# SIMPLE LINEAR REGRESSION NOTES pdf

## 1: MTH : Regression Analysis

*Simple Linear Regression To describe the linear association between quantitative variables, a statistical procedure called regression often is used to construct a model. Regression is used to assess the contribution of one or more "explanatory" variables (called independent variables) to one "response" (or dependent) variable.*

Correlation and simple regression formulas Linear regression analysis is the most widely used of all statistical techniques: Then the equation for computing the predicted value of Yt is: This formula has the property that the prediction for Y is a straight-line function of each of the X variables, holding the others fixed, and the contributions of different X variables to the predictions are additive. The slopes of their individual straight-line relationships with Y are the constants b1, b2, â€¦, bk, the so-called coefficients of the variables. That is, bi is the change in the predicted value of Y per unit of change in Xi, other things being equal. The coefficients and intercept are estimated by least squares, i. The first thing you ought to know about linear regression is how the strange term regression came to be applied to models like this. They were first studied in depth by a 19th-Century scientist, Sir Francis Galton. Galton was a self-taught naturalist, anthropologist, astronomer, and statistician--and a real-life Indiana Jones character. He was famous for his explorations, and he wrote a best-selling book on how to survive in the wilderness entitled "The Art of Travel: From the Practical to the Peculiar. They provide many handy hints for staying alive--such as how to treat spear wounds or extract your horse from quicksand--and introduced the concept of the sleeping bag to the Western World. Click on these pictures for more details: Galton was a pioneer in the application of statistical methods to measurements in many branches of science, and in studying data on relative sizes of parents and their offspring in various species of plants and animals, he observed the following phenomenon: The same is true of virtually any physical measurement and in the case of humans, most measurements of cognitive and physical ability that can be performed on parents and their offspring. Here is the first published picture of a regression line illustrating this effect, from a lecture presented by Galton in  The R symbol on this chart whose value is 0. Galton termed this phenomenon a regression towards mediocrity , which in modern terms is a regression to the mean. It is a purely statistical phenomenon. Unless every child is exactly as the same size as the parent in relative terms i. Return to top of page. Regression to the mean is an inescapable fact of life. Your children can be expected to be less exceptional for better or worse than you are. Your score on a final exam in a course can be expected to be less good or bad than your score on the midterm exam, relative to the rest of the class. The key word here is "expected. We have already seen a suggestion of regression-to-the-mean in some of the time series forecasting models we have studied: This is not true of random walk models, but it is generally true of moving-average models and other models that base their forecasts on more than one past observation. The intuitive explanation for the regression effect is simple: The best we can hope to do is to predict only that part of the variability which is due to the signal. Hence our forecasts will tend to exhibit less variability than the actual values, which implies a regression to the mean. Another way to think of the regression effect is in terms of selection bias. Suppose that we select a sample of professional athletes whose performance was much better than average or students whose grades were much better than average in the first half of the year. The fact that they did so well in the first half of the year makes it probable that both their skill and their luck were better than average during that period. In the second half of the year we may expect them to be equally skillful, but we should not expect them to be equally lucky. So we should predict that in the second half their performance will be closer to the mean. Meanwhile, players whose performance was merely average in the first half probably had skill and luck working in opposite directions for them. We should therefore expect their performance in the second half to move away from the mean in one direction or another, as we get another independent test of their skill. However, the actual performance of the players should be expected to have an equally large variance in the second half of the year as in the first half, because it merely results from a redistribution of independently random luck among players with the same distribution of skill as before. A nice discussion of regression to the mean in the broader context of social science research can be found here. Justification for regression assumptions Why should we assume that relationships between variables are

