

Counting

Frequency Tables, Cross Tabulations
Tables, Percentages

Frequency Tables

- A frequency table presents a count for every value of a variable in your database.
- Once you have collected and organized your data, creating frequency tables can:
 - Summarize concisely the distribution of values
 - Be used for validation: are the counts for any values unexpectedly low or high?
 - Be used to identify data entry errors: nonsense entries, misspellings, and stray marks will all show up as distinct values

Components of a Frequency Table

- Value: Each value is listed. If variables are not nominal, list from lowest (top) to highest.
- Frequency: The count of each value.
- Cumulative frequency: The count of each value, plus the counts of all lower values.
- Total: the sum of all cases in the table.
- Proportion: The frequency of a value divided by the total. Expressed as a number between 0 and 1.
- Cumulative Proportion: The count of each proportion, plus the proportions of all lower values.
- Percentage: Proportion expressed in terms of 100.

Grade Distribution in BUS 80.4W

Grade	$f(x)$	$cf(x)$	$p(x)$	$cp(x)$	$\%(x)$	$c\%(x)$
F	4	4	.16	.16	16%	16%
D	2	6	.08	.24	8%	24%
C	5	11	.20	.44	20%	44%
B	8	19	.32	.76	32%	76%
A	6	25	.24	1.00	24%	100%
Total	25		1.00		100%	

- A frequency table for the grade distribution of a hypothetical class.

Frequency Table-Additional Points

- A frequency table should always have a title and a clearly defined value column.
- The total is symbolized by n .
- Proportion and percentage allow researchers to compare the preponderance of values in two different datasets that have different n .
- The formula for proportion:

$$p(x) = \frac{f(x)}{n}$$

Missing Values

- Some observations in your data may not be complete. For example, a respondent may have skipped an item in a questionnaire. When this happens, it creates missing values.
- When constructing frequency tables using SPSS, you should use numeric data, which provides a valid percent as well as an overall percent.
- A valid percent (or proportion) is the frequency of observations for a given value divided by all non-missing values in the data set.
- Typically, we are interested only in valid percent, not a percentage that includes missing values.

Rankings and Percentiles

- It is important that values listed in the frequency table be listed in ascending order.
- This allows the easy calculation of cumulative statistics and percentile rankings.
- Percentile rank: the value of a variable below which a certain percent of observations fall.
- Median: the 50th percentile.
- Quartile: each 25% percentile.
- $Q_1 = 25^{\text{th}}$ percentile; $Q_2 = 50^{\text{th}}$ percentile; $Q_3 = 75^{\text{th}}$ percentile; $Q_4 = 100^{\text{th}}$ percentile.

Rankings and Percentiles

- If we are wondering what score would have a given percentile rank, we can apply this method, using the “Grade Distribution” table:
 - Find the cumulative percentage for a given value, and the cumulative percentage for the next lowest value; if the percentile falls in between, select.
 - The 50th percentile score is B.
 - The 88th percentile score is A.
 - The 12th percentile score is F.
 - The third quartile score is B.

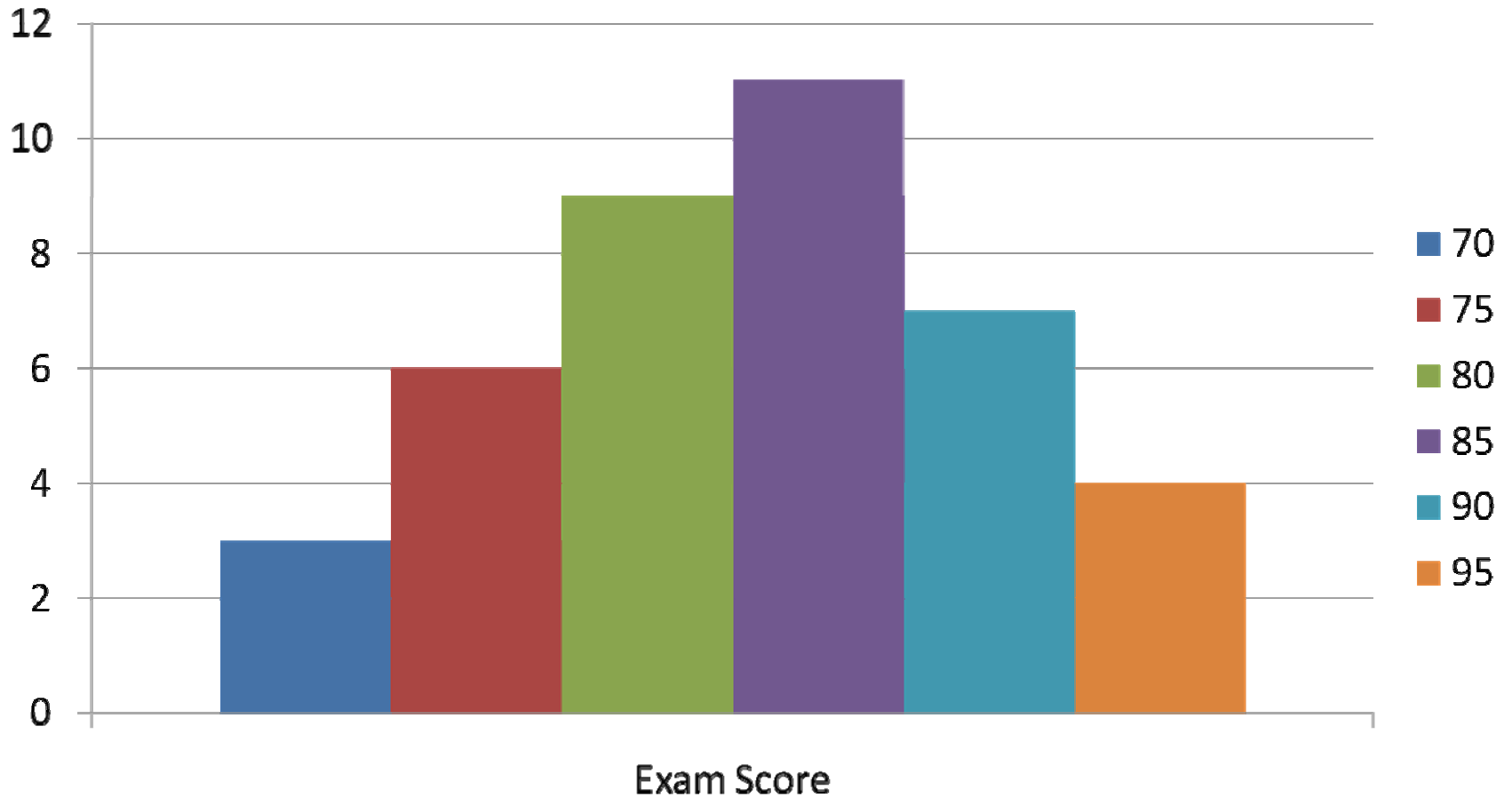
Determining Percentile Ranks

- Percentile Rank Formula:
$$pctrnk = \frac{cf + .5f}{N} * 100\%$$
 - Where *cf* refers to cumulative frequency of the next lowest value; *f* refers to the frequency of the value; and *N* is the total number of cases.
 - This formula allows us to work the opposite way: establishing a percentile rank when all we have are values.

Frequency Distribution Chart

- Frequency can be represented in a graph, where the counts for each value are presented.
- The frequency distribution chart indicates how concentrated (or dispersed) the distribution of scores in your sample is.
- The horizontal axis (the X-axis) lists the values. The vertical axis (the Y-axis) indicates the frequency count for each value. In the example below, the frequency distribution chart takes the form of a bar chart; the taller the bar, the higher the frequency of cases for that value.
- A frequency distribution chart can also take the form of a line graph.

Scores on Midterms, BUS80.4W



Cross Tabulations Tables

- A cross tabulation table displays the joint distribution of two or more variables.
 - Unlike a frequency table, which shows the distribution of values in a single variable.
- Each cell shows a unique combination of the values of two variables.
- Cells going across are called rows; at the end of each row is a row total.
- Cells going down are called columns; at the end of each column is a column total.
- The total number of observations in the table is called the grand total.

Cross Tabulations Tables

	Freshman	Transfer	Continuing	
Agree	15	11	18	44
Neutral	3	13	25	41
Disagree	8	14	35	57
	26	38	78	142

- This table tells us which kinds of students agreed with the statement “Brooklyn College cares about the happiness of its students”.
- 26 is the column total for the freshmen column. 44 is the row total for the “agree” row.
- The sum of the rows and the sum of the columns both equal the grand total, which is 142.

Cross Tabulations Tables

- Cross Tabulations Tables can also be summarized in terms of percentages:
 - Row percentages: the count of values in a cell as a percentage of the row total for that row.
 - Column percentages: the count of values in a cell as a percentage of the column total for that column.
 - Cell or Total percentages: the count of values in a cell as a percentage of the grand total.

Cross Tabulations Tables

	Freshman	Transfer	Continuing	Total
Agree	57.69%	28.95%	23.08%	30.99%
Neutral	11.54%	34.21%	32.05%	28.87%
Disagree	30.77%	36.84%	44.87%	40.14%
	100%	100%	100%	100%

- Above is a cross tabulations table with column percentages. Column percentages should add up to 100% at the bottom of the column.
- Standard practice is to choose the variable we think came first (type of student) as the column variable.
- This allows us to look for how the distribution of the row variable differed for different values of the column variable.
- For this reason, column percentages are preferred.

Central Tendency

Measures of Central Tendency

- Measures of Central Tendency tell us the "middle" or "typical" value of the data set.
 - The word “average” technically means any measure of central tendency; although it’s usually used to refer to one type: arithmetic mean.
- The purpose is to summarize the data with a single number.
- It can also be thought of as the “expected” number, where the degree of difference each number in the distribution is from the mean can be called the “error”.

Measures of Central Tendency

Most-Used Measures of Central Tendency:

- 1) The Mode: The most frequently occurring value in a distribution.
 - The mode can be determined for nominal, ordinal, interval, or ratio level variables.
 - Although the mode can tell a researcher what value is most common, it provides no information about the other values in the distribution.
- 2) The Median: The middle-ranked value in a distribution.
 - The median cannot be determined for nominal variables.
 - Unlike the mode, the median is affected by the other values in the distribution. If some of the lower values are taken out of the sample, the median will rise.
 - However, the median is not sensitive to outliers.
 - **Outliers** are unusually high or unusually low values in the distribution.

Measures of Central Tendency

- 3) The Arithmetic Mean: the sum of all values divided by the total number of cases in the distribution. Formula:

$$\bar{X} = \frac{\sum x}{n}$$

- The mean can only be calculated for interval and ratio level variables.
- Unlike the median, the mean is sensitive to outliers. In other words, a very high or low value can affect the mean, but will not affect the mode.
- This can be helpful: for example, when calculating a student's GPA, we would want to account for the student's best performances and worst performances when summarizing the distribution.
- Sometimes researchers prefer to use the median: i.e., when the outliers may mislead. Household income is an example.
- The mean is also useful in that it provides the basis for more advanced statistics, such as measures of dispersion and regression.

Measures of Central Tendency

4) The Weighted Mean

Any mean that is adjusted somehow can be referred to as a weighted mean. There are a number of common methods of weighting a mean:

- *Dropping Outliers*: typically, a researcher would drop the highest and lowest scores, and calculate the mean from the remaining scores. This is used in the Olympics for judged competitions, where a mean of all ratings is calculated excluding the highest and lowest scores, to avoid the impact of biased judging.

Measures of Central Tendency

- *b-Weighted Mean Formula*: adjusting the mean according to the proportion of cases in the distribution for some key variable.

$$\bar{x} = (x_1 p_1) + (x_2 p_2) + (x_3 p_3) + \textit{etc.}$$

- If the sample is thought not to be representative, researchers may choose to adjust the mean by the existing representation.
- Suppose the sample for a study of GPAs of undergraduates was thought not to be representative according to class standing. Means would be calculated separately for freshmen, sophomores, juniors, and seniors. Each mean would be multiplied by the proportion in the overall population (in this case, the proportion of all enrolled undergraduates), and added together.

Means Tables

- You may want to make comparisons between means in different distributions:
 - Between different variables.
 - Looking at the same variable at different points in time.
 - Looking at the same variable for different values of another variable.
- A means table is an effective way of doing (and presenting this sort of analysis)

Means Tables

This table summarizes mean SAT scores among counties in the NYC area. We can draw conclusions from this about which counties have higher or lower mean SAT scores.

County	Mean SAT
Bronx	990
Kings	1000
Nassau	1020
New York	1040
Queens	1050
Richmond	970
Westchester	1030

Dispersion

Measures of Dispersion

- Measures of Central Tendency tell us what a typical value in a distribution is...but say nothing about *how typical* that value is.
- Measures of Dispersion tell us how spread out values in a distribution are.
- Measures of Dispersion can only be calculated for continuous variables.
- They provide the basis for inferential statistics.

Measures of Dispersion

1. The Range: the difference between the lowest and highest value.
 - This can give us a ballpark figure of where the scores are found, but is not very precise.
 - Defined by the outliers.
2. Interquartile Range: the difference between the first and third quartile value.
 - Can give us a good sense of the spread of the central part of the distribution.
 - Unlike the range, it is not sensitive to outliers.
 - Nevertheless, it misses a lot of variance within the distribution.

Measures of Dispersion

3. The Mean Deviation (or Error).

- The difference between each value and the mean is called the deviation or error.
- We could find the mean of those errors. We might imagine that in a spread out distribution, the mean deviation would be higher, and in a concentrated distribution, it would be lower.
- This does not work, however...so we square the deviations instead, which leads us to more advanced measures.

Measures of Dispersion

4. Sum of Squared Errors (SSE)

- When we square a number, it is always positive.
- So we can add square all of the errors, and add them up. The larger the number is, the greater the deviation.
- If we are comparing two distributions with a different number of cases, the one with the larger number of cases is likely to be larger...even though the true degree of dispersion may not be.
- The SSE is not useful, but is a step in calculating statistics that solve this problem.

Measures of Dispersion

5. Variance: The SSE divided by $n-1$. You can think of it as the average of the squared errors in the distribution.
 - The variance is actually divided by the degrees of freedom ($n-1$) and not the sample size (n). This is because no variation is possible in one case; variation only becomes possible when we have 2 or more cases.
 - The variance is more useful than SSE, but it is expressed in numbers that do not reflect the original distribution.

Measures of Dispersion

6. Standard Deviation: The average deviation of values about the mean.
 - Calculated as the square root of the variance.
 - Avoids the problem of offsetting negative and positive deviations.
 - Utilizes the Sum of Squares concept, which is that the most accurate measures are those for which the average sum of squared errors is least.
 - We can compare two scores in separate distributions in terms of how many standard deviations they are from their mean; thus, the standard deviation is a way of standardizing numerical data.

Formulae

- SSE: $SSE = \sum (x - \bar{x})^2$

- Variance: $s^2 = \frac{\sum (x - \bar{x})^2}{n - 1}$

- Standard Deviation:

$$s = \sqrt{\sum \frac{(x - \bar{x})^2}{n - 1}}$$

Steps

- To find the Standard Deviation, do the following:
 1. Find the mean.
 2. Find the deviation of each value from the mean.
 3. Square each deviation.
 4. Add up the squared deviations (to get SSE).
 5. Divide the SSE by $n-1$ (to get the variance).
 6. Find the square root of the variance (this gives you the standard deviation).

Standardized Scores

Standardized Scores

- Standardized scores express a value in terms of the distribution of values.
 - Tell us where a value falls in relation to the mean and standard deviation of the distribution.
- If the GDP per capita of a country is \$3,200, is this high, medium, or low? We can use standardized scores to tell us where this value fits in relation to the mean.

Standardized Scores

- Standardized scores are also known as z-scores, and calculated using this formula:

$$z = \frac{x - \bar{x}}{s}$$

- Every value in the equation has a z-score (unlike the mean and standard deviation, which are calculated for the whole dist.
- x is the value you want to find a z-score for
- \bar{x} is the mean; s is the standard deviation.

Standardized Scores

- **Z is the number of standard deviations from the mean that a particular score falls.**
- Z-scores serve two main purposes:
 - they help us make comparisons between values in different distributions
 - they help us determine the likelihood that a hypothetical score might be found:
- Z scores greater than 2 and less than -2 are rare.
- The z-score distribution is the same for all continuous variables.

Example

- An example using SAT scores...how can z be used to make comparisons bw 2 distributions?

Year	1980	2000
1	880	910
2	950	940
3	1010	1030
4	1030	1100
5	1080	1150
6	1110	1210
7	1150	1220
8	1190	1290
mean	1050	1106.25
st dev	103.51	137.00

Example

- We can convert all of those scores into z-scores by applying the formula.

Year	1980	2000	1980	2000
1	880	910	-1.642	-1.432
2	950	940	-0.966	-1.214
3	1010	1030	-0.386	-0.557
4	1030	1100	-0.193	-0.046
5	1080	1150	0.29	0.319
6	1110	1210	0.58	0.757
7	1150	1220	0.966	0.83
8	1190	1290	1.353	1.341
mean	1050	1106		
st dev	103.5	137		

Example

- What can we conclude about the distribution?
 - The third lowest score in 1980 was 1010, which is lower than the third lowest score in 2000 (1030).
 - Their z-scores tell us that in terms of the distribution, a 1010 was a better score in 1980.
 - The second highest score in 1980 was 1150; in 2000, the second highest score was 1220.
 - But in terms of the distribution, the 1150 in 1980 was a better score than 1220 was in 2000.

Statistical Significance

Overview

- When we generate a new statistic from a sample, all we have is an estimate of the population.
- Our audience wants to know: how confident can we be that this statistic accurately reflects the population?
- Fortunately, statistics can also be used to assess how close our sample estimate is to the population.

Overview of Terminology

- *Descriptive Statistic*: a statistic calculated to describe the sample or population, such as a mean, proportion, or correlation coefficient.
- *Inferential Statistic*: a statistic calculated to assess how accurately a sample statistic matches the true population parameter. In other words, we try to determine whether we can **infer** that the population parameter is the same, or close to, the sample statistic.
- *Parameter*: a statistic, but exclusively used when talking about the population.

