - \bar{y}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

The critical value t_{df}^* depends on the particular confidence level and on the number of degrees of freedom.

Pooled t-Test and Confidence Interval for the Difference Between Means
 The conditions for the **pooled t-test** for the difference between the means of two independent groups are the same as for the two-sample t-test, with the additional assumption that the variances of the two groups are the same. We test the hypothesis

$$H_0: \mu_1 - \mu_2 = \Delta_0$$

where the hypothesized difference Δ_0 is almost always 0, using the statistic

$$t = \frac{(\bar{y}_1 - \bar{y}_2) - \Delta_0}{SE_{pooled}(\bar{y}_1 - \bar{y}_2)}$$

The standard error of $\bar{y}_1 - \bar{y}_2$ is

$$SE_{pooled}(\bar{y}_1 - \bar{y}_2) = s_{pooled} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

where the pooled variance is

$$s_{pooled}^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{(n_1 - 1) + (n_2 - 1)}$$

When the conditions are met and the null hypothesis is true, we can model this statistic's sampling distribution with a Student's t-model with $(n_1 - 1) + (n_2 - 1)$ degrees of freedom. We use that model to obtain a P-value for a test or a margin of error for a confidence interval.

The corresponding **pooled t confidence interval** is

$$(\bar{y}_1 - \bar{y}_2) \pm t_{df}^* \times SE_{pooled}(\bar{y}_1 - \bar{y}_2)$$

where the critical value t_{df}^* depends on the confidence level and is found with $(n_1 - 1) + (n_2 - 1)$ degrees of freedom.

Paired t-Test
 When the conditions are met, we're ready to test whether the mean paired difference is significantly different from a hypothesized value (called Δ_0). We test the hypothesis

$$H_0: \mu_d = \Delta_0$$

where the d 's are the pairwise differences and Δ_0 is almost always 0. We use the statistic

$$t_{n-1} = \frac{\bar{d} - \Delta_0}{SE(\bar{d})}$$

where s_d is the standard deviation of the pairwise differences.

Paired t-Confidence Interval
 When the conditions are met, we're ready to find the confidence interval for the mean of the paired differences. The confidence interval is

$$\bar{d} \pm t_{n-1}^* \times SE(\bar{d})$$

where the standard error of the mean difference is $SE(\bar{d}) = \frac{s_d}{\sqrt{n}}$

The critical value t^* from the student's t-model depends on the particular confidence level you specify and on the degrees of freedom, $n-1$, which is based on the number of pairs, n .

Wilcoxon
 $H_0: M_d = 0$ $H_a: M_d < 0$
 $H_0: M_d = 0$ $H_a: M_d \neq 0$

-Dependent Samples
 -Sample of Diff is NOT Normally Distributed

Mann-Whitney
 $H_0: M_1 - M_2 = 0$ $H_a: M_1 - M_2 > 0$
 $H_0: M_1 - M_2 = 0$ $H_a: M_1 - M_2 \neq 0$

-Independent Samples
 -One or both of the samples is NOT normally distributed

Required Sample Size (always round up)
 -proportion \hat{p}, \hat{q} and E are always entered as decimals. If \hat{p} and \hat{q} are not given use values of 0.5

-means (σ known) $n = Z_{\alpha/2}^2 \frac{\sigma^2}{E^2}$
 -means (σ unknown) $n = T_{\alpha/2}^2 \frac{s^2}{E^2}$

Binomial
 $H_0: p = 0.3$ $H_a: p < 0.3$
 $H_0: p = 0.4$ $H_a: p \neq 0.4$
 $P(X=x) = nC_r p^r q^{n-r}$
 n = # of trials (sample size)
 p = probability of a success
 q = probability of a failure
 $q = 1 - p$
 x = # of successes in "n" trials
 Reject H_0 if $p < \alpha$

The Kruskal-Wallis test
 The Kruskal-Wallis test extends the Wilcoxon rank-sum (Mann-Whitney) test in order to deal with more than two samples.

H_0 : All the samples come from the same distribution.
 H_a : At least one of the samples comes from a distribution that is shifted higher or lower than the others.