linear? Because linear relationships are the simplest non-trivial relationships that can be imagined hence the easiest to work with , and Because the "true" relationships between our variables are often at least approximately linear over the range of values that are of interest to us, and This is a strong assumption, and the first step in regression modeling should be to look at scatterplots of the variables and in the case of time series data, plots of the variables vs. And after fitting a model, plots of the errors should be studied to see if there are unexplained nonlinear patterns. This is especially important when the goal is to make predictions for scenarios outside the range of the historical data, where departures from perfect linearity are likely to have the biggest effect. If you see evidence of nonlinear relationships, it is possible though not guaranteed that transformations of variables will straighten them out in a way that will yield useful inferences and predictions via linear regression. And why should we assume that the effects of different independent variables on the expected value of the dependent variable are additive? This is a very strong assumption, stronger than most people realize. It implies that the marginal effect of one independent variable i. In a multiple regression model, the estimated coefficient of a given independent variable supposedly measures its effect while "controlling" for the presence of the others. However, the way in which controlling is performed is extremely simplistic: Many users just throw a lot of independent variables into the model without thinking carefully about this issue, as if their software will automatically figure out exactly how they are related. Even automatic model-selection methods e. They work only with the variables they are given, in the form that they are given, and then they look only for linear, additive patterns among them in the context of each other. A common practice is to include independent variables whose predictive effects logically cannot be additive, say, some that are totals and others that are rates or percentages. You need to collect the relevant data, understand what it measures, clean it up if necessary, perform descriptive analysis to look for patterns before fitting any models, and study the diagnostic tests of model assumptions afterward, especially statistics and plots of the errors. You should also try to apply the appropriate economic or physical reasoning to determine whether an additive prediction equation makes sense. Here too, it is possible but not guaranteed that transformations of variables or the inclusion of interaction terms might separate their effects into an additive form, if they do not have such a form to begin with, but this requires some thought and effort on your part. And why should we assume the errors of linear models are independently and identically normally distributed? This assumption is often justified by appeal to the Central Limit Theorem of statistics, which states that the sum or average of a sufficiently large number of independent random variables--whatever their individual distributions--approaches a normal distribution. Much data in business and economics and engineering and the natural sciences is obtained by adding or averaging numerical measurements performed on many different persons or products or locations or time intervals. Insofar as the activities that generate the measurements may occur somewhat randomly and somewhat independently, we might expect the variations in the totals or averages to be somewhat normally distributed. It is again mathematically convenient: This family includes the t distribution, the F distribution, and the Chi-square distribution. But here too caution must be exercised. Even if the unexplained variations in the dependent variable are approximately normally distributed, it is not guaranteed that they will also be identically normally distributed for all values of the independent variables. Perhaps the unexplained variations are larger under some conditions than others, a condition known as "heteroscedasticity". For example, if the dependent variable consists of daily or monthly total sales, there are probably significant day-of-week patterns or seasonal patterns. In such cases the variance of the total will be larger on days or in seasons with greater business activity--another consequence of the central limit theorem. It is also not guaranteed that the random variations will be statistically independent. This is an especially important question when the data consists of time series: A very important special case is that of stock price data, in which percentage changes rather than absolute changes tend to be normally distributed. This implies that over moderate to large time scales, movements in stock prices are lognormally distributed rather than normally distributed. A log transformation is typically applied to historical stock price data when studying growth and volatility. See the geometric random walk page instead. You still might think that variations in the values of portfolios of stocks would tend to be normally distributed, by virtue of the central limit theorem, but the central limit theorem is actually rather slow to bite on the lognormal distribution because it is so

