

<b>Part II</b>	<b>Exploring Relationships Between Variables</b>
<b>Chapter 7</b>	<b>Scatterplots, Association, and Correlation</b>
Scatterplot _____ is plotted on the x-axis. _____ is plotted on the y-axis.	Shows the relationship between two quantitative variables on the same cases (individuals). Explanatory ( <i>independent</i> /input) variable Response ( <i>dependent</i> /output) variable
Once we make a scatterplot, we describe association by telling about:	<b>1. Form:</b> straight, curved, no pattern, other? <b>2. Direction:</b> + or – slope? <b>3. Strength:</b> how much scatter {how closely points follow the form} <b>4. Unusual Features:</b> outliers, clusters, subgroups?
_____ is a deliberately vague term describing the relationship between two variables. If positive then _____	Association  increases in one variable generally correspond to increases in the other.
Correlation describes the _____ and _____ of the _____ relationship between two _____ variables, without significant _____	strength direction, linear  quantitative outliers.
3 conditions needed for Correlation:	1. Quantitative Variables 2. Straight Enough 3. Outlier
The correlation coefficient is found by _____.  It's value ranges from _____, it has no _____, and is immune to changes of _____	finding the average product of the z-scores (standardized values). $r = \frac{\sum z_x z_y}{n-1}$ -1 to +1 units. scale or order.
Perfect correlation $r = \pm 1$ , occurs only when _____.	$\pm 1$ the points lie exactly on a straight line. (you can perfectly predict one variable knowing the other)
No correlation $r = 0$ , means that knowing one variable gives you _____	0 no information about the other variable.
You should give the _____ and _____ of x and y along with the correlation because ...	Mean Standard deviation Correlation is not a complete description of two-variable data and the its formula uses means and standard deviations in the z-scores.
Scatterplots and correlation coefficients never prove _____	causation.
Lurking variable	A variable other than x and y that simultaneously affects both variables, accounting for the correlation between the two.
To add a categorical variable to an existing scatterplot _____	use a different plot color or symbol for each category.
<b>Chapter 8</b>	<b>Linear Regression</b>

Regression to the mean	Because the correlation is always less than 1.0 in magnitude, each predicted $\hat{y}$ tends to be fewer standard deviations from its mean than its corresponding $x$ was from its mean. ( $\hat{z}_y = rz_x$ )
Residual If positive If negative	Observed value – predicted value $y - \hat{y}$ Then the model makes an underestimate. Then the model makes an overestimate.
Regression line Line of best fit For standardized values For actual $x$ and $y$ values	The unique line that minimizes the variance of the residuals (sum of the squared residuals). $\hat{z}_y = rz_x$ $\hat{y} = b_0 + b_1x$
To calculate the regression line in real units (actual $x$ and $y$ values)	1. Find slope, $b_1 = \frac{rs_y}{s_x}$ 2. Find y-intercept, plug $b_1$ and point $(x, y)$ [usually $(\bar{x}, \bar{y})$ ] into $\hat{y} = b_0 + b_1x$ and solve for $b_0$ 3. Plug in slope, $b_1$ , and y-intercept, $b_0$ , into $\hat{y} = b_0 + b_1x$
3 conditions needed for Linear Regression Models: /* same as correlation */	1. Quantitative Variables 2. Straight Enough – check original scatterplot & residual scatterplot 3. Outlier (clusters) –points with large residuals and/or high leverage
$R^2$	The square of the correlation, $r$ , between $x$ and $y$ The success of the regression model in terms of the fraction of the variation of $y$ accounted for by the model. (XX% of the variability in $y$ is accounted for by variation in $x$ ) (differences in $x$ explain XX% of the variability in $y$ )
A high $R^2$	Does not demonstrate the appropriateness of the regression.
Looking at a _____ is a good way to check the Straight Enough Condition. It should be _____	a scatterplot of the residuals vs. the $x$ -values.  (appropriateness) boring: uniform scatter with no direction, shape, or outliers..
The ____ is the key to assessing how well the model fits (extracts the form).	variation in the residuals
Standard deviation of the residuals, $s_e$	Gives a measure of how much the points spread around the regression line.
$1 - R^2$	The fraction of the original variation left in the residuals. (The percentage of variability not explained by the regression line.)
Extrapolations	Dubious predictions of $y$ -values based on $x$ -values outside the range of the original data.
<b>Chapter 9</b>	<b>Regression Wisdom</b>
What can go wrong with regression:	1. Inferring Causation 2.Extrapolation 3.Outliers and Influential Points 4.Change in Scatterplot Pattern 5.Means (or other summaries) rather than actual data.
High leverage points With enough leverage the _____	Have $x$ -values far from $\bar{x}$ ( $(\bar{x}, \bar{y})$ is the fulcrum) and pull more strongly on the regression line. residuals

can appear deceptively small.	
Leverage and residual produce three flavors of outliers:	<ol style="list-style-type: none"> <li>1) Extreme Conformers: don't influence model but do inflate R<sup>2</sup></li> <li>2) Large Residuals: might not influence model much but aren't consistent with the overall form.</li> <li>3) Influential Points: those that distort the model</li> </ol>
Influential point [most menacing]	Omitting it from the data results in a very different regression model
Influential points are often difficult to detect because	They distort the model which causes their residual to be small.
The surest way to verify an outlier and its affects is to	Calculate the regression line with and without the suspect point.
A histogram of the residuals	Compliments a scatterplot of the residuals in the search for conditions, such as subsets, that may compromise the effectiveness of the regression model.
Consider comparing two or more regressions if you find	<ol style="list-style-type: none"> <li>1) Points with large residuals and/or high leverage.</li> <li>2) Change in Scatterplot Pattern as a result of changes over time or subsets that behave differently.</li> </ol>
Regressions based on summaries of the data _____ Because _____	Tend to look stronger than the regression on the original data. Summary statistics are less variable than the underlying data.
<b>Chapter 10</b>	<b>Re-expressing Data: Get It Straight!</b>
Re-expression	A means of altering the data to achieve the conditions/structure necessary to utilize particular summaries or models.
Several reasons to consider a re-expression:	<ol style="list-style-type: none"> <li>1. Make the form of a scatterplot straighter.</li> <li>2. Make the scatter in a scatterplot more consistent (not fan shaped).</li> <li>3. Make the distribution of a variable (histogram) more symmetric.</li> <li>4. Make the spread across different groups (boxplots) more similar.</li> </ol>
Ladder of Powers	Orders the effects that the re-expressions have on the data
A good starting point is _____ If all else fails _____	$  \begin{array}{cccccc}  2 & 1 & \frac{1}{2} & 0 & -\frac{1}{2} & -1 \\  y^2 & y & \sqrt{y} & \log y & -1/\sqrt{y} & -1/y  \end{array}  $ <p>taking logs. try whacking the data with two logs (log x and log y).</p>
Base 10 logs are roughly	One less than the number of digits needed to write the number.
Re-expression limitations:	<ol style="list-style-type: none"> <li>1. Can't straighten scatterplots that turn around.</li> <li>2. Can't re-express "-" data values with <math>\sqrt{\quad}</math> (+constant to shift &gt; 0)</li> <li>3. Minimal affect on data values far from 1-100. (-constant to shift)</li> <li>4. Can't unify multiple modes.</li> </ol>
When discussing the accuracy or confidence of the linear regression model be sure to comment on both the ___ & ___	<p>Appropriateness of the model as indicated by the residual plot</p> <p>Success of the model as indicated by R<sup>2</sup></p>