HYPOTHESIS TESTING	Two-Sided	One-Sided RIGHT	One-Sided LEFT	Decision or Conclusion	State of Nature H_0	State of Nature H_1
STEP 1: State Hypothesis for Parameter H_0 : Null Hypothesis H_a : Alternate Hypothesis	How we write it in this book No change, i.e., $H_0: \mu = \mu_0$ $H_a: \mu \neq \mu_0$	$H_0: \mu = \mu_0$ $H_a: \mu > \mu_0$	$H_0: \mu = \mu_0$ $H_a: \mu < \mu_0$		H_0	Correct Decision
STEP 2: Calculate the Test Statistic	How some people write it to spell out the details Practical example	Is the proportion of "up" days on the stock market different from the proportion of "down" days? 0.5	Is there a home field advantage? 0.5		H_1	Type I Error: E_1 $P[E_1] = \alpha$
STEP 3: Calculate the Critical Value					H_0	Type II Error: E_2 $P[E_2] = \beta$
STEP 4: Compare & Conclusion -Test Statistic vs Critical Value -P-value vs Level of Significance					H_1	Correct Decision

Confidence Interval for the Difference Between Two Proportions
 The confidence interval for the difference between two proportions is

$$(\hat{p}_1 - \hat{p}_2) \pm z^* \times SE(\hat{p}_1 - \hat{p}_2)$$

where z^* is the critical value and

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

Binomial Formula
 $P(x) = \frac{n!}{(n-x)!x!} p^x q^{n-x}$

Two-Proportion z-Test
 Testing whether the difference between two proportions is equal to a given number, K.
 In order to test $H_0: p_1 - p_2 = K$
 $H_a: p_1 - p_2 \neq K$
 we calculate the test statistic:
 $z = \frac{\hat{p}_1 - \hat{p}_2 - k}{SD(\hat{p}_1 - \hat{p}_2)}$

Two-Proportion z-Test for equal proportions
 Testing whether two proportions are equal.
 In order to test $H_0: p_1 - p_2 = 0$
 $H_a: p_1 - p_2 \neq 0$
 we calculate the test statistic:
 $z = \frac{\hat{p}_1 - \hat{p}_2}{SD(\hat{p}_1 - \hat{p}_2)}$

where $SD(\hat{p}_1 - \hat{p}_2) = \sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}$
 and $\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$ and $\hat{q} = 1 - \hat{p}$.

CHI SQUARE
Assumptions & Conditions
 -Counted Data Condition: The data must be counted for the categories of the categorical variable
 -Independence Assumption & Randomization Condition
 -Sample Size Condition \rightarrow Expected Cell Frequency ≥ 5

GOF VS Homogeneity (Independence)
 GOF: $df = r - 1$
 H_0 : The distribution is the same
 H_a : The distribution is different

How to Find Expected Values

In a contingency table, to test for homogeneity, we need to find the expected values when the null hypothesis is true. To find the expected value for row i and column j , we take

$$Exp_{ij} = \frac{\text{Total}_{i.} \times \text{Total}_{.j}}{\text{Table Total}}$$

Here's an example:
 Suppose we ask 100 people, 40 men and 60 women, to name their magazine preference. Sports Illustrated, Cosmopolitan, or The Economist, with the following result:

	Sports Illustrated	Cosmopolitan	Economist	Total
Men	25	5	10	40
Women	10	45	5	60
Total	35	50	15	100

Then, for example, the expected value under homogeneity for Men who prefer The Economist would be

$$Exp_{13} = \frac{40 \times 15}{100} = 6$$

Performing similar calculations for all cells gives the expected values:

	Sports Illustrated	Cosmopolitan	Economist
Men	14	20	6
Women	21	30	9

ANOVA - Analysis of Variance

One Way Analysis of Variance
 H_0 : All means are equal
 H_a : At least one mean differs from the others