asymmetrically long-tailed. A sum of 10 or 20 independently and identically lognormally distributed variables has a distribution that is still quite close to lognormal. Because the assumptions of linear regression linear, additive relationships with i. Its output contains no more information than is provided by its inputs, and its inner mechanism needs to be compared with reality in each situation where it is applied. Correlation and simple regression formulas A variable is, by definition, a quantity that may vary from one measurement to another in situations where different samples are taken from a population or observations are made at different points in time. In fitting statistical models in which some variables are used to predict others, what we hope to find is that the different variables do not vary independently in a statistical sense , but that they tend to vary together. In particular, when fitting linear models, we hope to find that one variable say, Y is varying as a straight-line function of another variable say, X. In other words, if all other possibly-relevant variables could be held fixed, we would hope to find the graph of Y versus X to be a straight line apart from the inevitable random errors or "noise". A measure of the absolute amount of variability in a variable is naturally its variance, which is defined as its average squared deviation from its own mean. Equivalently, we can measure variability in terms of the standard deviation, which is defined as the square root of the variance. The standard deviation has the advantage that it is measured in the same units as the original variable, rather than squared units. Our task in predicting Y might be described as that of explaining some or all of its variance--i. Why is it not constant? That is, we would like to be able to improve on the naive predictive model: More precisely, we hope to find a model whose prediction errors are smaller, in a mean square sense, than the deviations of the original variable from its mean. In using linear models for prediction, it turns out very conveniently that the only statistics of interest at least for purposes of estimating coefficients to minimize squared error are the mean and variance of each variable and the correlation coefficient between each pair of variables. The coefficient of correlation between X and Y is commonly denoted by rXY, and it measures the strength of the linear relationship between them on a relative i. The correlation coefficient is most easily computed if we first standardize the variables, which means to convert them to units of standard-deviations-from-the-mean, using the population standard deviation rather than the sample standard deviation, i. P is the Excel function for the population standard deviation. Here and elsewhere I am going to use Excel functions rather than conventional math symbols in some of the formulas to illustrate how the calculations would be done on a spreadsheet. Now, the correlation coefficient is equal to the average product of the standardized values of the two variables within the given sample of n observations: The average of the values in the last column is the correlation between X and Y. S in Excel, but the population statistic is the correct one to use in the formula above. If the two variables tend to vary on the same sides of their respective means at the same time, then the average product of their deviations and hence the correlation between them will be positive, since the product of two numbers with the same sign is positive.

*NOTES ON SIMPLE LINEAR REGRESSION 1. INTRODUCTION The purpose of these notes is to supplement the mathematical development of linear regression in Devore (). This.*

The estimated value based on the fitted regression model for the new observation at is: The prediction interval values calculated in this example are shown in the figure below as Low Prediction Interval and High Prediction Interval, respectively. The columns labeled Mean Predicted and Standard Error represent the values of and the standard error used in the calculations. Measures of Model Adequacy It is important to analyze the regression model before inferences based on the model are undertaken. The following sections present some techniques that can be used to check the appropriateness of the model for the given data. These techniques help to determine if any of the model assumptions have been violated. Coefficient of Determination R2 The coefficient of determination is a measure of the amount of variability in the data accounted for by the regression model. As mentioned previously, the total variability of the data is measured by the total sum of squares,. The amount of this variability explained by the regression model is the regression sum of squares,. The coefficient of determination is the ratio of the regression sum of squares to the total sum of squares. For the yield data example, can be calculated as: It may appear that larger values of indicate a better fitting regression model. However, should be used cautiously as this is not always the case. The value of increases as more terms are added to the model, even if the new term does not contribute significantly to the model. Therefore, an increase in the value of cannot be taken as a sign to conclude that the new model is superior to the older model. Adding a new term may make the regression model worse if the error mean square, , for the new model is larger than the of the older model, even though the new model will show an increased value of. In the results obtained from the DOE folio, is displayed as R-sq under the ANOVA table as shown in the figure below , which displays the complete analysis sheet for the data in the preceding table. These values measure different aspects of the adequacy of the regression model. For example, the value of S is the square root of the error mean square, , and represents the "standard error of the model. The values of S, R-sq and R-sq adj indicate how well the model fits the observed data. Residual Analysis In the simple linear regression model the true error terms, , are never known. The residuals, , may be thought of as the observed error terms that are similar to the true error terms. Since the true error terms, , are assumed to be normally distributed with a mean of zero and a variance of , in a good model the observed error terms i. Thus the residuals in the simple linear regression should be normally distributed with a mean of zero and a constant variance of. Residuals are usually plotted against the fitted values, , against the predictor variable values, , and against time or run-order sequence, in addition to the normal probability plot. Plots of residuals are used to check for the following: Residuals follow the normal distribution. Residuals have a constant variance. Regression function is linear. A pattern does not exist when residuals are plotted in a time or run-order sequence. There are no outliers. Examples of residual plots are shown in the following figure. Such a plot indicates an appropriate regression model. Such a plot indicates increase in variance of residuals and the assumption of constant variance is violated here. Transformation on may be helpful in this case see Transformations. If the residuals follow the pattern of c or d , then this is an indication that the linear regression model is not adequate. Addition of higher order terms to the regression model or transformation on or may be required in such cases. A plot of residuals may also show a pattern as seen in e , indicating that the residuals increase or decrease as the run order sequence or time progresses. This may be due to factors such as operator-learning or instrument-creep and should be investigated further. Example Residual plots for the data of the preceding table are shown in the following figures. One of the following figures is the normal probability plot. It can be observed that the residuals follow the normal distribution and the assumption of normality is valid here. In one of the following figures the residuals are plotted against the fitted values, , and in one of the following figures the residuals are plotted against the run order. Both of these plots show that the 21st observation seems to be an outlier. Further investigations are needed to study the cause of this outlier. This perfect model will give us a zero error sum of squares. Thus, no error exists for the perfect model. However, if you record the response values for the same