Overview of Symbols

Statistic	Sample	Population
Mean	\bar{x}	μ
Variance	s^2	σ^2
Standard Deviation	s	σ
Standard Error	$s_{\bar{x}}$	$\sigma_{\bar{x}}$

The Logic of Inferential Stats

- We don't know for certain that our sample statistic is a true reflection of the population...by chance, we could have chosen an unusual sample.
- The central limit theorem allows us to estimate the likelihood that our findings are due to chance, given:
 - How different the value is from the value given in our null hypothesis
 - Our sample size
 - The variation of scores in the distribution

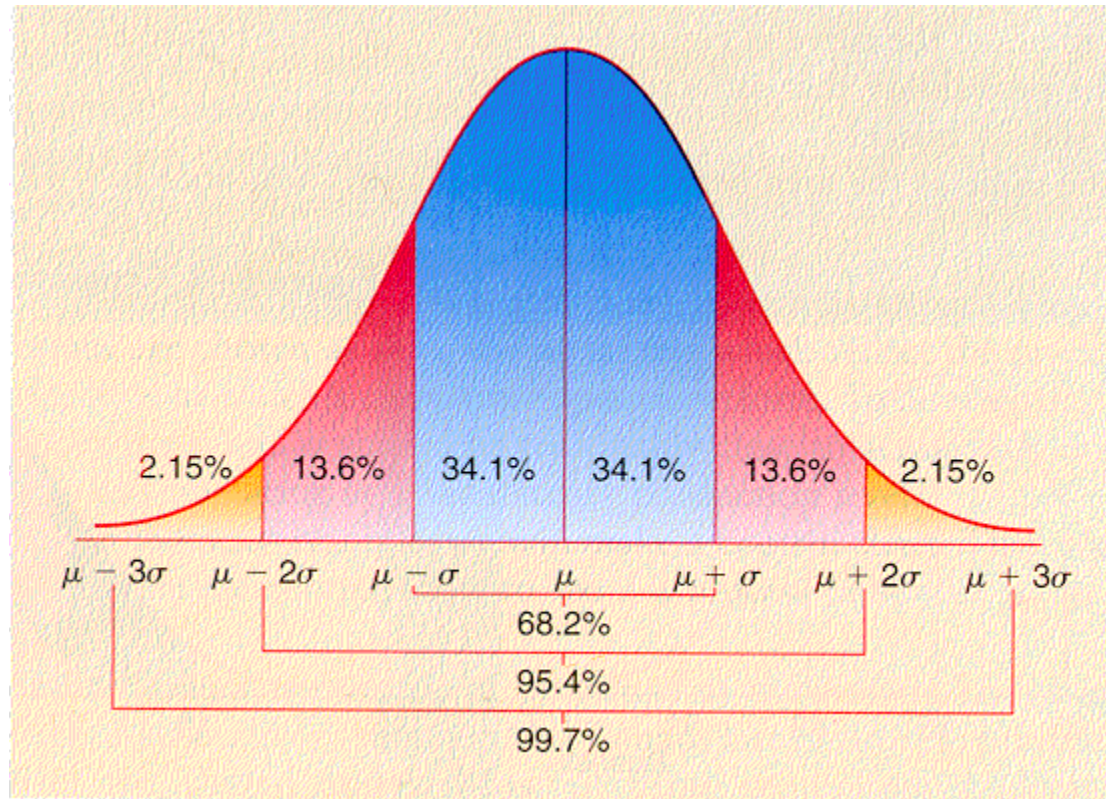
Frequency Distribution

- When we have a set of continuous values, we can graph these values in a frequency distribution graph, where:
 - The x-axis lists all possible values, from 0 (at the midpoint) to the highest value (at the right).
 - The y-axis shows counts for each value; such that a frequently occurring value would get a mark further from the midpoint (vertically).
- We can make marks on the graph for various combinations of value and frequency; these dots can then be summarized by a line.

Normal Distribution

- Most commonly, the frequency distribution resembles a normal distribution (or curve):
 - The most frequent values are the middle-ranked values.
 - The least frequent values are the highest- and lowest-ranked values.
 - The mean, median, and mode are the same.
 - The distribution of values is symmetrical.
 - The percent of values between any two points in the graph are always the same:

Normal Curve



Normal Curve-Explanation

- 50% of all scores fall above the mean, and 50% of all scores fall below the mean.
- 68.2% of all scores fall within one standard deviation above and below the mean.
- 95.4% of all scores fall within two standard deviations above and below the mean.
- 99.7% of all scores fall within three standard deviations above and below the mean.

Normal Curve-Explanation

- Working from reverse, if we know how the percentage of scores between two points, we can determine the z-score:
 - 95% of all scores fall between +1.96 and -1.96.
 - 99% of all scores fall between +2.58 and -2.58.
- If we can assume that values are normally distributed, then we can use z-scores to determine the percent of scores in a distribution above or below a value.

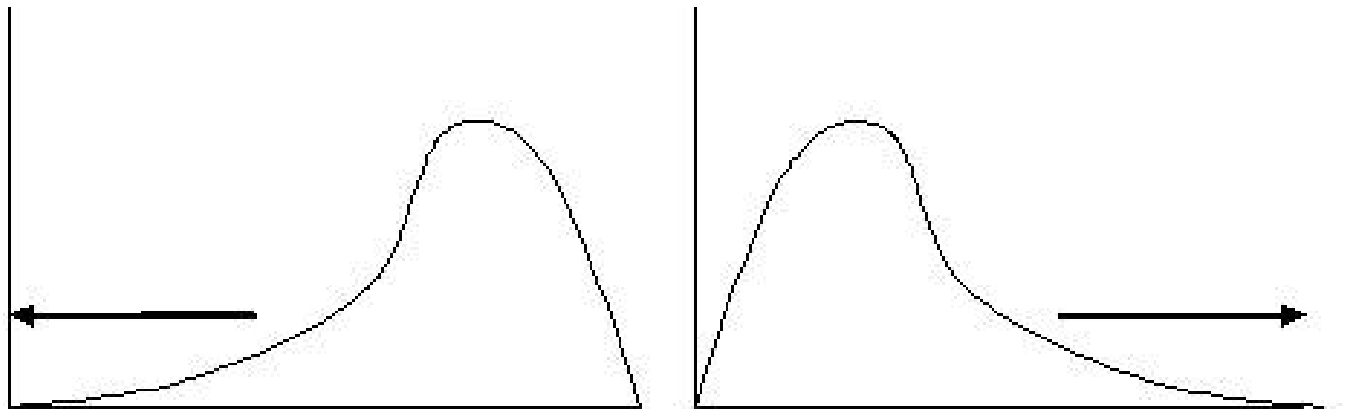
Non-Normal Distributions

- Some examples of variables distributed normally: IQ, by definition, is distributed normally, with 100 set as the mean. Human height has a close to normal distribution. Aptitude tests, such as the SAT, tend to have normal distributions.
- Not all variables are distributed normally. Age, for example, is not...there are generally more young people than old people. Income is not...there are a lot more poor people than rich.

Non-Normal Distributions

- Skewed Distribution: a frequency distribution where very low, or very high scores have the highest frequencies (as opposed to a normal distribution, where middle scores have the highest frequency).
- Positive skew: The right tail is longer; the mass of the distribution is concentrated on the left of the figure. The distribution is said to be right-skewed.
- Negative skew: The left tail is longer; the mass of the distribution is concentrated on the right of the figure. The distribution is said to be left-skewed.

Skewed Distributions



Negative Skew

Elongated tail at the **left**
More data in the left tail than
would be expected in a normal
distribution

Positive Skew

Elongated tail at the **right**
More data in the right tail than
would be expected in a normal
distribution

Central Limit Theorem

- There are three kinds of distributions:
 - Population Distribution: Distribution of all values in the population.
 - Sample Distribution: Distribution of all values in a sample drawn from the population.
 - Sampling Distribution: Distribution of means of samples successively drawn from the population.
 - The mean of the sampling distribution is $\mu_{\bar{x}}$
 - The st. dev. Is called the standard error $\sigma_{\bar{x}} = \frac{s}{\sqrt{n}}$

Central Limit Theorem

- The Central Limit Theorem states that the sampling distribution of any variable will be normal if it has a sufficient sample size (about 30 cases).
- Although we cannot assume that the population and any given sample are normally distributed, we can assume that the sampling distribution will be normally distributed.

Central Limit Theorem

- Assuming that the sampling distribution is normally distributed allows us to also assume:
 - The mean of any sample we drawn is more likely to be close to the mean than far.
 - We can determine the probability of a given value based on its standard error.
 - We can assess the probability that a hypothetical mean is the true population value based on how different it is from the sample mean.

Statistical Significance

- Five basic steps in statistical significance:
 1. State Alternate Hypothesis
 2. State Null Hypothesis
 3. Calculate Sample Statistic
 4. Calculate Inferential Statistic
 5. Draw Conclusion

Stat. Signif., Single Mean

1. Alternate Hypothesis: What you are expecting to find.
 - $H_a: \mu_x \neq 80$...here, my hypothesis is that the mean is not 80.
2. Null Hypothesis: The inverse of the alternate hypothesis.
 - $H_o: \mu_x = 80$...and I will attempt to disprove the notion that the mean is 80.

Stat. Signif., Single Mean

3. Calculate the Descriptive Statistic (the sample mean).
4. Calculate the Inferential Statistic:
 - t-test indicates the z-score, or the distance of the sample mean from the null hypothesis mean expressed in terms of standard errors.
 - p-value indicates the probability of getting such a mean if the null hypothesis is true.
 - SPSS provides #3 and #4 if you go to:
 - Analyze → Descriptive Statistics → Compare Means
→ One Sample t-test
 - Select your variable, and be sure to indicate the test value, which will be 80

In SPSS...

The screenshot displays the SPSS Statistics Data Editor interface. The main window shows a data table with two columns: 'exam1' and 'exam2'. The data points are as follows:

	exam1	exam2
1	50.00	64.00
2	53.00	65.00
3	55.00	68.00
4	60.00	75.00
5	62.00	78.00
6	64.00	79.00
7	67.00	80.00
8	68.00	84.00
9	71.00	86.00
10	72.00	87.00
11	75.00	88.00
12	79.00	89.00
13	81.00	89.00
14	84.00	90.00

The 'One-Sample T Test' dialog box is open, showing the following configuration:

- Test Variable(s): exam1
- Test Value: 80
- Buttons: OK, Paste, Reset, Cancel, Help, Options...

The taskbar at the bottom shows the following open applications: Lectures, Microsoft PowerPoi..., Connect to a network, Untitled1 [DataSet0]..., and *Output3 [Docume...]. The system clock indicates 9:49 PM.

Stat. Signif., Single Mean

One-Sample Statistics

	N	Mean	Std. Deviation	Std. Error Mean
exam1	14	67.2143	10.54790	2.81905

One-Sample Test

	Test Value = 80					
	t	df	Sig. (2-tailed)	Mean Difference	95% Confidence Interval of the Difference	
					Lower	Upper
exam1	-4.535	13	.001	-12.78571	-18.8759	-6.6955

Stat. Signif., Single Mean

5. Conclusion: How do we interpret output?
 - Mean is 67.21. We know it's not 80-but is it significantly different?
 - t is the t-test score, and the higher it is, the further the sample mean is from the hypothesized mean-and the more likely we are to reject H_0 . t is -4.355, which is fairly high (it is negative because our mean is lower than 80).
 - Sig. (2-tailed) is the p-value. It is .001. This means that the probability that we would draw a sample mean this different from the hypothesized mean because of mere chance is about .1%.

Stat. Signif., Single Mean

- Alpha Level: the maximum p-value we are willing to accept to reject H_0 .
- If we choose an alpha level of .05; we would reject H_0 if p-value is less than .05.
 - In this course, assume an alpha level of .05 unless I specify otherwise.
- If we choose an alpha level of .01; we would reject H_0 if p-value is less than .01.
- Here, it's .000, so we would reject H_0 in both circumstances.
- Our sample mean is significantly different from the hypothesized mean of 80.

Type I and Type II Errors

- There are off-setting risks when you choose an alpha level:
- Type I Error (False Rejection): When a true null hypothesis is rejected; or, we think we have statistical significance, but we don't.
 - More of a risk with a .05 alpha level.
- Type II Error (False Accept): When a false null hypothesis is accepted; or, we think we don't have statistical significance, but we do.
 - More of a risk with a .01 alpha level.

Hypothesis Test: Comparison of Means

- Alternate and Null Hypotheses:
 - Ha: $\mu_2 \neq \mu_1$
 - Ho: $\mu_2 = \mu_1$
- Using SPSS to get statistics:
 - Analyze → Descriptive Statistics → Compare Means → Paired Samples t-test
 - Select the two groups or variables you are comparing

In SPSS...

The screenshot displays the SPSS Statistics Data Editor interface. The main window shows a data grid with two columns, 'exam1' and 'exam2', and 27 rows. The first row is highlighted, showing values 50.00 for exam1 and 64.00 for exam2. A 'Paired-Samples T Test' dialog box is open in the center, with 'exam1' and 'exam2' listed in the 'Paired Variables' section. The 'Paired Variables' table is as follows:

Pair	Variable1	Variable2
1	[exam1]	[exam2]
2		

The dialog box also includes an 'Options...' button and 'OK', 'Paste', 'Reset', 'Cancel', and 'Help' buttons at the bottom. The taskbar at the bottom shows the Windows Start button, taskbar icons for 'Lectures', 'Microsoft PowerPoi...', 'Connect to a network', '*Output1 [Docume...', and 'StatisticalSignican...', along with the system tray showing 'SPSS Statistics Processor is ready' and the time '10:03 PM'.

Hypothesis Test: Comparison of Means

Paired Samples Statistics

		Mean	N	Std. Deviation	Std. Error Mean
Pair 1	exam1	67.2143	14	10.54790	2.81905
	exam2	80.1429	14	9.12189	2.43793

The means are 67.21 and 80.14.

Paired Samples Test

		Paired Differences				t	df	Sig. (2-tailed)	
		Mean	Std. Deviation	Std. Error Mean	95% Confidence Interval of the Difference				
					Lower				Upper
Pair 1	exam1 - exam2	-12.92857	3.02462	.80836	-14.67494	-11.18221	-15.993	13	.000

T is -15.993, and the p-value is .000. We again reject Ho.

Measures of Association

Cross Tabulations, Correlations, and
Causation

Multivariate Analysis

- *Univariate*: A statistical analysis consisting of one variable. Frequency distributions, means, and standard deviations are examples.
- *Multivariate*: A statistical analysis that considers the degree of association between two or more variables. Examples: correlation, regression, and path analysis.
- *Bivariate*: A kind of multivariate analysis where there are only two variables in use.

Multivariate Analysis

- Purpose: to establish associations between variables.
- To show that two variables have an association, we must:
 - Show that different values for one variable change from one case to another in a way that corresponds to differences in another variable.
 - Include only cases where a value for both variables has been reported (i.e., no missing values).

Showing Association

- Independent Variables: variables we suspect may have caused change in a dependent variable. Symbolized by X.
- Dependent Variables: variables we suspect may be affected by change in an independent variable. Symbolized by Y.
- Typically in multivariate analysis, we try to show that X causes Y: $X \rightarrow Y$

Forms of Association

- There are a number of ways we can describe a relationship between two variables.
- Significant vs. Insignificant
 - Is there an association at all?
 - A larger sample, and a greater consistency of values varying the same (or opposite way) increases the likelihood that we will find statistical significance.
- Strength of the Association
 - Is this a stronger or weaker association?

Forms of Association

- Nature of the Association?
 - Positive: When cases have higher values of X , we find the same cases to have higher values of Y .
 - Negative: When cases have higher values of X , we find the same cases to have lower values of Y .
- Categorical Variables vs. Continuous Variables
 - If we are using categorical variables, we must either work from a cross tabulations table; or convert to “dummy” variables to employ regression. X^2 , Φ , T_b
 - Continuous Variables: Use r , R^2 , or b (slope).

Forms of Association

- Linearity of Association
 - Linear relationship: X is consistently associated with higher values of Y or with lower values of Y.
 - An association may be non-linear if values of we find high values of X associated with both low and high (but not middle) values of Y.
 - We would call this a non-linear relationship.
 - If this is the case, we cannot use most of our conventional statistics, which are meant to detect linear relationships.

Cross Tabulations Tables

- A cross tabulation table displays the joint distribution of two or more categorical variables.
- How do we measure associations between two categorical variables?
 - Significance: X^2 (chi-squared)
 - Strength of Association: Φ (phi) or T_b (Kendall's Tau-b)
 - Nature of Association: T_b (Kendall's Tau-b)

Cross Tabulations Table

- How can we describe the association between two variables? First, assemble the data into a cross tabulations table:

Time Spent Studying	Female	Male	Total
Less Than 5 Hours/Week	13	11	24
5 to 9.9 Hours/Week	21	13	34
10 or More Hours/Week	25	14	39
Total	59	38	97

Cross Tabulations Table

- The cross tabulations table documents the association between gender (independent) and time spent studying (dependent).
- Do women study more than men? We can immediately tell that
 - Most of the people in the sample are female.
 - Most of the people study 10 hours/week or more.
 - We want to show that even though both genders tend to study more, the concentration of cases for studying more hours is greater for female than for male.

Chi-Square

- Chi-square compares the actual distribution of cases in the cells to the what the distribution would be if there were no relationship.
- The larger the sum of the differences between the observed cell counts and the expected cell counts (expected if there were no association), the greater the probability that the population matches the sample.
- We use hypothesis testing to assess statistical significance, and SPSS generates our statistics.

Chi-Square Hypothesis Testing

1. First, State Hypotheses.

- $H_a: \Phi \neq 0$
- $H_o: \Phi = 0$
- Phi is a measure of strength of association. If it is zero, there is no association. If it is significantly greater than zero, we reject H_o and conclude that there is an association.
- Chi-square is our “test statistic”, like t ; it only tells us that an association is significant.

Chi-Square Hypothesis Testing

2. Analyze the cross tabulations table in SPSS.
 - Go to Analyze → Descriptive Statistics → Crosstabs
 - Select “gender” as your column variable and “hours spent studying” as your row variable.
 - Independent variable should always be in the columns; dependent in the rows.
 - Click the “Statistics” button; you will get a menu of statistics to measure. Check “chi-square”, “Phi and Cramer’s V”, and “Kendall’s Tau-b”
 - Click continue, OK, and view your output.

Chi-Square Hypothesis Testing

The screenshot displays the SPSS Statistics Data Editor interface. The main window shows a data table with two variables: 'gender' and 'studytime'. The data is as follows:

gender	studytime
F	LessThan5Hours
F	LessThan5Hours
F	5to9.9Hours
F	5to9.9Hours
F	10orMoreHours
F	10orMoreHours
M	LessThan5Hours
M	5to9.9Hours
M	5to9.9Hours
M	10orMoreHours
M	10orMoreHours
M	10orMoreHours
F	LessThan5Hours
F	LessThan5Hours
F	LessThan5Hours
F	5to9.9Hours
F	5to9.9Hours
F	10orMoreHours
F	10orMoreHours
F	10orMoreHours
M	LessThan5Hours
M	LessThan5Hours
M	5to9.9Hours
M	5to9.9Hours
M	10orMoreHours
F	LessThan5Hours
F	LessThan5Hours
F	LessThan5Hours
F	5to9.9Hours

Two dialog boxes are overlaid on the data table:

- Crosstabs**: The 'Row(s)' field contains 'studytime' and the 'Column(s)' field contains 'gender'. The 'Display clustered bar charts' checkbox is checked.
- Crosstabs: Statistics**: The 'Chi-square' checkbox is checked. Under the 'Nominal' section, 'Phi and Cramer's V' and 'Kendall's tau-b' are checked. Under the 'Ordinal' section, 'Kendall's tau-c' is checked. The 'Nominal by Interval' section is empty. The 'Cochran's and Mantel-Haenszel statistics' section is also empty.

The taskbar at the bottom shows the following open applications: Lectures, Measures of Associ..., Counting [Compat..., Connect to a netw..., *Untitled1 [DataSet..., and *Output2 [Docume... The system clock shows 9:47 PM.

SPSS Output

studytime * gender Crosstabulation

- Tables:

Count		gender		Total
		F	M	
studytime	10orMoreHours	25	14	39
	5to9Hours	21	13	34
	LessThan5Hours	13	11	24
Total		59	38	97

Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)
Pearson Chi-Square	.635 ^a	2	.728
Likelihood Ratio	.629	2	.730
N of Valid Cases	97		

a. 0 cells (.0%) have expected count less than 5. The minimum expected count is 9.40.