Source	SS	DF	MS	F
Model	SST	I-1	MST	MST/MSE
Error	SSE	N-I	MSE	
Total	SSTotal	N-1		

I = # of sample groups, N = sum total of all sample sizes

ANOVA-BLOCKED Assumptions
 -See 1-Way ANOVA
 **Review BLOCKED ANOVA table in the book for detailed calculations.
2 Possible Tests
 $df_{block} = b - 1$ $df_{fac} = C - 1$ $df_{err} = (b-1)(C-1)$ $df_{tot} = N - 1$

The Model
 $Y_{ijk} = \mu + \alpha_i + \beta_j + \epsilon_{ijk}$

Example
 Suppose we are interested in testing the impact of gasoline types (regular, premium, ethanol) on gas mileage.
 Candidate for one way ANOVA as we have one factor - gasoline type - and one response variable - gas mileage.
 But a potentially confounding factor that impacts gas mileage is vehicle type.
 If we did a completely randomized design, our samples for each factor may have different types of vehicle which would bring into question any conclusion we might draw.
 One way of controlling for this is to divide the population of cars up by size and randomly assign one car from each size to each gasoline type (essentially creating a 2nd factor).
 This way we insure that the size of the car does not impact on the outcome.
 This is called a **randomized block design**.

Two Way Analysis of Variance (ANOVA)
A Test for Interaction

The first step is to see if there is any interaction between the two factors

H_0 : No interaction between factors A and B
 H_a : A and B interact

If H_0 is true then the individual measurements are modeled as

$$X_{ijk} = \mu + \alpha_i + \beta_j + \epsilon_{ijk}$$

Deriving a Test Statistic
 Our test statistic is therefore

$$F = \frac{MSAB}{MSE} \sim F_{(a-1)(b-1), ab(r-1)}$$

If we have determined that there is no significant interaction we can test for the main effects of each separately...

Step 3: If Interaction is not Present, then with caution, do the further analysis on the Factor Level Means to find out if they are all Equal, or some of them not Equal.

Meaningful Main Effects?

The main effects in our example were significant while the interaction term was not significant.

If the interaction had been significant then the main effects tests are less easily interpretable as they no longer represent the entire impact of each factor.

Source of Variation	Degrees of Freedom	Sum of Squares	Mean Sum of Squares	F-statistic
Factor A	a-1	SSA	MSA = SSA/dfa	MSA/MSE
Factor B	b-1	SSB	MSB = SSB/dfb	MSB/MSE
AB interaction	(a-1)(b-1)	SSAB	MSAB = SSAB/dfab	MSAB/MSE
Error	ab(r-1) or N - ab	SSE	MSE = SSE/dfe	
Total	N-1	SSTotal		

Bonferroni used to identify where the actual differences between means are.

Step 1 - Determine number of rows (r), number of columns (c) and number of observations in each samples being compared (n)
 Step 2 - Calculate "J" (the number of possible comparisons) $J = [(r)(c) - 1] / 2$ $s = \sqrt{MSE}$
 Step 3 - Calculate alpha/2J and the df = df_e
 Step 4 - Get t-values from table at a 1 tailed significance level of alpha/2J and df_e (pick closest value)
 Step 5 - Interval = $(x_{bar_1} - x_{bar_2}) \pm t_{\alpha/2J} * s * \sqrt{MSE(1/n_1 + 1/n_2)}$
 critical difference = margin of error = $t_{\alpha/2J} * s * \sqrt{MSE(1/n_1 + 1/n_2)}$

Rule - if the interval contains "0", there is no difference between the means being compared OR - if |actual diff| < |critical diff|, there is no difference between the means

-Used for 1-WAY ANOVA and 2-WAY ANOVA, do not use for BLOCKED ANOVA
 -For 2-WAY ANOVA there are 3 types of possible comparisons, row comparisons, column comparisons and cell comparisons, all steps will be exactly the same except step 1. The values of "r", "c" "n" are determined based on the type of comparison.