values of for a second time, in conditions maintained as strictly identical as possible to the first time, observations from the second time will not all fall along the perfect model. The deviations in observations recorded for the second time constitute the "purely" random variation or noise. The sum of squares due to pure error abbreviated quantifies these variations.

## 3: Introduction to linear regression analysis

*Simple Linear Regression Notes. Relationships. Estimating the Simple Linear Function. Measures of Variation. Assumptions. Assumption Checks. Slope. Estimate Averages.*

First, compute the estimate of the slope: Now the intercept may be computed: When you plot this line over the data points, the result looks like that shown in Figure 2. Some points have positive residuals they lie above the line ; some have negative ones they lie below it. If all the points fell on the line, there would be no error and no residuals. The mean of the sample residuals is always 0 because the regression line is always drawn such that half of the error is above it and half below it. This method of regression is called least squares. The slope is negative because the line slants down from left to right, as it must for two variables that are negatively correlated, reflecting that one variable decreases as the other increases. Confidence interval for the slope Example 1 What if the slope is 0, as in Figure 3? That means that y has no linear dependence on x, or that knowing x does not contribute anything to your ability to predict y. It is often useful to compute a confidence interval for a regression slope. If it contained 0, you would be unable to conclude that x and y are related. The test for this example will use an alpha of 0. Table 3 in "Statistics Tables" shows that t. An example of uncorrelated data, so the slope is zero. Because there is some error associated with your prediction, however, you might want to produce a confidence interval rather than a simple point estimate. What is a 90 percent confidence interval for the number of hours spent exercising per week if the exercise machine is owned 11 months? The first step is to use the original regression equation to compute a point estimate for y: For a 90 percent confidence interval, you need to use t. You have already computed the remaining quantities, so you can proceed with the formula: Figure 4 shows a violation of the second assumption. The errors residuals are greater for higher values of x than for lower values. Least squares regression is sensitive to outliers, or data points that fall far from most other points. You need to be wary of outliers because they can influence the regression equation greatly. Least squares regression is sensitive to outliers. It can be dangerous to extrapolate in regression—to predict values beyond the range of your data set. The regression model assumes that the straight line extends to infinity in both directions, which often is not true. According to the regression equation for the example, people who have owned their exercise machines longer than around 15 months do not exercise at all. Extrapolation beyond the data is dangerous.

## 4: Introduction to Linear Regression

*Simple linear regression is a statistical method that allows us to summarize and study relationships between two continuous (quantitative) variables. This lesson introduces the concept and basic procedures of simple linear regression. We will also learn two measures that describe the strength of the.*