Symmetric Measures

		Value	Asymp. Std. Error ^a	Approx. T ^b	Approx. Sig.
Nominal by Nominal	Phi	.081			.728
	Cramer's V	.081			.728
Ordinal by Ordinal	Kendall's tau-b	.071	.096	.735	.462
N of Valid Cases		97			

a. Not assuming the null hypothesis.

b. Using the asymptotic standard error assuming the null hypothesis.

Chi-Square Hypothesis Testing

3. Draw Conclusion.

- The first table is the cross tabulations table.
- The second table provides chi-square. We read across the row that says: “Pearson Chi-square”. The value is 0.635.
- Reading across, under “Asymp. Sig.” ...this is the p-value. It is .728.
- If the p-value is larger than .05, we cannot reject H_0 .

Strength of Association

- SPSS also reports two measures of strength of association:
 - Phi indicates the strength of association if at least one of the two variables is nominal.
 - If phi is closer to 1, it is a stronger association.
 - If phi is closer to 0, it is a weaker association.
 - If $\phi \geq .70$, it is a strong association.
 - If $\phi \geq .30$ and $< .70$, it is a moderate association.
 - If $\phi < .30$, it is a weak association.
 - Phi can never be negative, because there can be no negative associations among nominal variables.

Strength of Association

- Kendall's Tau-b indicates the strength of the relationship for two ordinal variables.
- Kendall's Tau-b is interpreted the same way as phi, except that if it has a negative sign (for example, $-.23$), we conclude that there is a negative relationship.
- In the SPSS Output, phi is $.081$ (weak). SPSS calculates Kendall's Tau-b, even though it is not appropriate; it is $.071$ (weakly positive).

Correlation

- To show that two continuous variables have an association, we do not need to draw a cross tabulations table.
- Pearson's r determines strength of assoc.
- To interpret Pearson's r :

$r=$	Meaning		$r=$	
0	No Assoc.		0	No Assoc.
.01 to .249	Weak, Positive		-.01 to -.249	Weak, Negative
.25 to .749	Moderate, Positive		-.25 to -.749	Moderate, Negative
.75 to 1.00	Strong, Positive		-.75 to -1.00	Strong, Negative

Pearson's r Hypothesis Testing

1. First, state hypotheses:

- $H_a: \rho \neq 0$
- $H_o: \rho = 0$
- Rho (ρ) refers to Pearson's r.
- We want to show that the association between X and Y is significantly different from 0.
- Our test statistic will be t.

Pearson's r Hypothesis Testing

2. Analyze the correlation in SPSS.

- Go to Analyze → Correlate → Bivariate
- Select the two variables you want to see correlations for...it is not necessary to specify which is independent or dependent.
- Click continue, OK, and view your output.

SPSS Output

Correlations

		advertexpense	revenue
advertexpense	Pearson Correlation	1	.804**
	Sig. (2-tailed)		.000
	N	30	30
revenue	Pearson Correlation	.804**	1
	Sig. (2-tailed)	.000	
	N	30	30

** . Correlation is significant at the 0.01 level (2-tailed).

- In either the top right or bottom left corner, there are three rows of numbers. The number cross from “Pearson Correlation” is Pearson’s r ; it is .804, a strong relationship.
- Across from Sig. (2-tailed) is the p -value.

Pearson's r Hypothesis Testing

3. Draw Conclusion.

- If the p-value is less than .05, we can reject H_0 and accept H_a . The association is statistically significant.
- SPSS also provides symbols to show that it is statistically significant:
 - One star means it is statistically significant at the .05 level; two means it is significant at the .01 level.
- We ignore the boxes where variables are correlated with themselves (always 1.00) and there is always a duplicate correlation, which we also ignore.

Multivariate Analysis and Regression

Multivariate Analysis

- Analysis of relationships with more than two variables.
- To establish cause:
 - Association
 - Time Order
 - Non-spuriousness
- Multivariate analysis is devoted to detecting the impact of other variables, esp. exposing spurious effects.

Establishing Time Order

- *Longitudinal Design*: A study that looks at changes in different points in time.
- *Cross-Sectional Design*: A study that looks at differences between two groups, where the measurements and comparisons take place at one point in time.
- Most marketing research is cross-sectional. This is due to the difficulties, costs, and need issues related to repeatedly doing the same (or similar) research. However, where possible, marketing researchers try to collect longitudinal research.

Cross Sectional Research

- How can we establish time order with cross-sectional data?
 1. Compare by stage in progression. (e.g., current young and old, current freshmen and seniors).
 2. Use non-statistical data to provide context.
 3. Express associations in terms of a theoretical construction.

Longitudinal Research

1. *Trend Studies*: measuring segments of the general public at different points in time. For example, if we were to conduct a survey of the political party affiliation of a randomly selected group of Americans in 2000, and then conduct a survey of the political party affiliation of another randomly selected group of Americans in 2008, this would be a trend study.
2. *Cohort Studies*: measuring changes in a group of people over time who share some common starting point. For example, we could conduct a study of changes in political attitudes of people born in 1946.

Longitudinal Research

3. *Panel Studies*: Measuring changes in the same cases over time.

- E.g., I would interview subjects today; 10 years later, I would contact these exact same subjects again and ask the same or similar questions.
- Ideally, I can compare changes over time for each case in the study.
- In reality, there is the problem of attrition: people who are part of the initial panel, but do not participate at later points in time (due to refusal, death, incarceration, or lost contact info).
- Marketing researchers can reduce this effect by offering incentives for participants who take part in every event in the panel study.
- Nonetheless, significant attrition can severely hamper the results of a panel study, and is a major risk despite the benefits of this design.

Spurious Associations

- A misleading correlation between two variables is produced through the operation of another causal variable.
- So we have three possible relationships:
 - A causes B
 - B causes A
 - **OR** *C causes* both A and B
- Multivariate analysis seeks to show that assoc between A and B is true by *controlling* for various possible C (third) variables.

Partial Correlations

- Measures the degree of association between two variables, with the effect of a set of controlling random variables removed.
- Same logic as bivariate correlation:
 - A score close to +1 is strong and positive.
 - A score close to -1 is strong and negative.
 - A score close to 0 is weak.
- If partial correlation is much smaller than bivariate correlation, we can conclude that third variable C “explains” the original association between A and B.

Partial Correlations

If the bivariate correlation between A and B is...	And the partial correlation between A and B, controlling for C, is...	We conclude that the relationship between A and B, controlling for C, is...
significant	significant	direct
significant	significant, but weaker	partially explained
significant	not significant	spurious
not significant	significant	suppressor

Partial Corr. with SPSS

The screenshot shows the SPSS Statistics Data Editor interface. The 'Analyze' menu is open, and the 'Correlate' option is selected, which has opened a sub-menu where 'Partial...' is highlighted. The main data grid shows a single variable named 'Ebayexpenses' with 15 rows of data. The taskbar at the bottom shows several open applications, including 'Lectures', 'Projections [...]', 'Linear Regre...', 'Correlation [...]', 'Microsoft P...', 'Connect to ...', 'Untitled1 [D...', and '*Output13 [...]'.

	Ebayexpenses
1	200.00
2	125.00
3	65.00
4	24.00
5	145.00
6	89.00
7	67.00
8	79.00
9	92.00
10	46.00
11	21.00
12	48.00
13	77.00
14	104.00
15	188.00
16	
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	

Partial Corr. with SPSS

The screenshot displays the SPSS Statistics Data Editor interface. The main window shows a data table with three variables: Ebayexpenses, Income, and Onlinehours. A dialog box titled "Bivariate Correlations" is open, showing the configuration for a partial correlation analysis. The dialog box includes options for selecting variables, choosing correlation coefficients (Pearson, Kendall's tau-b, Spearman), and testing for significance (Two-tailed, One-tailed). The "Flag significant correlations" option is checked.

	Ebayexpenses	Income	Onlinehours	var	var	var	var	var	var	var	var	var	var	var	var
1	200.00	145000.00	27.00												
2	125.00	189000.00	17.00												
3	65.00	77000.00	12.00												
4	24.00	30000.00	9.00												
5	145.00	110000.00	18.00												
6	89.00	96000.00	10.00												
7	67.00	122000.00	14.00												
8	79.00	74000.00	13.00												
9	92.00	99000.00	16.00												
10	46.00	56000.00	11.00												
11	21.00	91000.00	6.00												
12	48.00	64000.00	9.00												
13	77.00	103000.00	15.00												
14	104.00	116000.00	12.00												
15	188.00	87000.00	19.00												
16															
17															
18															
19															
20															
21															
22															
23															
24															
25															
26															
27															

Bivariate Correlations

Variables:

- Hours Spent Online Per ...
- Monthly Expenditure on ...
- Annual Household Inco ...

Correlation Coefficients:

- Pearson
- Kendall's tau-b
- Spearman

Test of Significance:

- Two-tailed
- One-tailed

Flag significant correlations

Buttons: OK, Paste, Reset, Cancel, Help

Partial Corr. with SPSS

- We find the bivariate correlation first: it is a moderate, positive correlation:

Correlations

		Monthly Expenditure on Ebay	Annual Household Income
Monthly Expenditure on Ebay	Pearson Correlation	1	.581 [*]
	Sig. (2-tailed)		.023
	N	15	15
Annual Household Income	Pearson Correlation	.581 [*]	1
	Sig. (2-tailed)	.023	
	N	15	15

*. Correlation is significant at the 0.05 level (2-tailed).

Partial Corr. with SPSS

The screenshot shows the SPSS Statistics Data Editor interface. The 'Analyze' menu is open, and the 'Correlate' option is selected, which has opened a sub-menu where 'Partial...' is highlighted. The data editor shows a single variable named 'Ebayexpenses' with 15 rows of data. The taskbar at the bottom shows several open applications, including 'Lectures', 'Projections [...]', 'Linear Regre...', 'Correlation [...]', 'Microsoft P...', 'Connect to ...', 'Untitled1 [D...', and '*Output13 [...]'.

	Ebayexpenses
1	200.00
2	125.00
3	65.00
4	24.00
5	145.00
6	89.00
7	67.00
8	79.00
9	92.00
10	46.00
11	21.00
12	48.00
13	77.00
14	104.00
15	188.00
16	
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	

Partial Corr. with SPSS

The screenshot displays the SPSS Statistics Data Editor interface. The main window shows a data table with three variables: Ebayexpenses, Income, and Onlinehours. A dialog box titled "Partial Correlations" is open, allowing the user to specify variables for the analysis.

Data Table:

	Ebayexpenses	Income	Onlinehours	var	var	var	var	var	var	var	var	var	var	var	var
1	200.00	145000.00	27.00												
2	125.00	189000.00	17.00												
3	65.00	77000.00	12.00												
4	24.00	30000.00	9.00												
5	145.00	110000.00	18.00												
6	89.00	96000.00	10.00												
7	67.00	122000.00	14.00												
8	79.00	74000.00	13.00												
9	92.00	99000.00	16.00												
10	46.00	56000.00	11.00												
11	21.00	91000.00	6.00												
12	48.00	64000.00	9.00												
13	77.00	103000.00	15.00												
14	104.00	116000.00	12.00												
15	188.00	87000.00	19.00												
16															
17															
18															
19															
20															
21															
22															
23															
24															
25															
26															
27															

Partial Correlations Dialog Box:

- Variables:** Monthly Expenditure on ..., Annual Household Inco...
- Controlling for:** Hours Spent Online Per ...
- Test of Significance:** Two-tailed One-tailed
- Display actual significance level

The taskbar at the bottom shows the following open applications: Lectures, Projections [...], Linear Regre..., Correlation [...], Microsoft P..., Connect to ..., Untitled1 [D...], and *Output13 [...]. The system clock indicates 10:59 PM.

Partial Corr. with SPSS

- We then calculate the partial correlation; the p-value is now much higher than .05, and is no longer significant.

Correlations

Control Variables			Monthly Expenditure on Ebay	Annual Household Income
Hours Spent Online Per Week	Monthly Expenditure on Ebay	Correlation	1.000	.145
		Significance (2-tailed)	.	.620
		df	0	12
	Annual Household Income	Correlation	.145	1.000
		Significance (2-tailed)	.620	.
		df	12	0

Partial Correlation

- A partial corr. controlling for one variable is called a “first order” partial correlation.
- You can control for many variables at once (second order, third order, etc.)
- Partial correlation is used to test for spuriousness; but it is not useful if we want to collect all of the associations we’ve observed to explain variation in one variable.
- For that, we must use statistical regression.

Regression

- A statistical technique used to model association between a dependent variable and one or more independent variables.
- A “model” is an equation that represents the average association of each independent variable with the dependent variable, controlling for other independent variables.
- This model can be applied such that when we are given a values of the independent variables, we can plug them into the equation and “predict” the value of the dependent variables.

Regression vs. Partial Corr.

- There are some important differences:
 - The results of a regression analysis are expressed in the form of an equation, not a single number (although there are single number summary statistics involved).
 - We must assign one variable the role of dependent variable (or “affected variable”) and the others the role of independent variables (or “causal variables”).
 - An equation can include any number of independent variables; it can only include one dependent variable.
 - This equation can be used to draw a line in a scatterplot, indicating the value of Y, on average, for given values of X.
 - We can use the regression equation to predict values of Y given a value of X.

Regression Equation

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots$$

- **Y** is the dependent variable. We are trying to explain variation in Y, or predict values of Y, given our Xs.
- **X** is the independent variable. A regression equation can actually have multiple independent variables. Each independent variable has a unique slope, and the purpose of the subscripts above is to identify each independent variable.

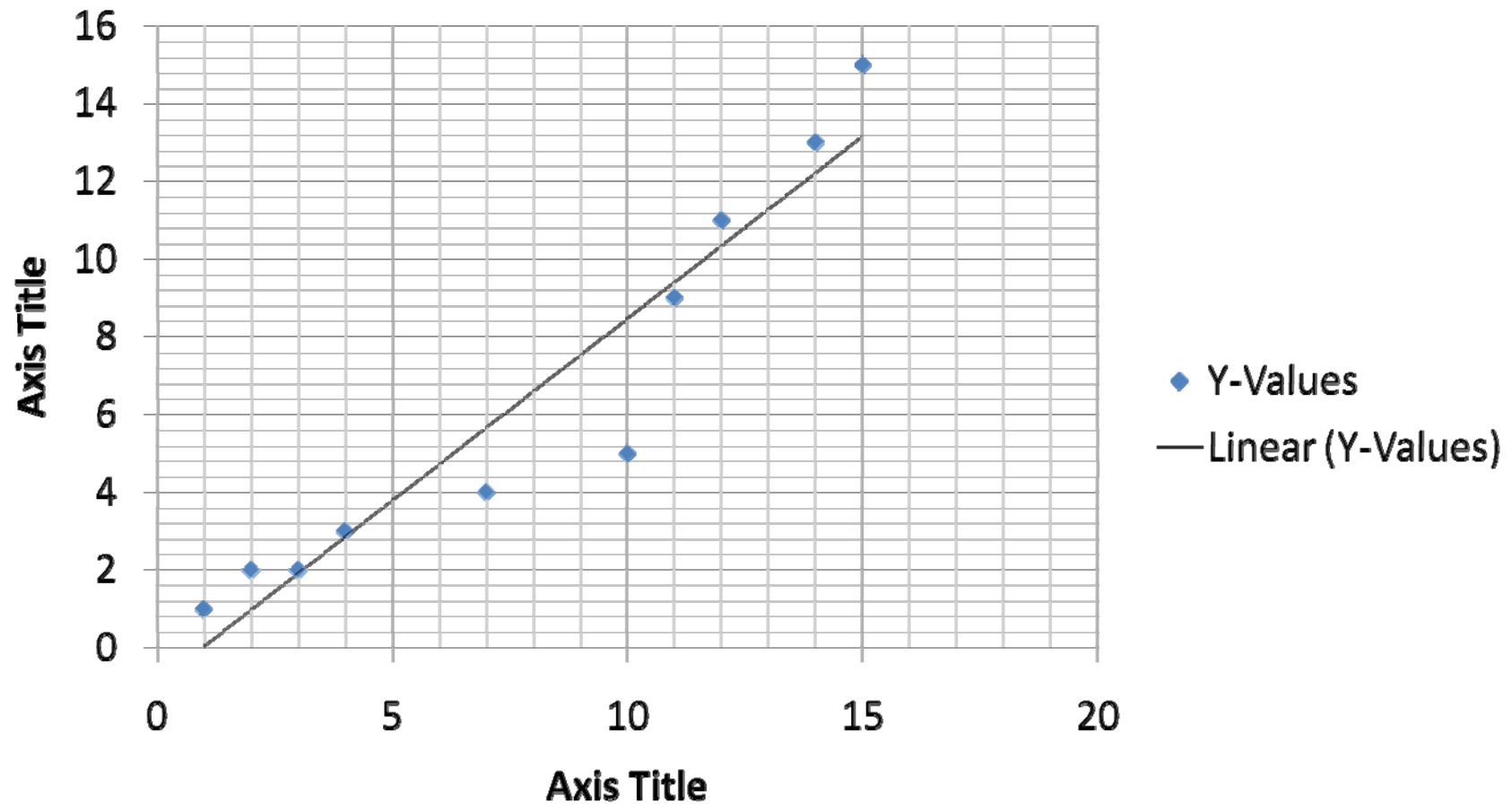
Regression Equation

- α is the Y-intercept. It is the value of Y when the regression line crosses the Y-axis.
- β is the slope for each variable. The slope tells us the average change in Y for every one unit change in X, controlling for all other variables.
- If the slope is positive, an increase in X is associated with an increase in Y.
- If the slope is negative, an increase in X is associated with a decrease in Y.