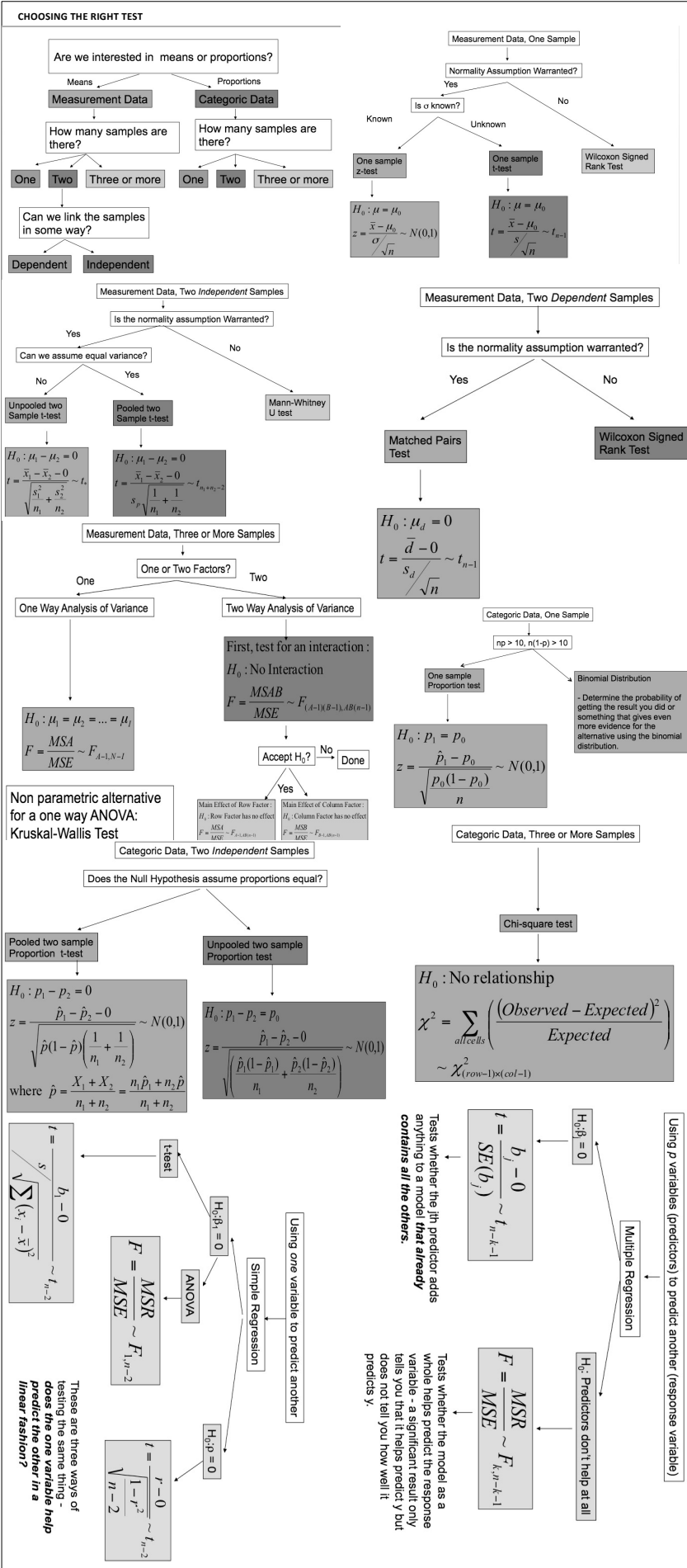
The Chi-Square Calculation
 Here are the steps to calculate the chi-square statistic:

- Find the expected values. These come from the null hypothesis model. Every null model gives a hypothesized proportion for each cell. The expected value is the product of the total number of observations times this proportion. (The result need not be an integer.)
- Compute the residuals. Once you have expected values for each cell, find the residuals, $Obs - Exp$.
- Square the residuals. $(Obs - Exp)^2$
- Compute the components. Find $\frac{(Obs - Exp)^2}{Exp}$ for each cell.
- Find the sum of the components. That's the chi-square statistic, $\chi^2 = \sum_{all \ cells} \frac{(Obs - Exp)^2}{Exp}$.
- Find the degrees of freedom. It's equal to the number of cells minus one.
- Test the hypothesis. Large chi-square values mean lots of deviation from the hypothesized model, so they give small P-values. Look up the critical value from a table of chi-square values, such as Table X in Appendix B, or use technology to find the P-value directly.