Numerous extensions have been developed that allow each of these assumptions to be relaxed i. Generally these extensions make the estimation procedure more complex and time-consuming, and may also require more data in order to produce an equally precise model. Example of a cubic polynomial regression, which is a type of linear regression. The following are the major assumptions made by standard linear regression models with standard estimation techniques e. This essentially means that the predictor variables x can be treated as fixed values, rather than random variables. This means, for example, that the predictor variables are assumed to be error-freeâ€"that is, not contaminated with measurement errors. Although this assumption is not realistic in many settings, dropping it leads to significantly more difficult errors-in-variables models. This means that the mean of the response variable is a linear combination of the parameters regression coefficients and the predictor variables. Note that this assumption is much less restrictive than it may at first seem. Because the predictor variables are treated as fixed values see above , linearity is really only a restriction on the parameters. The predictor variables themselves can be arbitrarily transformed, and in fact multiple copies of the same underlying predictor variable can be added, each one transformed differently. This trick is used, for example, in polynomial regression , which uses linear regression to fit the response variable as an arbitrary polynomial function up to a given rank of a predictor variable. This makes linear regression an extremely powerful inference method. In fact, models such as polynomial regression are often "too powerful", in that they tend to overfit the data. As a result, some kind of regularization must typically be used to prevent unreasonable solutions coming out of the estimation process. Common examples are ridge regression and lasso regression. Bayesian linear regression can also be used, which by its nature is more or less immune to the problem of overfitting. In fact, ridge regression and lasso regression can both be viewed as special cases of Bayesian linear regression, with particular types of prior distributions placed on the regression coefficients. This means that different values of the response variable have the same variance in their errors, regardless of the values of the predictor variables. In practice this assumption is invalid i. This is to say there will be a systematic change in the absolute or squared residuals when plotted against the predictive variables. Errors will not be evenly distributed across the regression line. Heteroscedasticity will result in the averaging over of distinguishable variances around the points to get a single variance that is inaccurately representing all the variances of the line. In effect, residuals appear clustered and spread apart on their predicted plots for larger and smaller values for points along the linear regression line, and the mean squared error for the model will be wrong. Typically, for example, a response variable whose mean is large will have a greater variance than one whose mean is small. In fact, as this shows, in many casesâ€"often the same cases where the assumption of normally distributed errors failsâ€"the variance or standard deviation should be predicted to be proportional to the mean, rather than constant. Simple linear regression estimation methods give less precise parameter estimates and misleading inferential quantities such as standard errors when substantial heteroscedasticity is present. However, various estimation techniques e. Bayesian linear regression techniques can also be used when the variance is assumed to be a function of the mean. It is also possible in some cases to fix the problem by applying a transformation to the response variable e. This assumes that the errors of the response variables are uncorrelated with each other. Actual statistical independence is a stronger condition than mere lack of correlation and is often not needed, although it can be exploited if it is known to hold. Bayesian linear regression is a general way of handling this issue. Lack of perfect multicollinearity in the predictors. For standard least squares estimation methods, the design matrix X must have full column rank p; otherwise, we have a condition known as perfect multicollinearity in the predictor variables. This can be triggered by having two or more perfectly correlated predictor variables e. It can also happen if there is too little data available compared to the number of parameters to be estimated e. At most we will be able to identify some of the