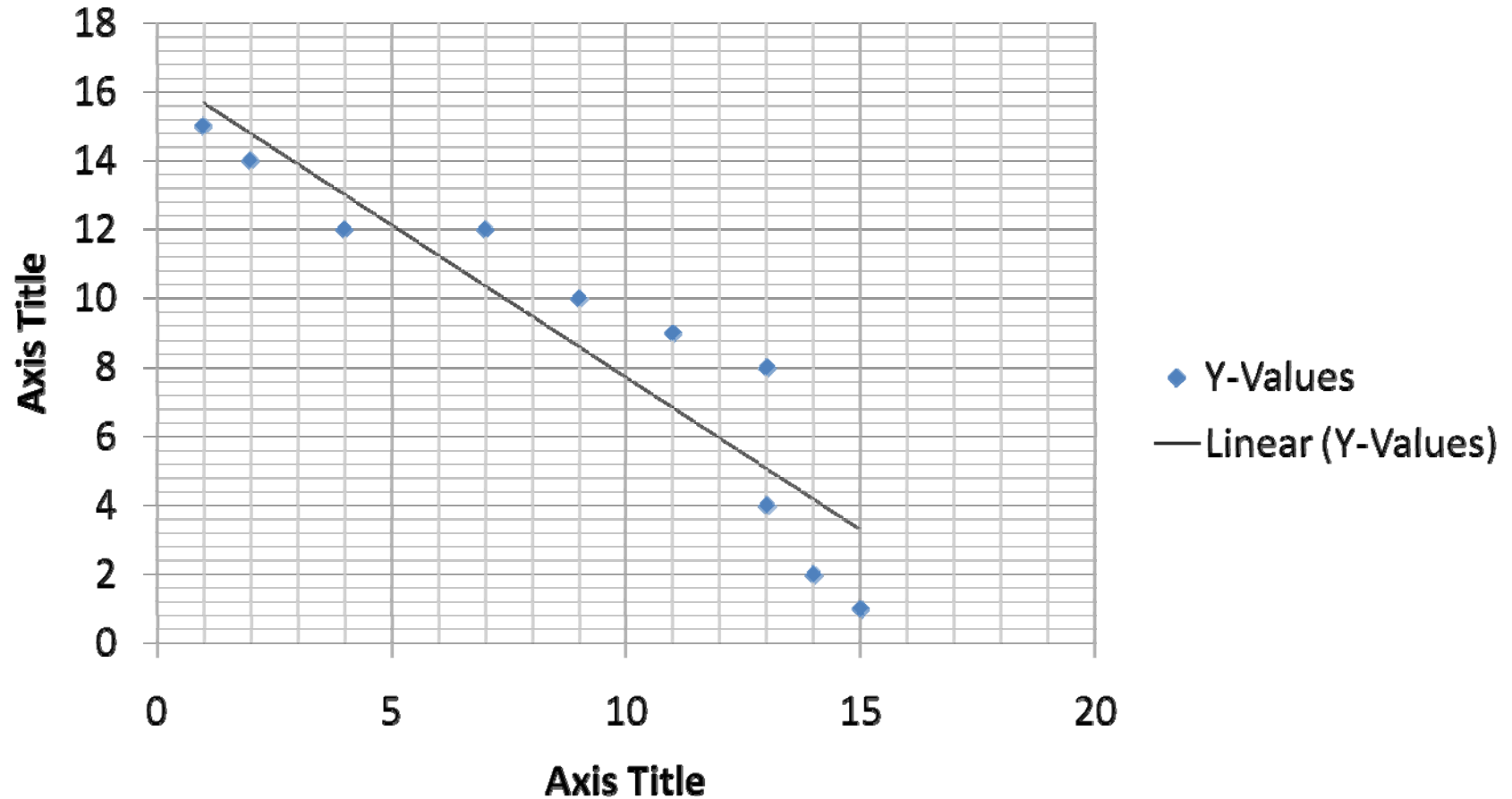
Scatterplot

- A scatterplot plots the values of X and Y for each observation in a graph.
- If there is a statistically meaningful relationship between X and Y, we should see a pattern in how the dots appear in the graph.
- The regression line bisects the dots in the graph...it is the one line that has the least squared deviation (error) from the line.

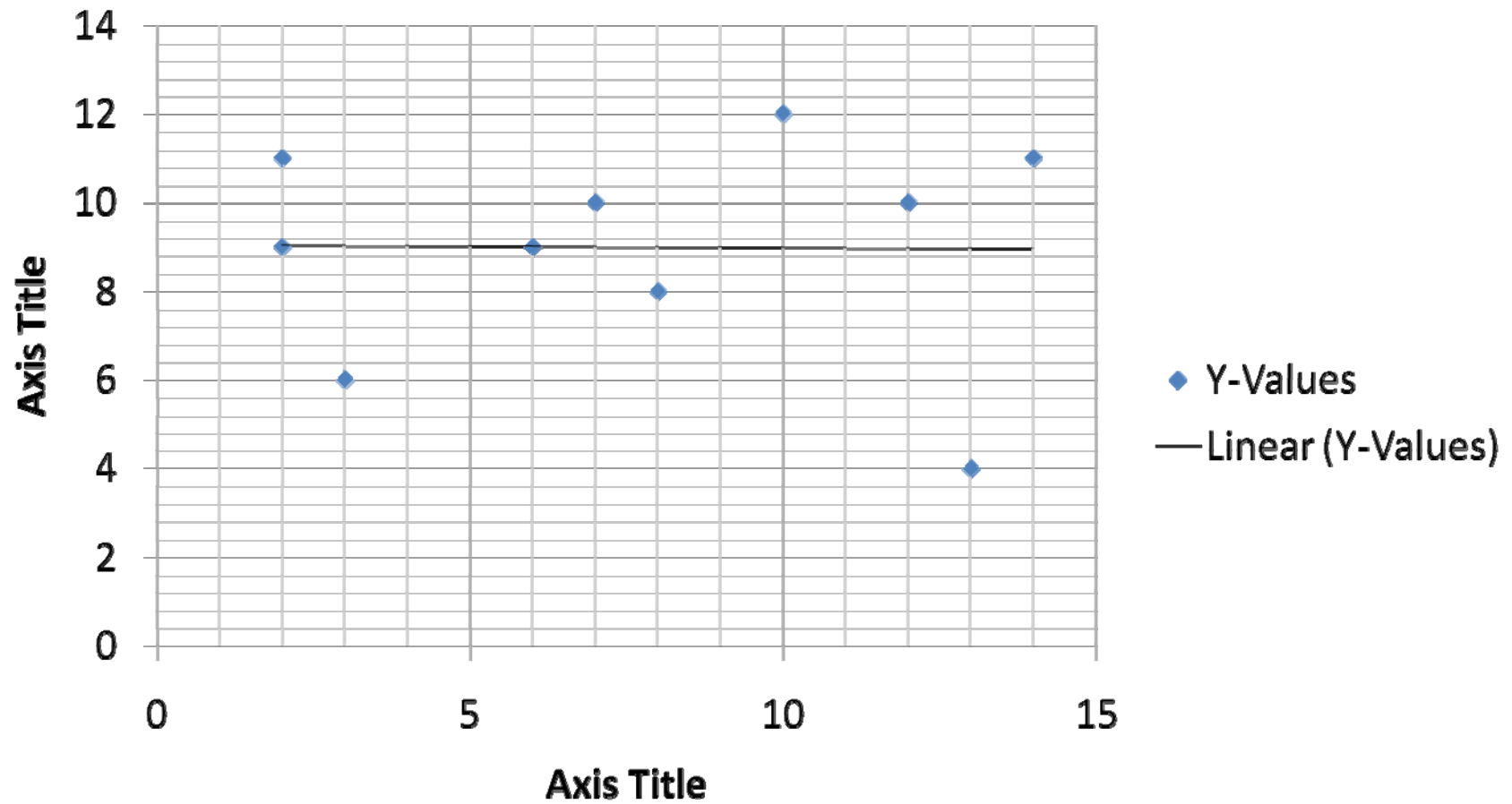
Scatterplot (Positive Rel.)



Scatterplot (Negative Rel.)



No Relationship



Analyzing Regression

- There are two questions to answer in a regression model:
 1. How well does the model fit the reality?
 - R^2 tells us the goodness of fit for the whole model
 - F tells us the statistical significance of R^2
 2. How well do the predictors (Xs) predict Y?
 - β tells us the impact of each X on Y
 - t and p -value tell us whether each X is statistically significant

Analyzing Regression

- R^2 tells us the percent of variation in Y explained by all X s.
- If R^2 is closer to 1, the model fits reality well.
- If R^2 is closer to 0, the model fits reality poorly.
- A high R^2 suggests a useful model...but equally important are p-values for each indep var
- Parsimony: a model should contain only the important X s for explaining variation in Y .

Regression with SPSS

The screenshot shows the SPSS Statistics Data Editor interface. The main window displays a dataset with one variable, 'Ebayexpenses', and 15 rows of data. The 'Analyze' menu is open, showing the 'Regression' option selected. The 'Regression' submenu is also open, showing various regression methods. The 'Data View' tab is active at the bottom.

	Ebayexpenses
1	200.00
2	125.00
3	65.00
4	24.00
5	145.00
6	89.00
7	67.00
8	79.00
9	92.00
10	46.00
11	21.00
12	48.00
13	77.00
14	104.00
15	188.00

The 'Analyze' menu is open, showing the following options:

- Reports
- Descriptive Statistics
- Tables
- Compare Means
- General Linear Model
- Generalized Linear Models
- Mixed Models
- Correlate
- Regression
 - Linear ...
 - Curve Estimation...
 - Partial Least Squares...
- Loglinear
- Classify
 - Binary Logistic...
 - Multinomial Logistic...
- Dimension Reduction
- Scale
- Nonparametric Tests
- Forecasting
- Survival
- Multiple Response
- Quality Control
- ROC Curve...

The 'Regression' submenu is open, showing the following options:

- Linear ...
- Curve Estimation...
- Partial Least Squares...
- Binary Logistic...
- Multinomial Logistic...
- Ordinal...
- Probit...
- Nonlinear ...
- Weight Estimation...
- 2-Stage Least Squares...
- Optimal Scaling (CATREG)...

Regression with SPSS

- SPSS first provides statistics for the whole model:

Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.905 ^a	.820	.790	24.63851

a. Predictors: (Constant), Hours Spent Online Per Week, Annual Household Income

- R^2 tells us that 82% of variation in our dep var is explained by the indep vars.
- Adjusted R^2 accounts for model size; the more parsimonious the model, the less difference between R^2 and adjusted R^2 .

Regression with SPSS

- Next, SPSS provides statistical signif for the whole model.

ANOVA^b

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	33104.661	2	16552.331	27.267	.000 ^a
	Residual	7284.672	12	607.056		
	Total	40389.333	14			

a. Predictors: (Constant), Hours Spent Online Per Week, Annual Household Income

b. Dependent Variable: Monthly Expenditure on Ebay

- The higher F is, the more likely it will be significant. If p is less than .05, we can conclude that the model is significant. Typically, that is not our most important concern.

Regression with SPSS

- Next, SPSS provides β and t-scores for the regression coefficients.

Coefficients^a

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	-43.250	20.430		-2.117	.056
	Annual Household Income	.000	.000	.077	.509	.620
	Hours Spent Online Per Week	8.943	1.579	.858	5.664	.000

a. Dependent Variable: Monthly Expenditure on Ebay

- The regression equation:

$$Y = -43.250 + .000X_1 + 8.943X_2$$

Regression with SPSS

- Controlling for income, every hour spent online per week increases by \$8.94 the amount of money spent on ebay per month.
- Controlling for hours online, every \$ in annual household income increases by \$0.00 the amount of money spent on ebay per month.
- (Actually, it's just a very small number-this often happens because of scale. To avoid this, I could express income in thousands).

Building Better Models

- Three Methods:

1-Review of Prior Literature: Reviewing past research on a subject to determine what previous analysis has found to be significant independent variables. Two caveats:

- Variables that are not statistically significant in one model may be statistically significant in others.
- The goal of new research is to expose new important variables and relationships...so researchers do not want to rely too much on prior research.

Building Better Models

2-Stepwise Regression Method: also known as the trial and error method.

- Start with a basic equation: $Y = a + bX$.
- Add a variable each time, and see if it:
 - Increases R^2 , and
 - The p-value is less than .05
- If it satisfies both criteria, we include it in the model.
- We continue until we have a parsimonious model with a high R^2 .
- In the example before, we would drop annual household income, because it was not st. significant.

Building Better Models

- *c-Data Mining Software*: software that will create a “decision tree” that shows how variables in a data set are related.
- This tells us which variables we should include in our regression equation.
- Examples: CHAID (Chi-squared Automatic Interaction Detector), SPSS Answer Tree.

Regression Table

- A regression table documents the coefficients and p-values in the analysis.
- Allows for the quick comparison of the usefulness of variables.
- Also useful for comparing equations using different combinations of the variables in your study.

Regression Table

Table 1.3. Regressions on Newspaper Circulation

<i>Category</i>	1	2	3	4
GNP per capita, averaged over 1991–95	1.12*** (13.6)	.80*** (8.24)	.76*** (7.58)	.64*** (6.89)
Illiteracy rate, averaged over 1991–95		-.03*** (-6.89)	-.03*** (-6.15)	-.02*** (-5.7)
Ethnic diversity	-.88** (-2.88)		-.50* (-1.70)	.19 (.75)
Africa				-.94*** (-5.05)
Constant	-5.17*** (-6.73)	-2.11** (-2.46)	-1.57* (-1.77)	-.70 (-.89)
R ²	.78	.80	.81	.84
Number of observations	96	79	76	76

* Significant at the 10 percent level.

** Significant at the 5 percent level.

*** Significant at the 1 percent level.

Source: GNP: compiled from World Bank databases; illiteracy rate: UNESCO (1999); ethnic diversity fractionalization index: Taylor and Hudson (1972); state ownership of newspapers: Djankov and others (2001).

Non-Linear Models

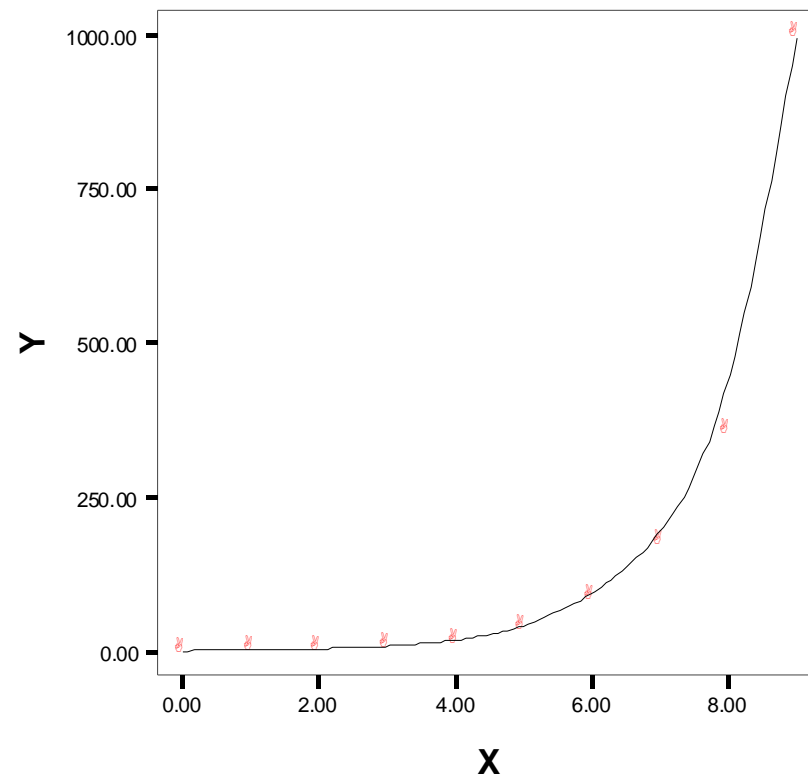
- Basic regression modeling assumes the association is linear.
- This means that an increase in X produces a consistent increase in Y.
- What if the increase in X produces a small increase in Y at low levels, and a large increase in Y at high levels?
- The equation then must be written as a nonlinear equation.

3 Types of Nonlinear Models

- Quadratic Model: Low levels of X produce small increases in Y; high levels of X produce high increases in Y.
 - Equation: $Y = a + bX^2$
- Logarithmic Model: Low levels of X produce large increases in Y; high levels of X produce small increases in Y.
 - Equation: $Y = a + b \log X$
- Parabolic Model: Relationship between X and Y is different at middle levels of X than at extremes.
 - Equation: $(y-k)^2 = 4a(x-h)$

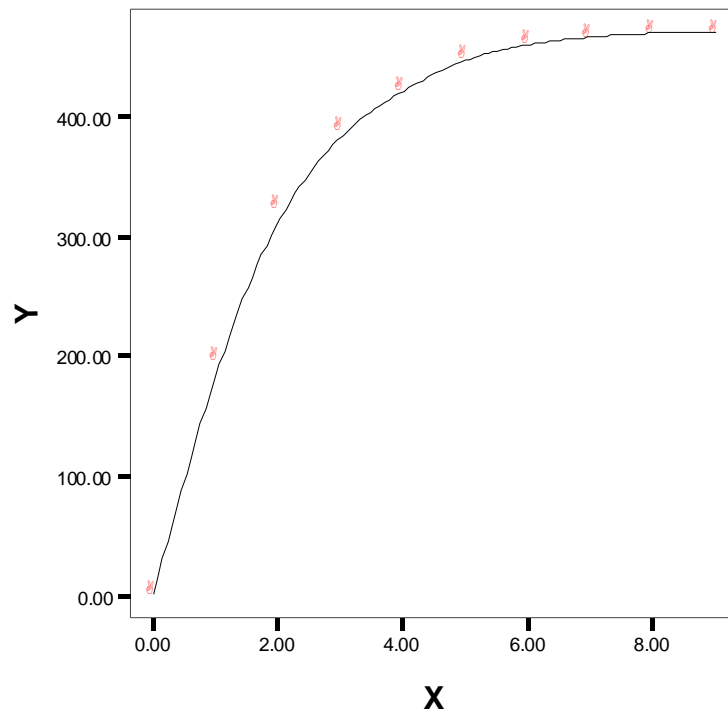
Quadratic Relationship

- $Y = \alpha + \beta X^2$



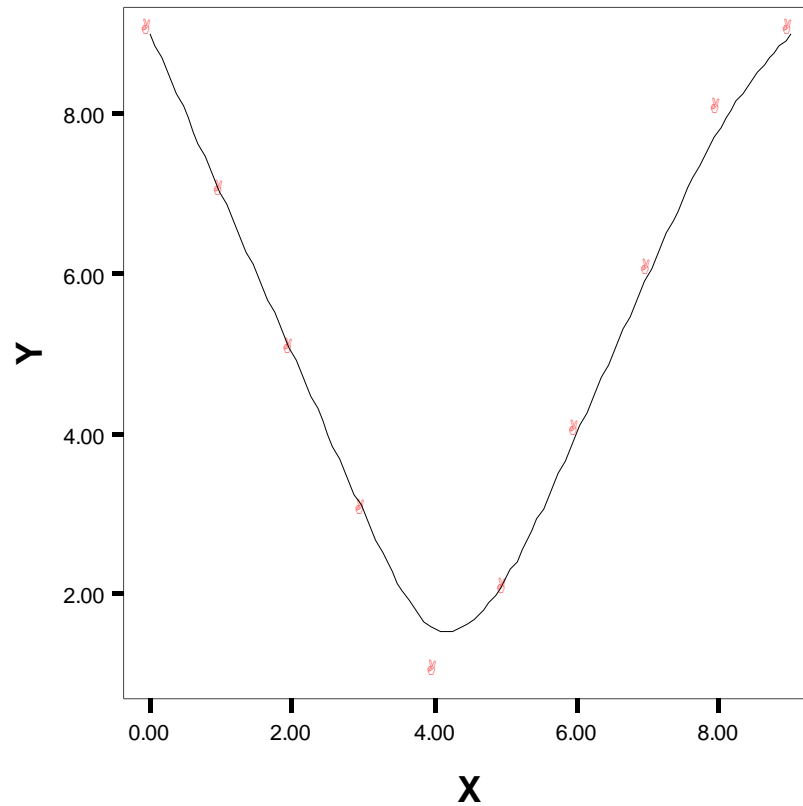
Logarithmic Relationship

- $Y = \alpha + \beta \log X$



Parabolic Relationship

- $(y-k)^2 = 4a(x-h)$



Indirect Associations

- Causal chain: a series of variables related to one another, that fall into a sequential time order.
 - Antecedent variable: the initial causal variable; it is directly related to the intervening, but not the dependent variable.
 - Intervening variable: directly related to the antecedent and dependent variables.
 - Dependent variable: the variable that is influenced by the antecedent (indirectly) and the intervening (directly).