The steps of the chi-square calculations are often laid out in tables, as in Table 16.3. Use one row for each category, and columns for observed counts, expected counts, residuals, squared residuals, and the contributions to the chi-square total:

	Observed	Expected	Residual - (Obs - Exp)	Component $\frac{(Obs - Exp)^2}{Exp}$
Monday	192	193.369	-1.369	0.0097
Tuesday	189	202.582	-13.582	0.9106
Wednesday	202	203.695	-1.695	0.0141
Thursday	199	200.607	-1.607	0.0129
Friday	218	199.747	18.253	1.6880

Table 16.3 Calculations for the chi-square statistic in the trading days example.



REGRESSION

Recognizing Violations of the Assumptions

- Non-linearity: The residual plot will have a curve to fit.
- Non-constant variance (heteroscedasticity): The spread of the residuals will change with the fitted value.
- Non-normality of the errors:
 - If there are more outliers than would be expected by the empirical rule then by removing them.
 - If there is skewness, try variance-stabilizing transformations.

The regression equation is $\hat{y} = 42.9 + 0.325 \text{ midr}$

Predictor	Coef	SE Coef	T	P
Constant	42.907	4.087	10.50	0.000
midr	0.32528	0.06942	5.38	0.000

$R^2 = 11.4724$, $R\text{-Sq} = 20.3\%$, $R\text{-Sq(Adj)} = 19.6\%$

ANOVA test to determine if slope of the population regression line is non-zero.

Source	DF	SS	MS	F	P
Regression	1	3815.1	3815.1	28.99	0.000
Residual Error	114	15004.2	131.6		
Total	115	18819.3			

Analysis of Variance

Unusual Observations

Obs	midr	fx	Fir	SR	Fir	Residual	St Resid
46	77	40.00	67.85	1.27	-27.85	-2.44R	
54	73	33.57	66.76	1.17	-33.19	-2.91R	
59	60	29.29	62.42	1.11	-33.14	-2.90R	
60	74	19.29	67.12	1.20	-47.84	-4.19R	
109	74	42.14	67.12	1.20	-24.98	-2.15R	
111	69	41.43	65.32	1.09	-23.89	-2.09R	
114	18	52.86	48.59	2.06	4.17	0.28 X	
115	51	33.57	59.53	1.37	-25.96	-2.28R	
116	17	44.29	48.33	5.13	-4.04	-0.37 X	

An Example of Data Transformations

- Often if the true relationship between X and Y is non-linear it is possible to "linearize" the data by transforming either the X or the Y or both.
- For example, consider $Y = a_0(a_1)^x$
- If we take the log of both sides we get $\log(Y) = \log(a_0(a_1)^x) = \log(a_0) + X \log(a_1) + \log(\epsilon)$

and thus the log of Y is linearly related to X

Simple Linear Regression

Population Model: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

Y_i = value of the dependent variable on i^{th} observation
 X_i = known value of the independent variable on the i^{th} observation
 ϵ_i = r.v. or random error such that $\mu(\epsilon_i) = 0$ (equally likely to be above as below the regression line)
 $\sigma^2(\epsilon_i) = \sigma^2 \forall i$ (constant variance)
 ϵ_i 's are independent, normally distributed

We want to estimate β_0 and β_1 to get the "best" regression line possible.

Idea: Have an independent variable X that you believe can help predict some other variable Y. What we're going to do is estimate β_0 and β_1 in order to get an estimate of the regression line

$$Y_i = \beta_0 + \beta_1 X_i$$

This is an estimate of what the true regression line would be if we could measure the whole population. This in turn is an approximation of the true interaction between X and Y.