parameters, i. See partial least squares regression. Methods for fitting linear models with multicollinearity have been developed; [5] [6] [7] [8] some require additional assumptions such as "effect sparsity"â€"that a large fraction of the effects are exactly zero. Note that the more computationally expensive iterated algorithms for parameter estimation, such as those used in generalized linear models , do not suffer from this problem. Beyond these assumptions, several other statistical properties of the data strongly influence the performance of different estimation methods: The statistical relationship between the error terms and the regressors plays an important role in determining whether an estimation procedure has desirable sampling properties such as being unbiased and consistent. This illustrates the pitfalls of relying solely on a fitted model to understand the relationship between variables. A fitted linear regression model can be used to identify the relationship between a single predictor variable $x_j$ and the response variable $y$ when all the other predictor variables in the model are "held fixed". This is sometimes called the unique effect of $x_j$ on $y$. In contrast, the marginal effect of $x_j$ on $y$ can be assessed using a correlation coefficient or simple linear regression model relating only $x_j$ to $y$; this effect is the total derivative of $y$ with respect to $x_j$. Care must be taken when interpreting regression results, as some of the regressors may not allow for marginal changes such as dummy variables , or the intercept term , while others cannot be held fixed recall the example from the introduction: It is possible that the unique effect can be nearly zero even when the marginal effect is large. This may imply that some other covariate captures all the information in $x_j$, so that once that variable is in the model, there is no contribution of $x_j$ to the variation in $y$. Conversely, the unique effect of $x_j$ can be large while its marginal effect is nearly zero. This would happen if the other covariates explained a great deal of the variation of $y$, but they mainly explain variation in a way that is complementary to what is captured by $x_j$. In this case, including the other variables in the model reduces the part of the variability of $y$ that is unrelated to $x_j$, thereby strengthening the apparent relationship with $x_j$. The meaning of the expression "held fixed" may depend on how the values of the predictor variables arise. If the experimenter directly sets the values of the predictor variables according to a study design, the comparisons of interest may literally correspond to comparisons among units whose predictor variables have been "held fixed" by the experimenter. Alternatively, the expression "held fixed" can refer to a selection that takes place in the context of data analysis. In this case, we "hold a variable fixed" by restricting our attention to the subsets of the data that happen to have a common value for the given predictor variable. This is the only interpretation of "held fixed" that can be used in an observational study. The notion of a "unique effect" is appealing when studying a complex system where multiple interrelated components influence the response variable. In some cases, it can literally be interpreted as the causal effect of an intervention that is linked to the value of a predictor variable. However, it has been argued that in many cases multiple regression analysis fails to clarify the relationships between the predictor variables and the response variable when the predictors are correlated with each other and are not assigned following a study design. Simple and multiple linear regression[ edit ] Example of simple linear regression , which has one independent variable The very simplest case of a single scalar predictor variable $x$ and a single scalar response variable $y$ is known as simple linear regression. Nearly all real-world regression models involve multiple predictors, and basic descriptions of linear regression are often phrased in terms of the multiple regression model. Note, however, that in these cases the response variable $y$ is still a scalar. Another term, multivariate linear regression, refers to cases where $y$ is a vector, i. General linear models[ edit ] The general linear model considers the situation when the response variable is not a scalar for each observation but a vector, $y_i$. Conditional linearity of $E$.

## 5: Correlation & Simple Linear Regression

*Simple Linear Regression We have been introduced to the notion that a categorical variable could depend on different levels of another variable when we discussed contingency tables.*

Simple linear regression is a statistical method that allows us to summarize and study relationships between two continuous quantitative variables: One variable, denoted x, is regarded as the predictor, explanatory, or independent variable. The other variable, denoted y, is regarded as the response, outcome, or dependent variable. The other terms are mentioned only to make you aware of them should you encounter them in other arenas. Simple linear regression gets its adjective "simple," because it concerns the study of only one predictor variable. In contrast, multiple linear regression, which we study later in this course, gets its adjective "multiple," because it concerns the study of two or more predictor variables. Here is an example of a deterministic relationship. Note that the observed x, y data points fall directly on a line. As you may remember, the relationship between degrees Fahrenheit and degrees Celsius is known to be: Here are some examples of other deterministic relationships that students from previous semesters have shared: For each of these deterministic relationships, the equation exactly describes the relationship between the two variables. This course does not examine deterministic relationships. Instead, we are interested in statistical relationships, in which the relationship between the variables is not perfect. Here is an example of a statistical relationship. The response variable y is the mortality due to skin cancer number of deaths per 10 million people and the predictor variable x is the latitude degrees North at the center of each of 49 states in the U. You might anticipate that if you lived in the higher latitudes of the northern U. The scatter plot supports such a hypothesis. There appears to be a negative linear relationship between latitude and mortality due to skin cancer, but the relationship is not perfect. Indeed, the plot exhibits some "trend," but it also exhibits some "scatter. Some other examples of statistical relationships might include:

## 6: Simple linear regression - Wikipedia

*Regression analysis is the art and science of fitting straight lines to patterns of data. In a linear regression model, the variable of interest (the so-called "dependent" variable) is predicted from k.*