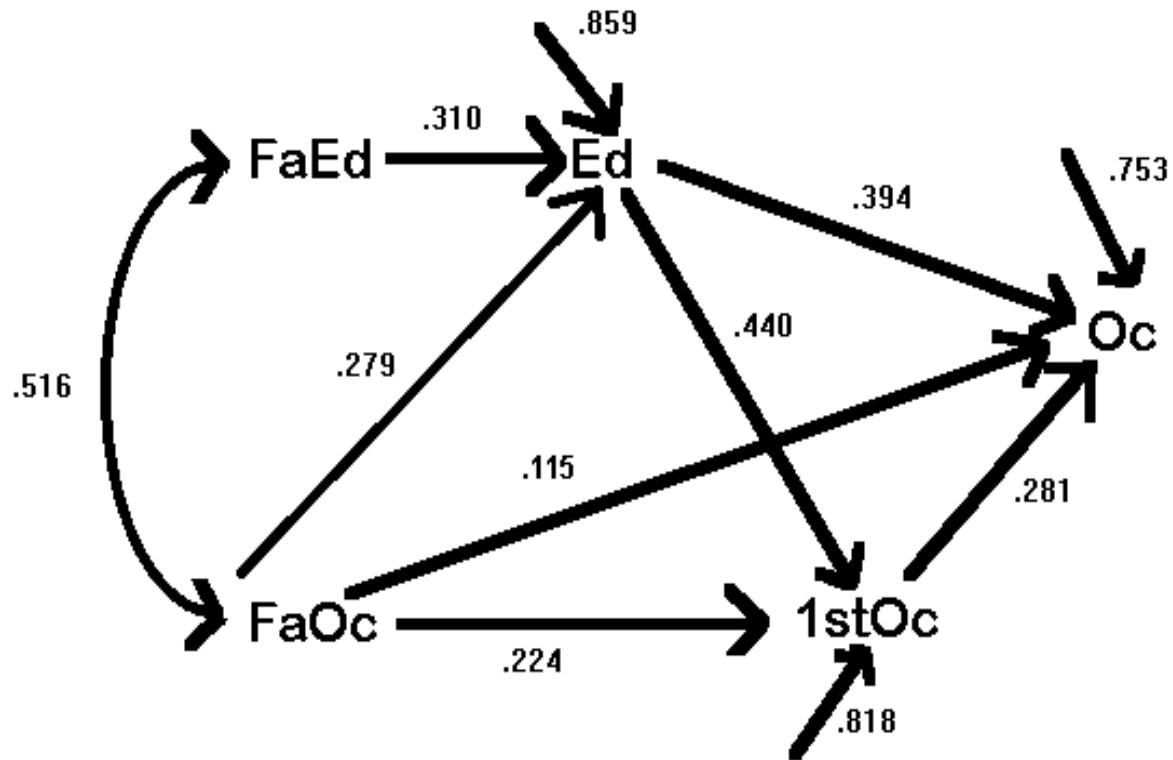
Path Analysis

- Path analysis is a method of modeling and analyzing a series of direct and indirect associations between variables.
- Variables are arranged in a sequence, and arrows are drawn depicting all theorized associations in the model.
- Regression coefficients are calculated for each arrow, and are provided to show degree of association in different parts of the model.

Path Analysis

- Coefficients are calculated using Simultaneous Equation Modeling Software:
 - Any variable that is “caused” in the model is called an endogenous variable. Any variable not “caused” is exogenous.
 - An equation is written for each endogenous var.
 - SEM software calculates regression coefficients for each path in the model, controlling for the effects of all other variables in the model. I.e., it calculates coefficients for all equations simultaneously.
 - LISREL, AMOS are two popular programs.

Path Diagram



Statistical Projections

(Forecasting)

Statistical Projections

- Point estimates that use longitudinal data to identify a likely future value of a variable.
- A projection is best thought of as the average value that would result, given past values.
- Projections are not predictions...we realize that the specific estimate is unlikely.
- Purpose: minimize risk by identifying a value that is closest to what the true value will be.

Strategic Plan

- Projections are an important part of creating realistic, aspirant goals in a strategic plan.
- A strategic plan entails:
 1. A set of **general goals** that represent key improvements in areas management has deemed important.
 2. A set of **specific targets**, for which achievement of such goals can be empirically measured.
 3. An **Outcomes Assessment Strategy**: a description of how achievement of targets will be measured. This may be more complicated than it at first seems.

Strategic Plan: Targets

- Specific targets are often built on projections.
- Targets are optimistic versions of the projections:
 - projections are what company x should achieve at its current level of performance
 - targets are what company x could achieve if it is maximally effective.
- The accuracy of projections is therefore extremely important.
- Budget and revenue targets may be built on statistical projections, and if they are incorrect, it could cause a budget shortfall in the coming year.

Types of Projections

1. Prior Term: The most recent figure is simply applied to the next term.

- $Y_t = Y_{t-1}$
- Does not account for possibility that Y_{t-1} is an outlier.
- Does not account for long term trend.
- Does not identify emerging trends.

Types of Projections

2. Simple Averages: The mean of some set of previous terms is applied to the next term.

$$\blacksquare Y_t = \frac{(Y_{t-1} - Y_{t-2} - Y_{t-3} \dots)}{n}$$

- Prior terms are included to minimize the impact.
- If Y_{t-1} is an outlier, then, its influence on the projection will be limited.
- Does not account for long term trend.
- Does not identify emerging trends.

Types of Projections

3. Weighted Average: The mean of some set of previous terms is taken, with some (typically, the most recent) weighted more heavily.

$$\blacksquare Y_t = \frac{3Y_{t-1} - 2Y_{t-2} - Y_{t-3}\dots}{n_{3 \times 2 \times 1 \dots}}$$

- Because most recent results are weighted, recent trends are accounted for.
- Influence of Y_{t-1} is greater, less limited by prior events.
- Does not identify emerging trends.

Types of Projections

4. Moving Averages: In each year, a new mean of a consistent series of prior terms is calculated.
 - Best used when we have a long series of averages in a spreadsheet.
 - Most common method is the 3-year Moving Average: each year, we calculate the average of the prior three years.
 - Allows the calculation of a long series of projections to minimize erratic fluctuations.

Types of Projections

5. Regression

- $Y_t = a + b_1 Y_{t-1} + b_2 X_2 + b_3 X_3 + \dots$
- We take the value for the prior term and treat it as a predictive variable.
- In addition, we identify certain independent variables to serve as leading indicators
 - i.e., variables that tend to predict changes in Y .

Types of Projections

- Y_t is measured for the current term; Y_{t-1} , X_2 , X_3 , and all other indep variables are measured for the prior term.
 - The only method that can use leading indicators to identify emerging trends.
 - Because regression is the average association between Y and its independent variables, if Y_{t-1} is an outlier, prior data will suppress its influence.
 - Long term trends are accounted for.

Final Notes on Projections

- The most common projection method, ARIMA, is beyond the scope of this course.
- ARIMA (Auto-Regressive Integrated Moving Average) builds on the moving average and regression methods discussed above.
- Forecasting is the science of developing, testing, and utilizing statistical projections.
- Forecasting has become a key part of how companies manage uncertainty, and therefore is a skill that is very much in demand.

Experimentation

Experiment

- A research method in which researchers expose their subjects to some cause, with the intention of measuring a resulting effect.
- There are numerous types of experiment, but in each, the researcher actively introduces the cause, with an interest in measuring the effect.
- Experiments are the best method of establishing cause.

Classical Experiment

- A “true” experiment, where the researcher has maximum possible control.
- Components of a classical experiment:
 1. Stimulus: a causal variable that the researcher uses to stimulate an effect among the subjects of the experiment. (Independent Variable).
 2. Dependent Variable: The focus of the study; the researcher applies the stimulus to see if there is change in the dependent variable.

Classical Experiment

3. Extraneous Variables: aka Control Variables-any variables, other than the independent, that might affect the dependent variable.
4. Subjects: The people or beings who participate in the experiment.
 - They are typically “blind”, i.e., they do not know whether they are part of the experimental or control group.

Classical Experiment

5. Pre-Test: A measure of levels of the dependent variable prior to the application of stimulus.
6. Post-Test: A measure of levels of the dependent variable following application of stimulus.
7. Experimental Group: The stimulus is applied only to one group of subjects.
8. Control Group: A group of subjects, otherwise identical to the experiment group, who are not exposed to the stimulus.

Classical Experiment

- A classical experiment uses the following procedure:

Stage	Experimental	Control
Pretest	Yes	Yes
Stimulus	Yes	No
Posttest	Yes	Yes

- If there is an effect, the difference in scores from the pretest to the posttest should be different for the experimental group, as compared with the control group.

Classical Experiment

- *A Controlled Setting*: a setting where researchers can restrict environmental factors, preventing any such factors from influencing the behavior of the control or experimental groups.
- The participants must have identical representation of extraneous variables.
 - Any characteristics of the participants that the researcher identifies as possible extraneous variables must be represented equally in the experimental and control groups.

Classical Experiments

- The Placebo Effect: a phenomenon where the results of a any stimulus-not just the specific treatment given to the experimental group-produces a change in the dep. Variable.
- Extremely important to look out for when comparing among multiple actions, such as the effectiveness of two or more advertising campaigns.
- Control groups should be designed to isolate the placebo effect.

Interpreting the Results

$$\Delta = (E_2 - E_1) - (C_2 - C_1) \quad \dots\text{where}$$

E_2 refers to the summary statistic (for example, a mean) value for the **experimental** group on the **posttest**.

E_1 refers to the summary statistic value for the **experimental** group on the **pretest**.

C_2 refers to the summary statistic value for the **control** group on the **posttest**.

C_1 refers to the summary statistic value for the **control** group on the **pretest**.

Δ refers to the total experimental effect.

Interpreting the Results

- The greater the value of Δ , the greater the experimental effect.
- Statistical significance in experimental design is conducted by doing a hypothesis test for a difference in means or proportions.
- More complicated experiments with multiple experimental groups use ANOVA (Analysis of Variance), which is a difference of means test for more than two groups.

Types of Experiment

- Classical Experiment: Any experiment with:
 - A controlled setting
 - A Pre-Test and a Post-Test
 - An experimental group and a control group
- Quasi-Experiment: Any experiment that falls short of one of those three criteria.
 - Lacks a controlled setting; has only one test; or does not have a control group.

Double-Blind Tests

- Neither the subjects nor the researchers know which subjects belong to the control or experimental group.
- Only after all the data have been recorded (and in some cases, analyzed) do the researchers learn which individuals are which.
- Purpose: to lessen impact of researcher prejudices and unintentional physical cues on the results.
- The key that identifies which group each subject belongs to is kept by a third party and not given to the researchers until the study is over.

Advantages of Experiments

- The best way to determine causality.
- Researcher can isolate specific variables for attention.
- Researcher can Isolate components of complex relationships.
- They can be replicated easily.
- Researcher has flexibility with application of stimulus.

Disadvantages of Experiments

- External validity: how generalizable are experiments to the real world?
- Cost: must pay for laboratory space.
- Attrition: if participants need to leave midway, it can bias the posttest.
- Maturation/Learning: Experience of Pretest can bias Posttest

Field Experiments

- Any experiment conducted in a natural setting, as opposed to a laboratory experiment, conducted in an artificial setting.
- This can include the workplace, where people shop, a public square, or a park.
- Experiment must take place where the participants would be likely to go were they not participating in the experiment.
- The strength of laboratory experiments is that they come closest to eliminating the effects of all of the extraneous variables on the dependent variable.

Field Experiments

- The weakness of laboratory experiments is that they may not be generalizable to the real world (i.e., they lack external validity).
- Field experiments often do not conform to the rigors of the classical experimental design. They are much more likely to lack pretests or control groups because:
 - They take place in natural settings, and thus it is difficult to control all potential extraneous factors.
 - The effort to establish control and experimental groups, and take a pretest, may interfere with the “naturalness” of the participants’ behavior, thus removing the justification for doing a field experiment in the first place.
- Test marketing is an example of a field experiment.

Qualitative Research

What is Qualitative Research?

- The examination, analysis and interpretation of observations for the purpose of
 - discovering underlying meanings and patterns of relationships,
 - including classifications of types of phenomena and entities
 - in a manner that does not involve mathematical models.

Benefits of Qualitative Research

- Allows more in-depth examination of research subject.
- It can investigate more thoroughly the why and how of decision making (not just managerial, but consumer).
- More useful for establishing social processes, including the relationships between different actors, offices, and symbols.
- More amenable to follow-up (probing) questions that can extract greater detail from respondents.

Weaknesses of Qualitative

- Time consuming
- Costly
- High subjective; high level of interpretation.
- Small sample sizes
- Greater ethical risk
- Difficult to summarize

5 Types of Qualitative Research

- Interviews
- Focus Groups
- Observation
- Projective Techniques
- Physiological Research

Interviews

A researcher either asks questions for a response, or suggesting ideas to a respondent for commentary. Two general formats:

1. Open-Ended Items on a Questionnaire:

- Instead of response options, respondent provides responses as he/she sees fit to questionnaire item.
- Researcher reviews responses for key terms, themes, phrasing, and opinions.
- Especially useful when we cannot anticipate an exhaustive range of possible responses.

Interviews

2. *Depth-Interviews*: The researcher actively responds to the subject's answers, probing for more information and providing counterpoints.

The researcher starts with a general list of questions, but does not stick to the script. For each response, the researcher can respond spontaneously with:

- probing questions
- significant events or facts that relate to the response
- requests for comments on related topics
- suggestions of possible solutions to problems.

Interviews

- Can fall along a spectrum from a fairly straightforward question and answer session to something that resembles a free-form conversation between interviewer and subject.
- *Laddering*: establishing a linkage from product attributes to consumer values through probing questions.
 - A useful technique for probing in marketing research.

Interviews: Adv & Disadv

Advantages of Interviews:

1. Compared with closed-ended items in a questionnaire, they allow the research an opportunity to get greater detail from the respondent.
2. The researcher may get information that he or she was not planning to get (through probing and dialog). This takes away the burden (for the researcher) of knowing everything relevant to the subject in advance of the questionnaire.

Disadvantages of Interviews:

1. It is highly subjective: we rely on the interpretations of researchers.
2. It may be difficult and possibly quite time consuming to do the analysis, which involves a detailed review of a potentially large amount of text.
3. The risk of bias is very high any time the interviewer has the chance to interact with the respondent.

Focus Groups

A small group involved in an unstructured, spontaneous discussion for the purpose of gaining information relevant to the research process.

A facilitator leads the discussion:

- Asking questions.
- Providing background information.
- Answering respondents' questions as needed.
- Maintaining decorum and assuring that the discussion is respectful.
- The Facilitator's task is to ensure the discussion is "focused" on some general area of interest.

Developed by Robert K. Merton 1941.

Information from focus groups can be used to generate ideas, to learn the respondents' vocabulary when relating to a certain types of product, or to gain some insights into basic needs and attitudes.

Focus Group Medium

1. Face-to-face: Participants get together in a room. Facilitator suggests questions.
 - Record discussion with note takers or recording device.
 - Unethical to tape participants without their knowing.
2. Online: Participants join a newsgroup, message board, or some other environment.
 - Benefits: all opinions can be recorded, people can think about what they want to say, and people can participate from anywhere at times they find convenient.
 - Disadv: requires internet familiarity, we may want spontaneous discourse, and infrequent participants may forget to keep up with the group.

Two types of Focus Groups

1. Traditional focus group:

- Facilitator provides a list of basic questions for participants to discuss with one another.
- Facilitator abstains from involving him/herself in the discussion, except as needed.
- Facilitator remains neutral to avoid bias.

2. Nontraditional focus group:

- Facilitator is actively involved in the discussion, expressing an opinion and defending arguments.
- Facilitator hopes to inspire heated debate.
- In exploratory research, bias less of a concern.

Focus Groups: Adv & Disadv

Advantages of Focus Groups:

1. Generate fresh ideas that might be suggested or alluded to by participants
2. Allows clients to observe participants
3. May be directed at understanding a wide variety of issues, such as reactions to a new food product or logo
4. They allow fairly easy access to special respondent groups

Disadvantages of Focus Groups:

1. Rarely use representative samples-so generalizability is low
2. Facilitator may bias responses, even if he or she is
3. The cost per participant is high, though the total spent on focus group research is actually a fraction of what may be spent on quantitative research.

Observation

- Researchers watch how people behave in their ordinary environment and record what they see.
- Appropriate Conditions for Use of Observation:
 - The event must begin and end within a reasonably short time span.
 - The public behavior occurs in a setting in which researcher can readily observe actions
 - It is best used when consumers cannot recall their behaviors (“faulty recall”), such as knowing how many different web pages they accessed.

Forms of Observation

- Undisguised vs. Disguised: Are you going to tell your subjects that you are watching them (**undisguised**) or are you going to try to keep them from knowing (**disguised**)?
- Structured vs. Unstructured: Are you going to identify in advance specific behaviors you intend to record, and how you will code them (**structured**), or are you going to let things happen and decide as you go which elements are to be recorded and how (**unstructured**)?
- Direct vs. Indirect: Are you going to observe behavior as it occurs (**direct**), or are you going to wait until afterward, to review the results of their behavior (**indirect**)?
- Human vs. Mechanical: Are you going to record the information yourself (**human**), or will you do so with a **mechanical** device: a camera, audio device, or people meter (Nielsen uses this for their ratings methodology)?

Observation: Adv & Disadv

1. Advantages of Observational Data

- Allows researchers to see what subjects actually do rather than relying on their own self-reporting-reduces bias.
- Exposes minute details of human behavior people may not otherwise volunteer to report.
- Gets at behavior that people are embarrassed to report.

2. Disadvantages of Observational Data

- Very time consuming.
- Low sample size, harming generalizability.
- Highly subjective; researcher interpretation high.
- Researcher cannot determine consumer's motives, attitudes, and intentions-since they cannot ask questions.
- Subjects may also feel that their privacy is being violated, which is an important ethical problem.

Projective Techniques

- Participants are placed in (projected into) simulated activities in the hopes that they will divulge things about themselves that they might not reveal under direct questioning.
- Word-Association Test: involves reading words to a respondent, who then answers with the first word that comes to his or her mind
- Sentence Completion Test-respondents are given incomplete sentences and asked to complete them in their own words
- Picture Test-respondents are instructed to describe their reactions by writing a short story about the picture

Projective Techniques

- Cartoon or Balloon Test- a line drawing with an empty balloon above the head of one of the actors; respondents are instructed to write what the actor is saying or thinking.
- Role-Playing Activity: participants are asked to pretend that they are the third person such as a friend or neighbor, and to describe how they would act in a certain situation or to a specific statement.
 - By reviewing their comments, the researcher can spot latent reactions, positive or negative, conjured up by the situation.

Projective Tech. Adv & Disadv

1. Advantages of Projective Techniques

- Allows researcher to measure unconscious factors affecting human behavior, including things people may not know about themselves.
- It may get at minute details of human behavior that people may not report.
- It may also get at behavior that people are embarrassed to report.

2. Disadvantages of Projective Techniques

- The interpretation of observed behavior is subjective.
- There is some debate among psychologists about the usefulness of efforts to tap into the unconscious mind.