Multiple Regression

- a single dependent variable Y
- multiple independent predictor (or explanatory) variables X_1, X_2, \dots, X_k

The Model

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik} + \epsilon_i$$

Assumptions are the same as before:

- Y_i 's are normally distributed, independent, and have constant variance σ^2 .
- Y is linearly related to the predictors

Two Tests

- We used Minitab to determine the appropriate sample regression line $y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$
- Minitab also gives 2 types of tests
 - F-test (ANOVA table) that tests the complete model
 - T-tests (one for each predictor) that tests the usefulness of each predictor given that all the other predictors are already in the model.

Hypothesis Testing - Is there a Linear Relationship?

$H_0: \beta_1 = 0$ What does this mean?

- no linear relationship between X and Y
- X is not a useful linear predictor of Y

Test Statistic: $t = \frac{b_1 - 0}{SE(b_1)}$

$H_a: \beta_1 > 0$ Reject if $t > t_{\alpha, n-2}$
 $H_a: \beta_1 < 0$ Reject if $t < -t_{\alpha, n-2}$
 $H_a: \beta_1 \neq 0$ Reject if $t > t_{\alpha/2, n-2}$ or $t < -t_{\alpha/2, n-2}$

Hypothesis Test

$H_0: \beta_1 = 0$ (X has no explanatory info for predicting Y)

Test Statistic: $F = \frac{MSR}{MSE}$ has a F-distribution with 1 and $n-2$ degrees of freedom - if H_0 is true.

Reject H_0 if $F > F_{\alpha, 1, n-2}$ or use P-values. (Note: $F_{1, n-2} = t_{n-2}^2$)

ANOVA Table

Source	SS	df	MS	F
Model	SSR	k	MSR	MSR/MSE
Error	SSE	n-k-1	MSE	
Total	SSTotal	n-1		

F-Test (ANOVA)

- The F-test looks very much like an ANOVA and compares the full model against a model that includes only the constant $y = \beta_0$
- In other words, it is testing the hypothesis $H_0: \beta_1 = \dots = \beta_k = 0$
 $H_a: \text{At least one } \beta_i \text{ is not zero}$

Two Key Intervals

- C.I. for the mean of Y at a given value of X
- Prediction interval (P.I.) for a future observation of Y at a given value of X

For a given x^* , which interval will be bigger?

- P.I. since the mean varies less than a given observation.

Same point estimate for both: $\hat{y}_x = b_0 + b_1 x^*$

C.I. for $\mu(Y|x^*)$: $(b_0 + b_1 x^*) \pm t_{\alpha/2, n-2} \times s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$

P.I. for $Y|x^*$: $(b_0 + b_1 x^*) \pm t_{\alpha/2, n-2} \times s \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$

T-test (for each predictor)

- The T-test (for the first predictor) is testing the full model $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$ against a model that contains all the predictors except predictor one $y = \beta_0 + \beta_2 x_2 + \dots + \beta_k x_k$
- Thus, it is testing whether predictor x_1 improves the predictive power of a model that already contains all the others.

Coefficient of Determination

$R^2 = \text{"Coefficient of Determination"}$

= Proportion of total variation explained by the model

$$R^2 = \frac{SSR}{SSTotal} \quad \text{or} \quad 1 - \frac{SSE}{SSTotal}$$

(In simple linear regression, r^2 - the correlation coefficient - is equal to R^2)

Computation of R^2_{adj}

$$R^2_{adj} = 1 - (SSE/[n - (k + 1)]) / (SSTotal/(n - 1))$$

$$= 1 - (MSE) / (SSTotal/(n - 1))$$

$$= 1 - (\text{Error Variance} / \text{Total Variance})$$

Multicollinearity

- **Multicollinearity** occurs when at least one predictor is or is close to being a linear combination of the other predictor variables.
- **Pairwise collinearity** occurs when two predictors are highly correlated (as in previous slide)
- Multicollinearity will affect our interpretation of the coefficients but does not affect our prediction.
- We test using the Variation Inflation Factors (VIF) and say that any VIF bigger than 10 constitutes significant multi-collinearity.
- In such a case, we would most likely get rid of one or more predictors.