Physiological Measurements

- Monitoring a respondent's involuntary responses to marketing stimuli by the use of electrodes and other equipment.
- The pupilometer is a device that attaches to a person's head and determines interest and attention by measuring the amount of dilation in the pupil or the eye.
- The galvanometer is a device that determines excitement levels by measuring the electrical activity in the respondents' skin.
- Physiological measurement is often uncomfortable to the respondents, and thus is fairly uncommon.

Presentation

Importance

- Your presentation is how you introduce your work to others.
- It is typically the first--and often the last--thing people see about your research.
- Even the best marketing research design won't have a strong impact if it is presented poorly.
- Regardless of how much work you did, your presentation is the part that people will see, and what they will remember.

Warning

- Your research may be used to analyze programs, marketing campaigns, or product ideas in which **your co-workers have invested time and energy.**
- If your results are perceived as threatening, **they will attempt to diminish your work.**
- The easiest way to do this is by finding errors; and **any error can be used to dismiss your work.**
- Anything you produce, even for your own company, becomes **part of your reputation.**
- So effective presentation is extremely important.

Report vs. Presentation

- The presentation is typically a quick overview given in limited time; the report provides the detailed analysis that can be referenced later.
- Presentation: emphasize visuals, conclusions, and memorable phrases.
- Report: emphasize detail, support for all assertions, and sources.

Report Organization

- Front Matter: Title, Abstract, Contents, Authorization
 - Abstract (Executive Summary): brief review of subject, method, and findings of report.
 - Contents: list of all chapters, sections, and tables by page.
 - Letter of Authorization: a market research company's permission to perform research on a private company.
 - Letter of Transmittal: a release form identifying who is allowed to see the report.
- Body: Discussion of methodology, results, conclusions, future directions.

Report Organization

- End Matter: typically provides backing information and details for the body.
 - Appendices: includes all tables, graphs, and charts not included in the body of the report.
 - Typically, detailed collections of data that lack visual appeal are located here.
 - Bibliography: references cited in text are included here. Audience may be interested in source of supporting data.

Visuals

- Refers to any table, graph, diagram, or illustration that is used to present the results.
- Purpose: to quickly and effectively convey the findings of the research to your audience.
 - Most people are not comfortable with statistics.
 - Audience is likely not as immersed in the details of the subject matter as are you, the researcher.
 - Presentation should seek to “tease” audience about information that they could follow in more detail in your report.

Tables

- Diagrams that simply list values, frequencies, means, or other statistical information.
- Univariate tables are *frequency tables*; multivariate tables are *cross tabulations tables*.
- Tables are not visually appealing; the mass of numbers can cause eyes to glaze over.
- Best used in appendices or in more technical or academic reports.

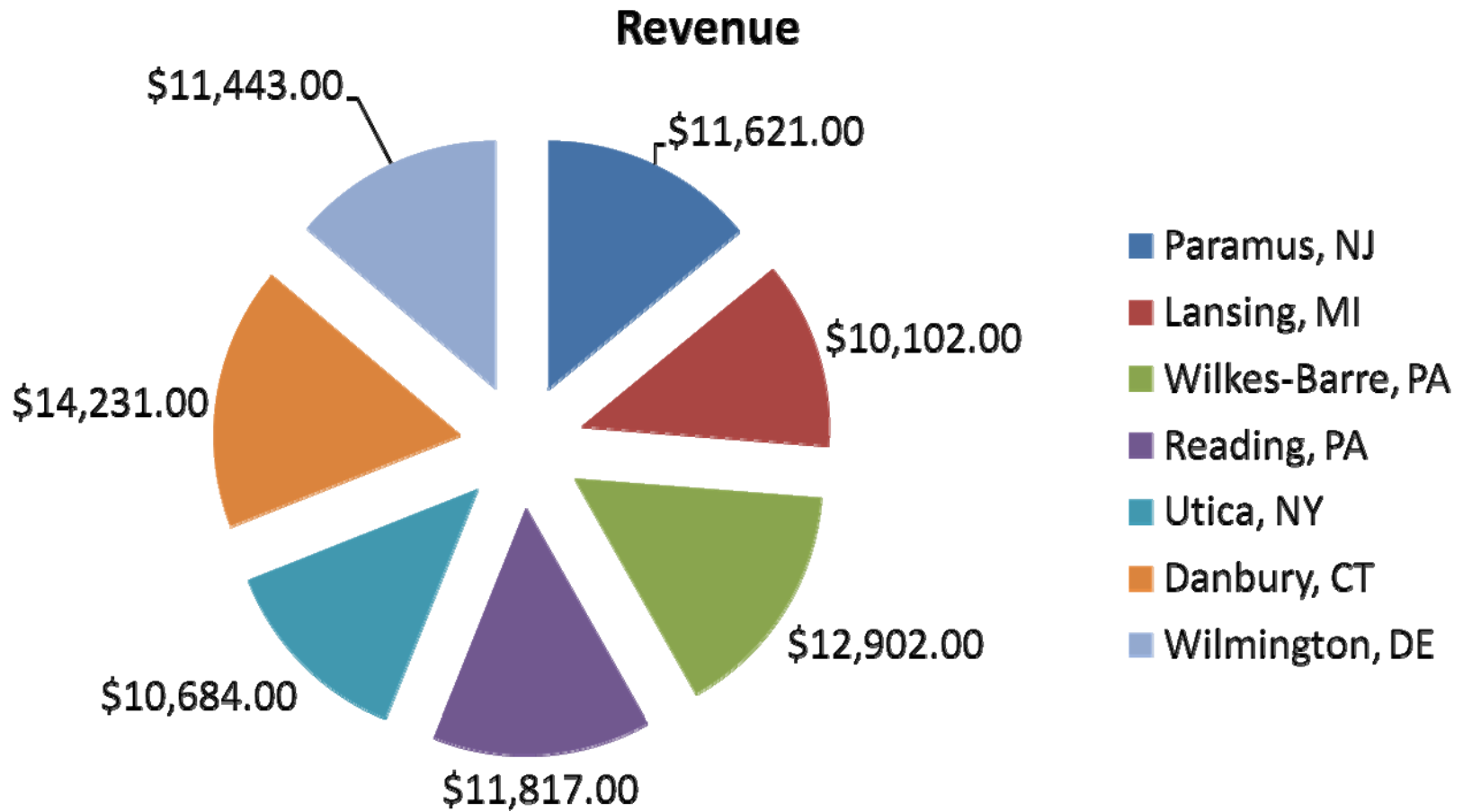
Tables

Store Location	March 2009 Revenue (in \$K)	March 2009 Expenses (in \$K)	March 2009 Difference (in \$K)
Paramus, NJ	\$13,050	\$11,621	\$1,429
Lansing, MI	\$9,795	\$10,102	(\$307)
Wilkes-Barre, PA	\$14,624	\$12,902	\$1,722
Reading, PA	\$12,932	\$11,817	\$1,115
Utica, NY	\$16,995	\$10,684	\$6,311
Danbury, CT	\$17,212	\$14,231	\$2,981
Wilmington, DE	\$11,813	\$11,443	\$370

Pie Charts

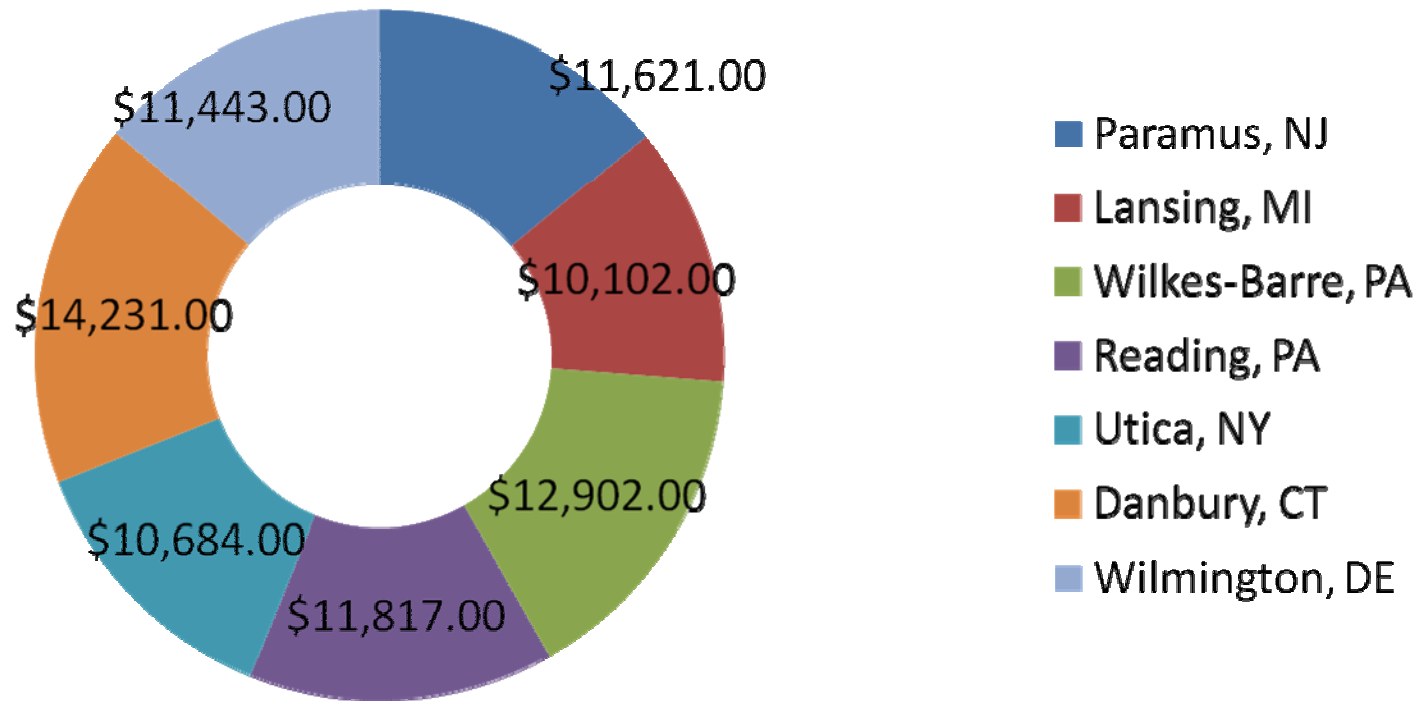
- Data is presented in a circular format, where the size of each “slice” or “wedge” indicates the proportion of all cases for a given value.
- Most useful for showing that a particular value represent a high or low **percentage** of the total.
- Not to be used when there are missing values (or, missing values should be excluded from the chart, so that it only uses valid percentages).

Pie Chart



Doughnut Chart

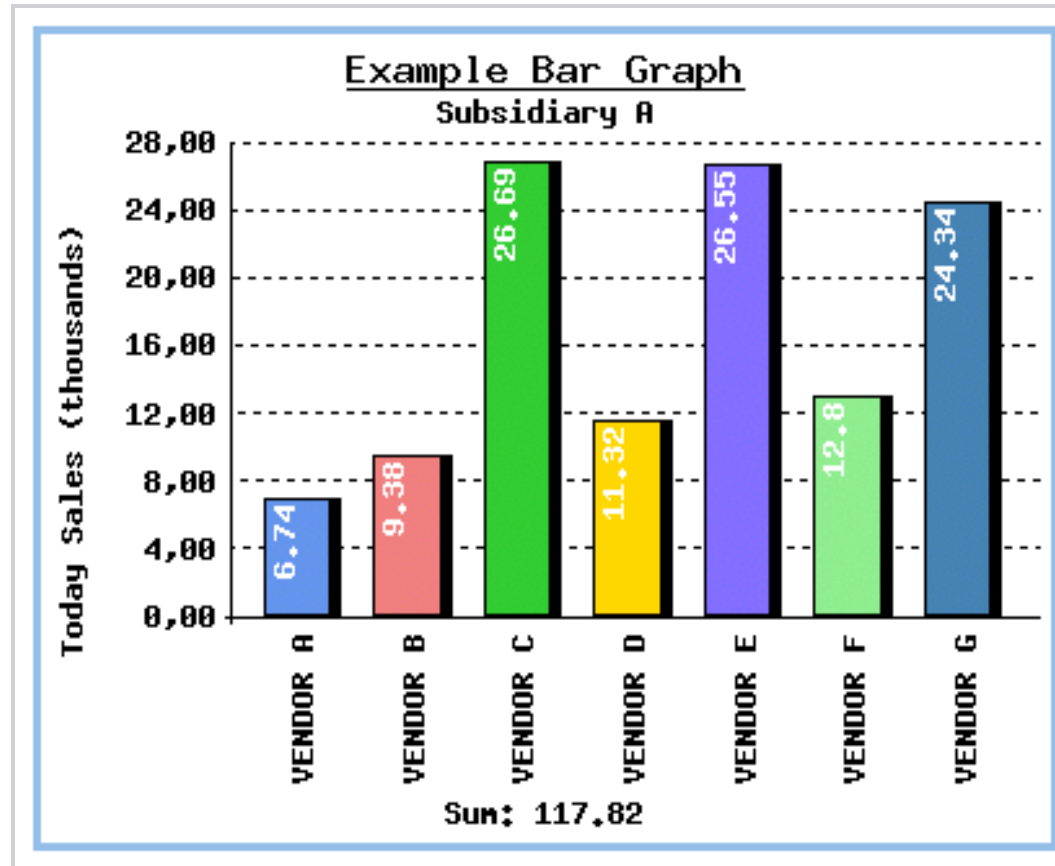
Expenses



Bar Chart

- Express counts or percentages in a series of bars, where the higher the bar, the greater the frequency or percentage.
- Unlike pie charts, it is not necessary to include 100% of the cases in the chart.
- E.g., I could use a bar chart to compare the number of freshmen vs. the number of sophomores, even if the counts of juniors and seniors are missing.

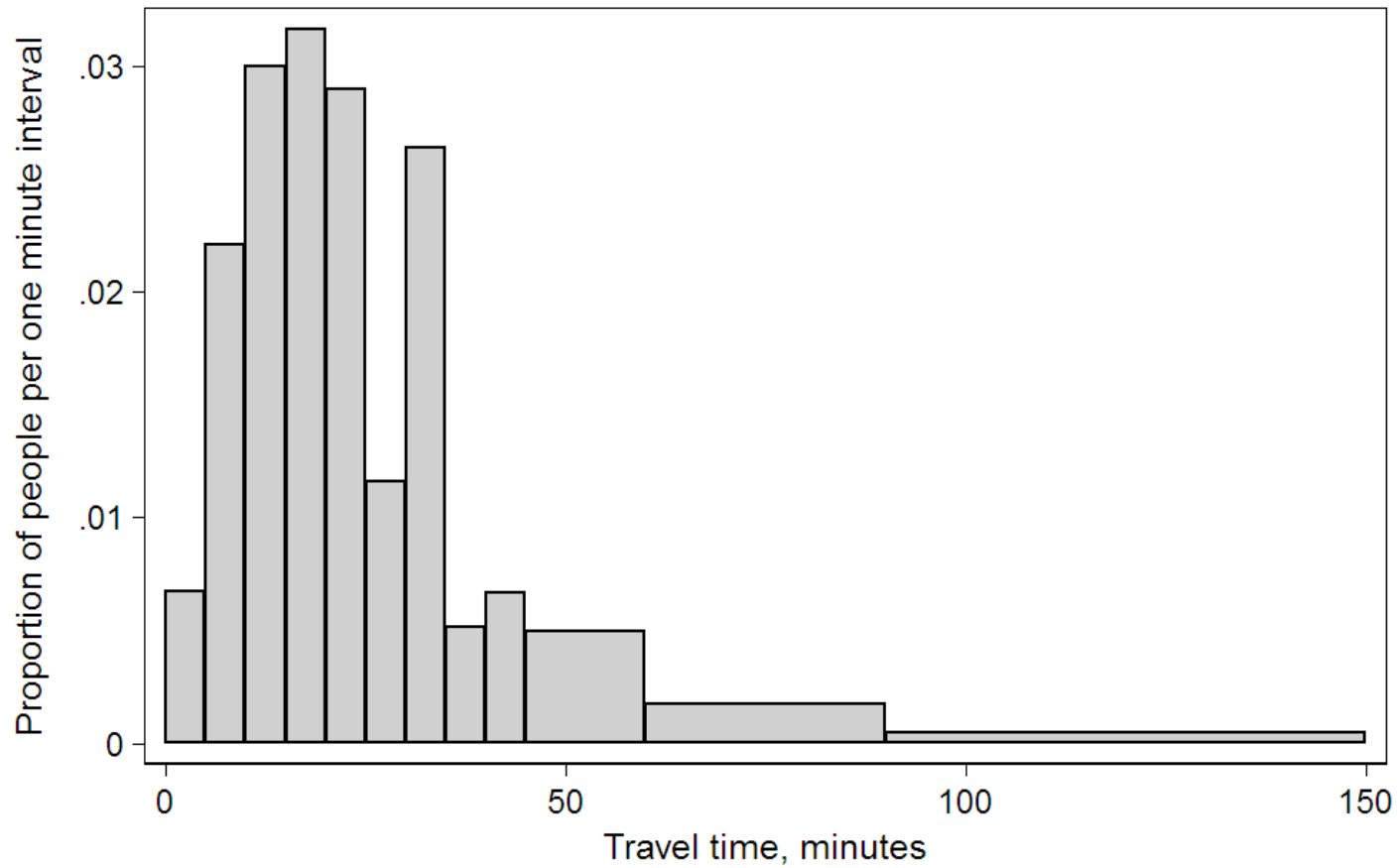
Bar Chart



Histogram

- A type of bar chart where the width of the bars, and not just the height, are used to present results.
- If you were comparing the count of people who were between 5 ft. 3 in. and 5 ft. 5 in. tall to people who were between 5 ft. 6 in. and 5 ft. 11 in. tall:
 - The range with the higher count would be taller, as with a bar chart.
 - But the range 5'3"-5'5" group would be two spaces wide; and the 5'6"-5'11" would be five.
- A histogram should only be used when we are measuring groupings of continuous level variables.

Histogram



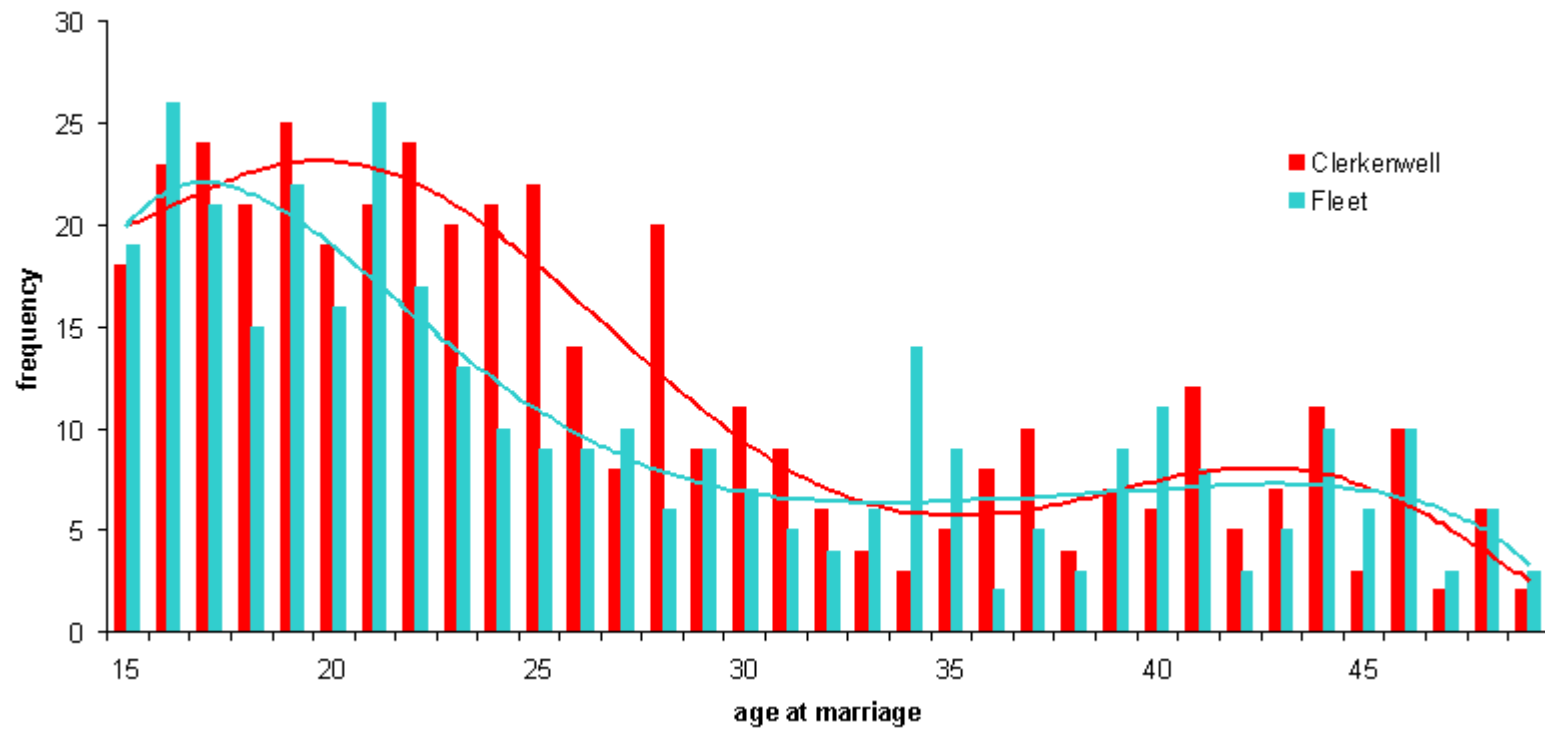
Line Graphs

- A line connects points, usually working from left to right in relation to the X-axis.
- Typically used when dealing with continuous variables; continuous values are always, at least, implied.
- Most useful for describing differences in some quality over time, or under changing conditions.

Line Graphs: 3 Types

- *Frequency Distribution Graph*: displays the counts for various levels of a variable. A univariate chart.
 - The x-axis represents values of a variable, ascending from low scores at the left to high scores at the right. The y-axis represents frequency, with low frequencies at the bottom, and high frequencies at the top.
 - At each value of x , we can enter a high dot for high frequencies and a low dot for low frequencies. Then, we draw a line from the left to the right, which shows us which values of x had the highest frequencies.
 - A normal curve is a frequency distribution graph.

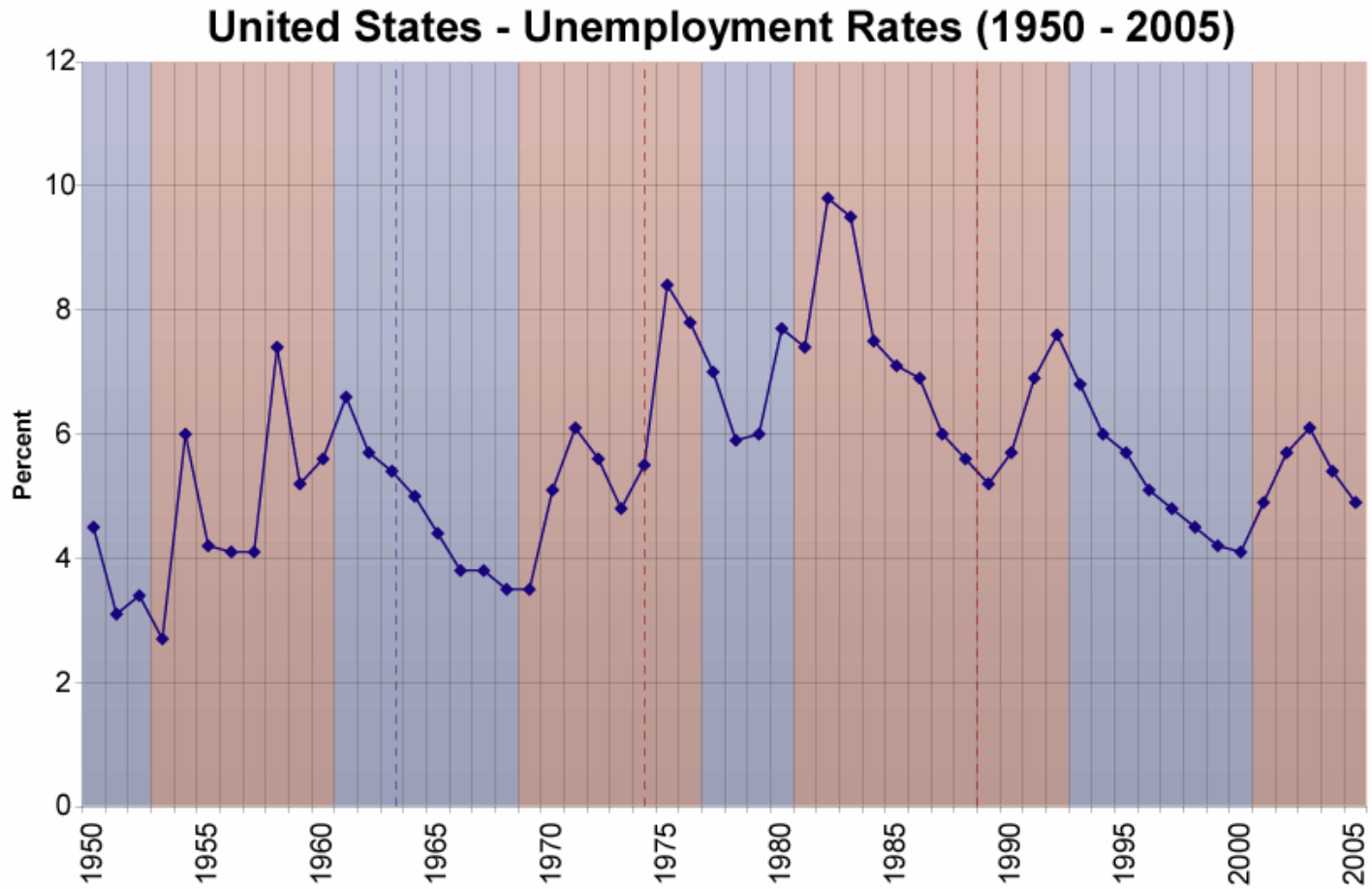
Frequency Distribution Graph



Line Graphs: 3 Types

- *Longitudinal Graph*: displays variables at various points in time.
 - The x-axis represents time, from earliest date (left) to most recent (right).
 - The y-axis represents values of a variable, from low (bottom) to high (top).
 - This is an effective tool for showing how a particular variable has changed over time.
 - This is also a univariate graph.

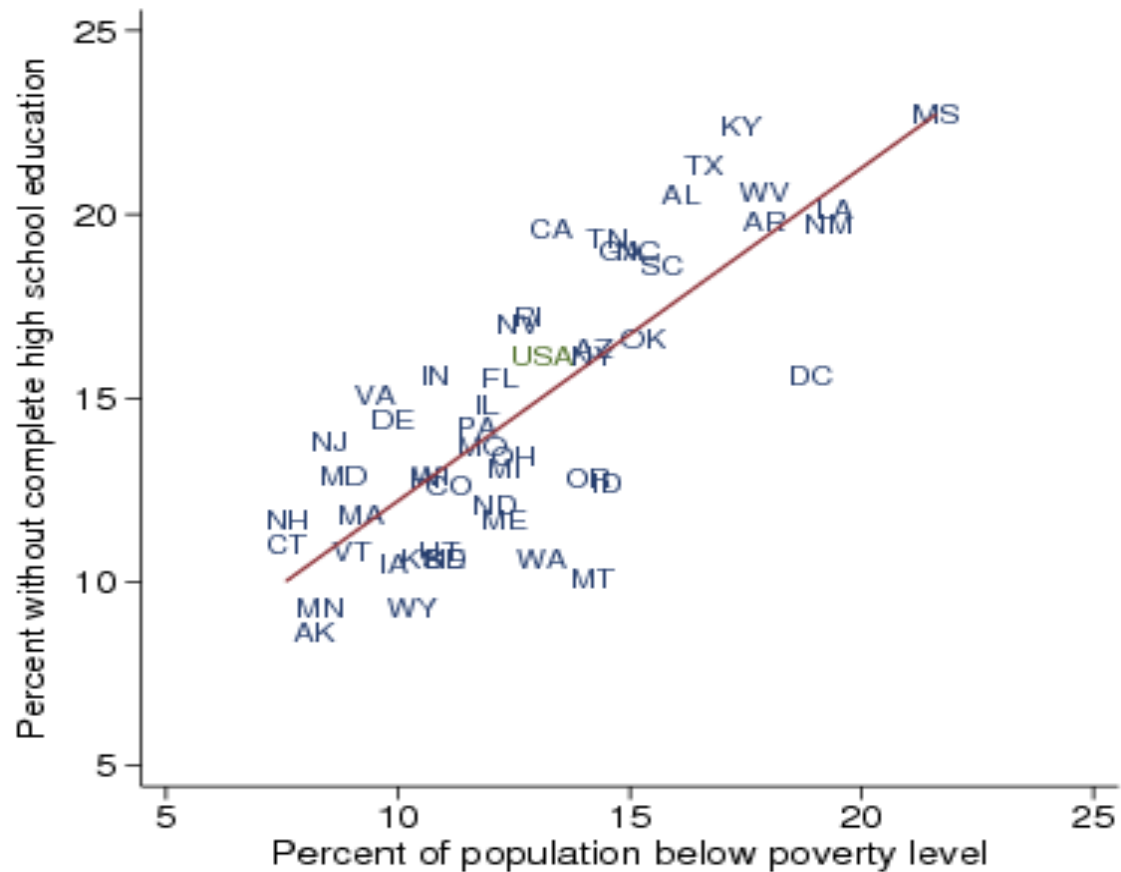
Longitudinal Graph



Line Graphs: 3 Types

- *Scatterplot with Summary Line*: A bivariate graph, where a single line summarizes the overall distribution of dots in a graph.
 - The x-axis represents the independent variable; the y-axis represents the dependent variable.
 - The scatterplot itself is a series of dots representing combinations of values of x and y for each case. The summary line is determined by regression analysis, and is sometimes known as a regression line.
 - Very commonly used to show that two continuous variables are related.

Scatterplot/Regression Line

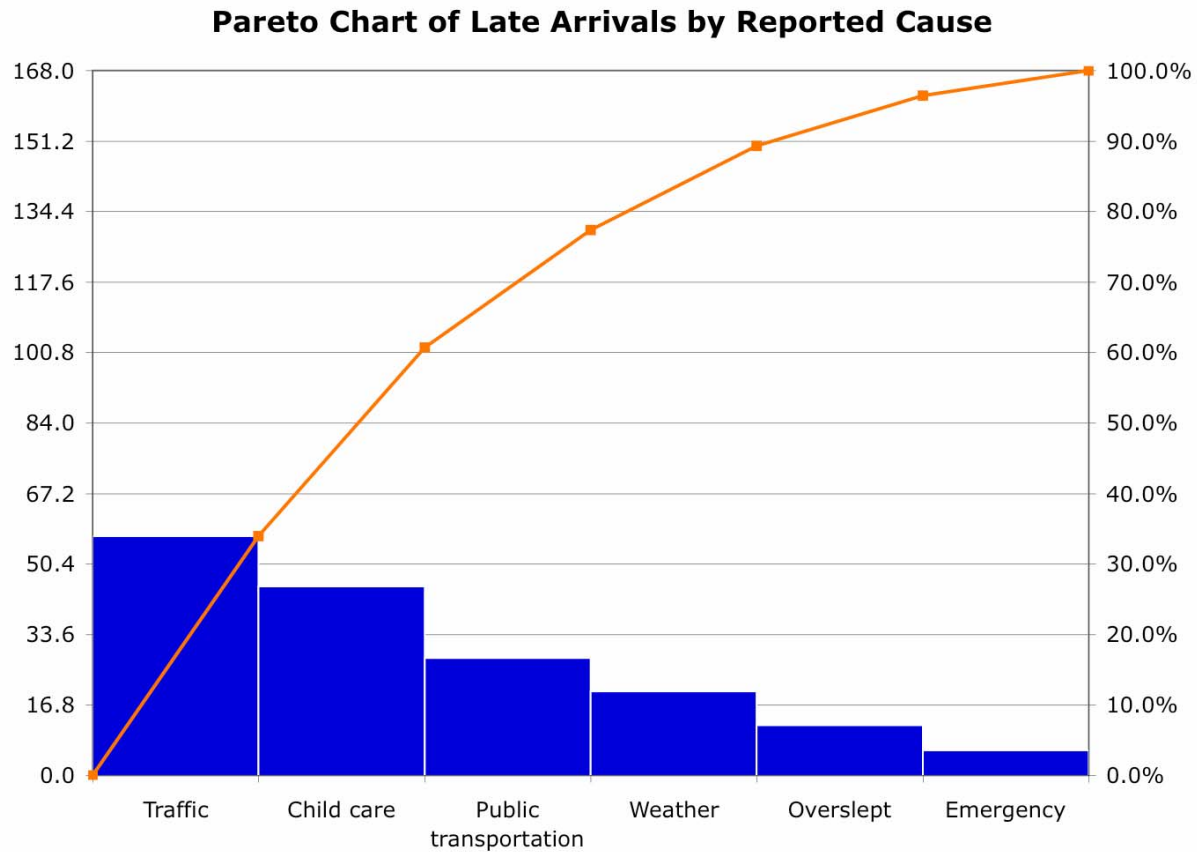


Friedrich Huebler, September 2005, huebler.blogspot.com

Pareto Chart

- A combined bar chart and line graph.
- The bar chart plots the values in order of descending frequency.
- The graph is accompanied by a line graph which shows the cumulative totals of each category, left to right.
- The line graph is typically cumulative.
- The purpose is to highlight the most important among a (typically large) set of factors.

Pareto Chart



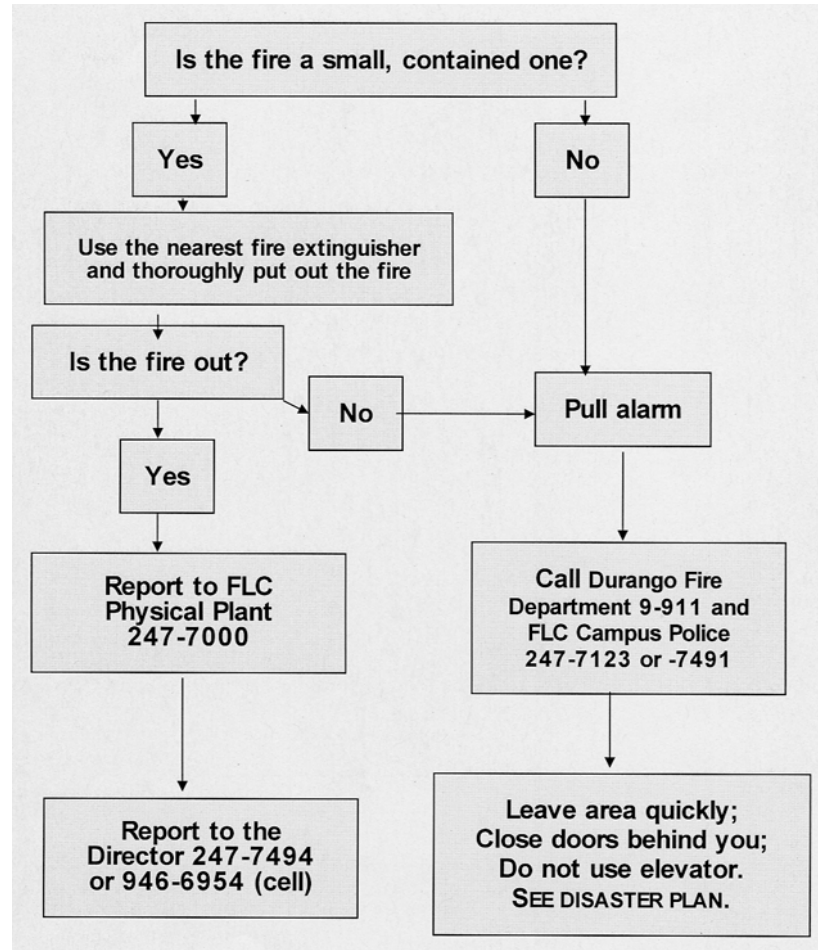
Process Charts

- Charts that depict a succession of elements connected (typically) by arrows.
- Generally, the earlier elements are at left, and the more recent are at right.
- These charts describe a process in terms of a sequence of events.
- Managers or researchers can then assess where in the chart an element might need to be added, removed, or changed.
- 2 Types: Flow Chart and Path Diagram.

Flow Chart

- Documents a succession of events, tasks, or offices in the production of some good or service.
- A flow chart might show how a check is processed, from the store where a transaction is made, to Accounts Payable, to a bank.
- Useful for assessing the role of each office in a specific business process.

Flow Chart



Path Diagram

- Documents a succession of causal associations in a model explaining a series of processes.
- Instead of offices, individuals, or actions, “stops” in a path diagram refer to variables in a model.
- Arrows identify correlations between variables, including direct and indirect effects.
- Path diagrams are used to model a complex series of causes & effects in an ongoing process.
- SPSS doesn't do this type of analysis; LISREL or AMOS are programs that do.