$$VIF_j = \frac{1}{1 - R_j^2}$$

Introducing Categorical Data

- Suppose you have a categorical predictor with k levels
- Introduce **k-1 dummy variables** (see slides on qualitative predictors)
- Essentially creating k separate regression lines where the assumption is that the impact of the numeric predictors is independent of the level of the categorical predictor
- If interaction is suspected, need to include interaction predictors as in the example in the slides on qualitative predictors

Multiple Regression or Simple Regression Assumptions

- Error Terms Norm Dist.
- Constant Variance
- Error Terms Ind.
- Linear Model

Test #1-Significance of Model

Ho: $\beta_1 = \beta_2 = \beta_3 = \dots = \beta_k = 0$
 Ha: At least one $\beta_i \neq 0$

$F_{stat} = MS_{reg} / MS_{err}$

Reject Ho if $|F_{stat}| > |F_{crit}|$ or $p < \alpha$.

Test #2-Significance of Variable in Model

Ho: $\beta_i = 0$, $T_{stat} = \frac{b_i - 0}{s_{b_i}}$, $df = df_e$
 Ha: $\beta_i \neq 0$

Reject Ho if $|T_{stat}| > |T_{crit}|$ or $p < \alpha$

-All numbers come from minitab.
 -This test can be repeated for each independent variable.
 CI (coefficient) = $b_i \pm t_{\alpha/2} \times s_{b_i}$

CI = $\hat{fit} \pm T_{\alpha/2} * SE_{fit}$ (interval for average of all observations)
 PI = $\hat{fit} \pm T_{\alpha/2} * SQRT(SE_{fit}^2 + S^2)$ (interval for 1 observation)

-Regression Equation Given in Output as $\hat{y} = c + B_1 x_1 + B_2 x_2 + \dots + B_k x_k$
 -Residual = actual value - value predicted by model by subbing independent variables into the regression equation.

-The best model will be the one with the highest R^2 or R^2_{adj} , the lowest S and the lowest number of ind. variables.

Variance Inflation Factor (VIF)

Multicollinearity is a problem caused by highly correlated independent variables. It is a problem if $VIF > 10$.

$$VIF = 1 / (1 - R^2)$$

- Relatively few predictors, to keep the model simple
- A relatively high R^2 indicating that much of the variability in y is accounted for by the regression model
- A relatively small value of se , the standard deviation of the residuals, indicating that the magnitude of the errors is small
- Relatively small P-values for the F - and t -statistics, showing that the overall model is better than a simple summary with the mean and that the individual coefficients are relatively different from zero
- No cases with extraordinarily high leverage that might dominate and alter the model
- No cases with extraordinarily large residuals, and Studentized residuals that appear to be nearly Normal. Outliers can alter the model and certainly weaken the

Stepwise Regression

1. Find the simple regression model that maximizes $|t_{stat}|$, provided its p-value $< \alpha$;
2. Add a second variable if it has the highest $|t_{stat}|$, provided its p-value $< \alpha$;
3. Having added a variable, consider dropping an existing variable with the smallest $|t_{stat}|$, if its p-value $> \alpha$;
4. Keep adding variables until no new variable meets the p-value $< \alpha$ criterion.

Mallows' Cp Statistic

- We choose the set of p predictors that give the smallest Cp statistic subject to the constraint that C_p is approximately d.

ANOVA and REGRESSION

$R^2 = 1 - (SS_{err} / SS_{tot})$
 $R^2_{adj} = 1 - (MS_{err} / MS_{tot})$
 $s = \text{SQRT}(MSE)$ $MS = \frac{SS}{df}$
 $SS_e = \text{SUM}[(n_i - 1) * s_i^2]$

Manual Simple Regression Calculations:

Slope = $b_1 = r * (s_y / s_x)$ r = coefficient of correlation
 Intercept = $b_0 = \bar{y} - b_1 \bar{x}$ s_x = standard dev. of x
 Regression equation: $\hat{y} = b_0 + b_1 \bar{x}$ s_y = standard dev. of y
 \bar{x} = mean of x values
 \bar{y} = mean of y values

No cases with extraordinarily large residuals, and Studentized residuals that appear to be nearly Normal. Outliers can alter the model and certainly weaken the