Path Diagram

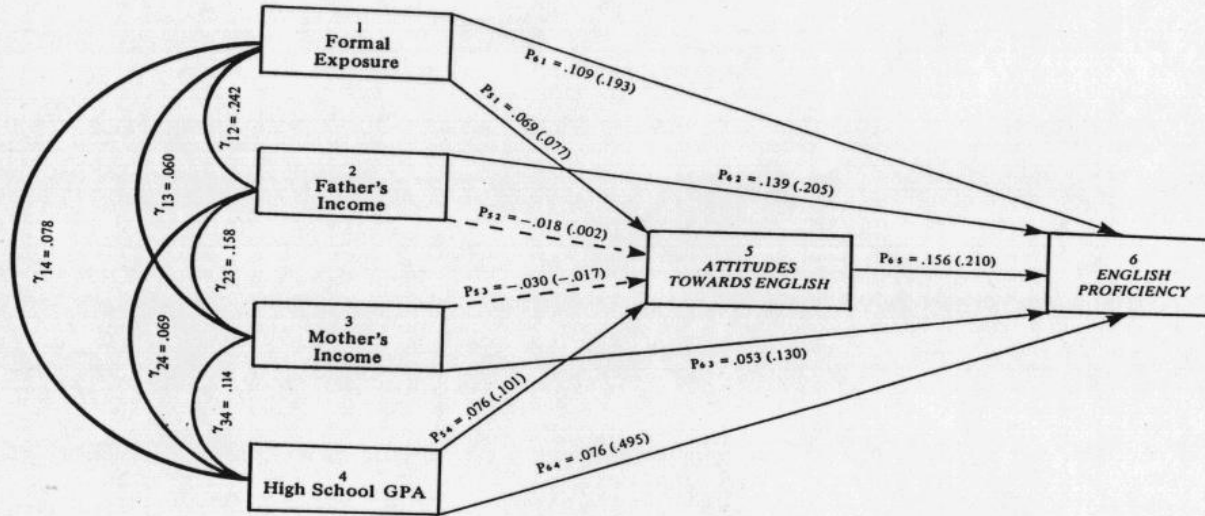
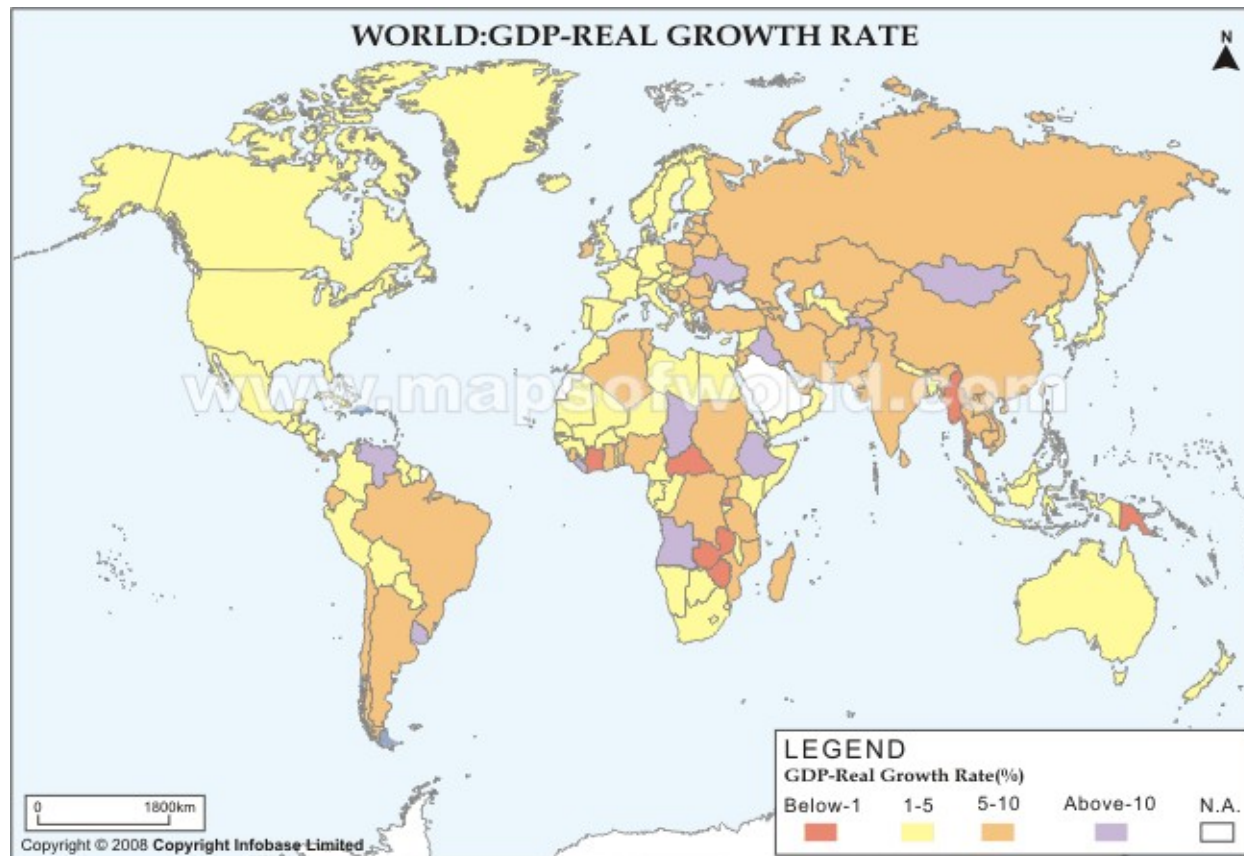


Figure 1
The path analysis model

GIS Mapping

- GIS (Geographic Information System) is software that allows researchers to display results of analysis in the form of a map.
- Useful for making comparisons between geographic areas on some key variable.
- SPSS doesn't do this; ArcView is the most popular software (provided by ESRI, one of the largest marketing research companies).

Map Used to Make a Point



Map Created by ArcView GIS

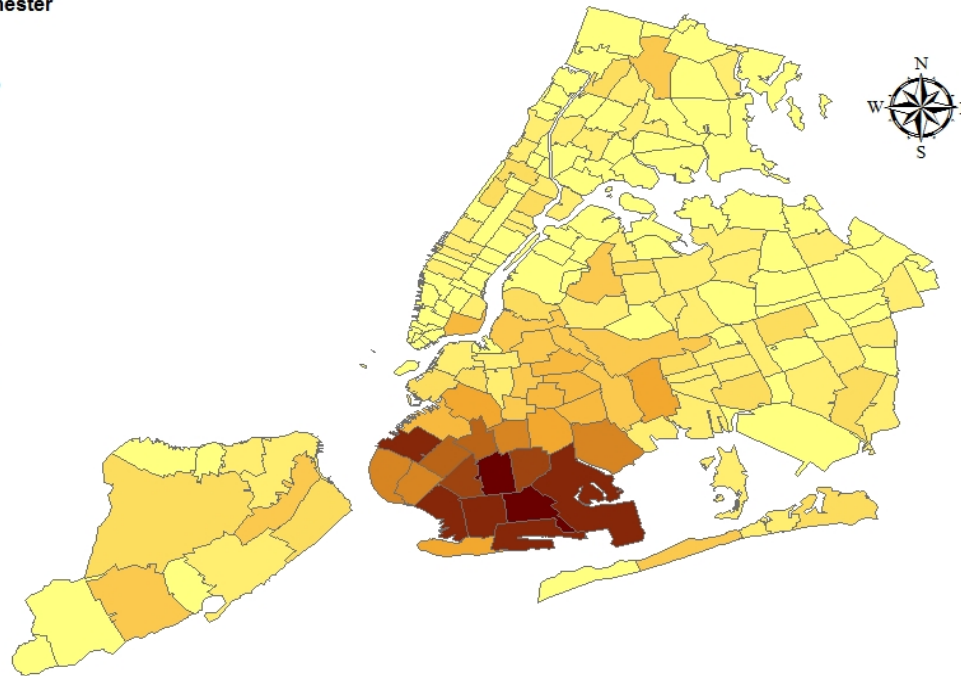
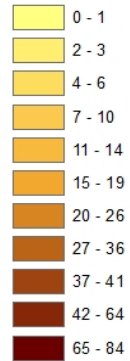
Enrollment at Brooklyn College, Count By Zip Code

Fall 2007 Semester

Legend

NYMetroZIP

enroll



Source: Office of the AVP for Finance, Budget, and Planning / Comptroller, Show Registration Files
Updated: June 16, 2008

Written Report: Rubric

1. **Basic grammar (spelling, punctuation, phrasing)**

- Proofread extensively. Proofreading means reviewing text and tables and charts to make sure they are all correct.
 - Check spelling and punctuation using a dictionary, grammar book. Another option is to consult online sources.
 - Proofread for not only grammar but also the accuracy of the numbers and findings you present. Check for inconsistencies between your text and the tables and charts you've provided. They cannot contradict.

Written Report: Rubric

1. Basic grammar (continued)

- **Avoid passive voice** (“was done by”; “was used for”).
 - Passive voice: “The marketing research was conducted by Judith”.
 - Active voice: “Judith conducted the marketing research”.
- Eliminate extra words: efforts should be made to express yourself in fewer words whenever possible.
- Use the appropriate verb tense consistently.

Written Report: Rubric

2. Format (correct use of citations, bibliography, quoting, paragraphs)

- Paragraphs should have a single main idea, and should start with a topic sentence.
- Paragraphs should not exceed nine printed lines, if possible.

3. Presentation (organization of the report, clarity, sufficient detail)

- Paper should tell a story; start with a problem, and then show all of the decisions and discoveries that went into solving that problem.

Written Report: Rubric

3. Presentation (continued)

- Use Headings and Subheadings: act as signals and signposts to serve as a roadmap for a long report.
- Adhere to any existing form standards (margins, spacing, font size).
- All of the required sections must be included.
- Use of jargon is recommended for this project.

Written Report: Rubric

3a. Quality of the appearance of your graphs and tables.

- Are titles and labels clear, simple, and accurate?
- Are all of the required tables included? Does the text match the tables in content and meaning?
- Did you take the time to convert your SPSS printout into a series of attractive and easily readable charts and tables?

Written Report: Rubric

4. Understanding of the pertinent information related to your topic.

- Are the terms you've used appropriate, accurate, and used correctly?
- All of the required sections must be covered in sufficient detail. Not just length of discussion, but substance, are important.
- Did you apply concepts from the course in your research? Is this clear from the text of the term paper?

Written Report: Rubric

- 5. Ability to critically evaluate the major issues pertaining to your topic.**
 - Does your narrative discussion include concepts from the course? Does it show how these concepts were used to assess results and build conclusions?
 - What alternative choices could have made, given the availability of sufficient funds?
 - Is there a thorough discussion of what methods would be the best possible methods for answering the questions you've tried to answer with your survey?

Oral Report: Before the Presentation

- Create an outline of all of the information you expect to mention.
- From that outline, select a series of key points for emphasis.
- Consider visuals: which information is appropriate for which diagrams, tables, or graphs?
- Select a medium for delivering your presentation: Power Point slides; video; audio.
 - What technology will be needed to assist the presentation? Make sure it is available and in working form in advance.
- Practice the execution of the speech in advance.
 - Eye contact, elocution, clarity, mechanics.
- Professionalism: prepare to have a professional appearance (grooming and attire).
- Arrive early.

Oral Report: Rubric

1. Organization (Weak-Strong)

- Audience cannot understand presentation because there is no sequence of information
- Audience has some difficulty following presentation because student jumps around.
- Student presents information in a logical sequence which audience can follow.
- Student presents information in a logical, interesting sequence which audience can follow.

Oral Report: Rubric

2. **Subject Knowledge (Weak-Strong)**

- Student does not have grasp of information; student cannot answer question about subject.
- Student is somewhat uncomfortable with information and is able to answer only rudimentary questions.
- Student is at ease with expected answers to all questions, but fails to elaborate.
- Student demonstrates full knowledge (more than required).

Oral Report: Rubric

3. Visual Aid (Weak-Strong)

- Student does not refer to or does not have a visual aid.
- Student occasionally uses the visual aid to support presentation.
- Student's visual aid and use of it clearly relates to presentation and improves audience understanding.
- Student's visual aid and use of it is exceptional.

Oral Report: Rubric

4. **Mechanics (Weak-Strong)**

- Presentation has many spelling or grammatical errors.
- Presentation has a few misspellings and/or grammatical errors.
- Presentation has one or two misspellings and/or grammatical errors.
- Presentation has no misspellings and/or grammatical errors.

Oral Report: Rubric

5. Eye Contact (Weak-Strong)

- Student reads all of report with no eye contact.
- Student occasionally uses eye contact, but still reads most of the report.
- Student maintains eye contact most of the time but frequently returns to notes.
- Student maintains eye contact with audience, seldom returning to notes.

Oral Report: Rubric

6. **Elocution (Weak-Strong)**

- Student mumbles, incorrectly pronouncing terms, and speaks too quietly for all of the student hear.
- Student's voice is too low. Student incorrectly pronounces terms frequently. Audience members have difficulty hearing presentation.
- Student's voice is clear. Student pronounces most words correctly. Most audience members can hear presentation.
- Student uses a clear voice and correct pronunciation of terms so that all audience members can hear presentation.

Oral Report: Style Points

- Speak loudly and enunciate. Always use a microphone in a large room.
- Don't speak in a monotone. Also avoid the "valley girl" tendency to speak sentences as if they are questions.
- Make eye contact with the audience.
- Giggling, making jokes, or otherwise goofing around is really a bad idea. Be serious.
- Lettering on slides should be large and legible. Don't put too much information on one slide.
- Language: avoid unprofessional talk like slang and swearing.
- Diagrams, tables, and graphs should be colorful and clearly labeled.

Oral Report: Style Points

- Be sure to cite references when appropriate. Nothing upsets a person in the audience more than to see their work used without credit.
- Dress neatly and professionally. Err on the side of being conservative; the people who care how you dress expect you to dress conservatively.
- If English is not your first language, have a native speaker proof-read your slides, and listen to your talk. Extra practice. Take advantage of the ESL center on campus.
- Think about the first sentence and last sentence of your presentation. If you think you will be nervous, write out the first and last sentence on an index card. Start with a bang, and end with a bang!

Ethics

Ethical Research

- Ethics are norms for conduct in research that distinguish between or acceptable and unacceptable behavior.
- Marketing researchers can expose their employers to risk of legal action or public relations disasters by engaging in unethical research.

Purpose of Ethics in Research

1. Ethics promote the aims of research, such as knowledge, truth, and avoidance of error.
 - For example: prohibitions against fabricating, falsifying, or misrepresenting research data promote the truth and avoid error.
2. Ethical standards promote the values that are essential to collaborative work, such as trust, accountability, mutual respect, and fairness.
 - Research often involves a great deal of cooperation and coordination among many different people in different disciplines and institutions.
 - For example, many ethical norms in research, such as guidelines for authorship, copyright and patenting policies, data sharing policies, and confidentiality rules in peer review, are designed to protect intellectual property interests while encouraging collaboration.
 - Most researchers want to receive credit for their contributions and do not want to have their ideas stolen or disclosed prematurely.

Purpose of Ethics in Research

3. Ethics help to ensure that researchers can be held accountable to the public.
 - Federal policies on research misconduct, on conflicts of interest, on the human subjects protections, and on animal care and use are necessary in order to make sure that researchers who are funded by public money can be held accountable to the public.
4. Ethics in research also help to build public support for research.
 - People more likely to fund research project if they can trust the quality and integrity of research.

Purpose of Ethics in Research

5. Ethics promote a variety of other important moral and social values, such as social responsibility, human rights, animal welfare, compliance with the law, and health and safety.
 - Ethical lapses in research can significantly harm to human and animal subjects, students, and the public.
 - A researcher who fabricates data in a clinical trial may harm or even kill patients, and a researcher who fails to abide by regulations and guidelines relating to radiation or biological safety may jeopardize his health and safety or the health and safety and staff and students.

Ethical Philosophy

- Ethical codes are meant to protect research, researchers, and research subjects from harm.
- Some feel that ethical rules can constrain research that may have important benefits to society.
- Ethical philosophy refers to the degree of importance one places on the risk of harm to the subject vs. benefit to society.

Ethical Philosophy

1. Teleology: Emphasis is placed on the benefits of research to society.
 - Although there may be some harm to research subjects, this research also adds to the collective knowledge of our society, which can lead to social progress.
 - Some harm to subjects may be acceptable if the possible benefits are truly great.
2. Deontology: Emphasis is placed on the rights and protections of the subjects.
 - Researchers should protect subjects from physical or emotional harm at all times no matter what the social benefits of the research.

Ethical Philosophy

- People generally favor deontology.
- Marketing researchers, in particular, lean heavily on the side of deontology.
- Marketing research is not intended to have broad social benefit; and researchers work for companies concerned with public image.
- Academic researchers tend to lean closer to teleology.

Nuremberg Code

- A set of research ethics principles for human experimentation set as a result of the Nuremberg Trials at the end of the Second World War.
- Specifically, they were in response to Nazi human experimentation carried out during the war by individuals such as Dr. Josef Mengele.
- The Nuremberg Code was established to formally denounce and ban potentially injurious research on human subjects without their consent.

Nuremberg Code

1. Respect for persons -- the right to informed consent
 - Informed Consent – Disclosure of who you are, who you are affiliated with, purpose of research, and if and when results might be made available...participants should be willing, and aware.
 - Anonymity – no identifying information about respondent will be solicited; indiv cannot be tied to results
 - Confidentiality – research collects identifying information, but guarantees that no subjects will be individually identifiable, that all your public disclosure of results (tables, reports, and publications) will only discuss findings in the aggregate.

Nuremberg Code

2. Beneficence -- minimization of harm and maximization of benefits.
 - Using coercion to force people to participate in research is not acceptable.
 - Coercion means the use of not only the threat of force, but also the withholding of basic needs and the alleviation of punishment to obtain compliance with the study.
 - Researchers should aim to create research that benefits society.
3. Justice -- equitable distribution of benefits and burdens
 - Risks and benefits in research must be applied equally.
 - Representative samples are not only good research, but ethical-if research entails some level of risk, one group should not be the sole group to bear that risk.

Vulnerable Populations

- *Prisoners*: Under no circumstances is a researcher allowed to even suggest that by participating in research, it will reflect favorably on their record when they go up for parole or early release.
- *Children*: A legal guardian or parent must grant you written permission to conduct your research.
 - Authorities in charge of any grounds or locations which involve your research (such as schools, daycare, or recreation centers) must also give their permission for you to conduct your research.

Informed Consent

Informed Consent can be obtained using a consent statement. Ideally, this would include the following:

- A brief description of the purpose and procedure of the research, including the expected duration
- A statement of any risks, discomforts, or inconveniences associated with participation
- A guarantee of anonymity or at least confidentiality, and an explanation of both.

Informed Consent

- The identification, affiliation, and sponsorship of the research as well as contact information
- A statement that participation is completely voluntary and can be terminated at any time without penalty
- A statement of any alternative procedures that may be used
- A statement of any benefits to the class of subjects involved
- An offer to provide a free copy of a summary of the findings

Getting IRB Approval

- An **institutional review board** (IRB) is a committee that has been formally designated to approve, monitor, and review research with human subjects.
- The purpose of IRB review is to protect the rights and welfare of research subjects.
- In the US, the Office for Human Research Protections (part of HHS) regulations have empowered IRBs to approve, require modifications in planned research prior to approval, or disapprove research.

Getting IRB Approval

- IRB Approval is needed when:
 - Research will actively involve human subjects.
 - Research participants are identifiable.
 - Special permissions are needed (such as parent, principal, or warden).
 - Research will be used for purposes other than internal accounting.

Getting IRB Approval

- Academia: Colleges and universities have powerful IRBs, and a detailed research plan must be submitted and approved before research can begin.
- Even if research does not require IRB approval, IRB should be made aware of the research, and indicate that it is not needed.
- Business World: An IRB may or may not exist. Some larger companies are sensitive to negative press, and thus employ a IRB officers.

Sugging and Frugging

- Sugging and Frugging
- *Sugging*: “Selling Under the Guise of Survey Research”. Using a survey as part of an advertisement or sales pitch.
- *Frugging*: “Fundraising Under the Guise of Survey Research”. Using a survey as part of an effort to solicit contributions to a non-profit company.
- Both are **prohibited by the American Marketing Association code of ethics**.
 - These tactics are detrimental to the general public, who are flooded with junk mail and internet spam.
 - These tactics are also detrimental to the efforts of marketing researchers, as it contributes to the sense among the general public that they are oversurveyed.
 - This may contribute to the general decline in response rates to marketing research that has been observed over the years.