The line of best fit would be the straight line that connects the predicted values. Simple linear equations like the one described above always have the following components: If one extends the line of prediction back done to an "X" value of 0, it would cross the Y-axis at this value. This tells us the the additional savings that is predicted for each additional week. The parameters for intercept and slope enable us to create an equation that describes the relationship between time and saving, and this allows us to make predictions about total savings at various points in time. In this case, the equation would be: If we wanted to predict total savings after 10 weeks, we simply use 10 as the value for "X" in the equation as follows: XLS worksheet on correlation and linear regression allows you to change the values of "X" and "Y", and it will then recalculate the slope, the Y-intercept, and the correlation coefficient. Additional Helpful Notes The correlation coefficient "r" The correlation coefficient, "r," gives an indication of how tightly the observed data points conform to the line of perfect prediction. The value of "r" ranges from 1 to An "r" value of 1 indicates a perfect positive correlation, as depicted in this graph. In contrast, the next graph has points that are more widely scattered about the line of "best fit," but the slope of the line is still upward, and the "r" value is only 0. The next graph has an "r" value of Finally, this next graph also shows a negative correlation, but the individual points fit much more closely to the line of prediction, and the r-value is The closer "r" gets to a value of 0, the poorer the fit between the data and the line of prediction. An "r" value of 0 or close to 0 would indicate virtually no apparent correlation between the "X" variable and the "Y" variable. For example, you may observe a reasonable correlation between the observed data and the line of prediction, but r2 values less than 1 suggest that there are other independent variables that have an impact on the dependent variable that you are trying to predict. The Slope The slope of the line is a measure of its steepness. Another way of thinking about slope is that it indicates, on average, how much the "Y" variable changes for each incremental change in the "X" variable. Non-linear Relationships Many relationships are not linear. They can follow many patterns. For more information on correlation and regression methods, see the online learning module on Correlation and Linear Regression.

## 7: Quiz: Simple Linear Regression

*Regression is a method for studying the relationship between two or more quantitative variables Simple linear regression (SLR): One quantitative dependent variable.*

## 8: Simple Linear Regression

*For one independent variable (Simple Linear Regression): The mean value of one variable, Y, depends on the values of another, X. For example, the average starting salary depends on a student's GPA. There are many relationship functions but the simplest is a straight line.*

## 9: Lesson 1: Simple Linear Regression | STAT

*Simple linear regression is a statistical method that allows us to summarize and study relationships between two continuous (quantitative) variables: One variable, denoted x, is regarded as the predictor, explanatory, or independent variable.*

*Carolina curiosities Disapproving the extension of fast-track procedures to bills to implement trade agreements entered into a Deleuzes Nietzsche Jon Roffe Credit risk analysis case study filetype Care for a pet chimpanzee Syntactic Theory and First Language Acquisition Uf0d8/tKeeping on Biogas lamp burning for 8 hours Introduction to Windows and Graphics Programming with Visual C .Net Bar snacks : food for drink What do libraries really do with electronic resources? : the practice in 2003 Jay Shorten Exercise programming for older adults How we found America Lizard to Start With (Workbook) First 25 years, 1952-1977 Urban Growth Centres Strategy in the greater golden horseshoe Epidemiology of the global pandemic Sten H. Vermund and Sheetal Khedkar Bloodletting in Appalachia Ethics in law enforcement India My beloved charioteer Shashi Deshpande I Cant Come Down, Ive Been Set Free Hammond Barnhart dictionary of science Java application architecture modularity patterns with examples using osgi Open seasons for game and fish Out of breath : how can I slow down? Muscle Energy Techniques with DVD-ROM Lautreamont and Sade (Meridian: Crossing Aesthetics) Pt. 1 Gray, E.D. The question constitutionally considered. Distributing and receiving, inside, and outside. Regulatory Chemicals Handbook (Chemical Industries, V. 80) Guide to plant families of southern africa Truth about homosexuality Mla citation practice worksheet Methods of arterial and venous assessment Peter Gorman, Mario De Nunzio, Richard Donnelly The Fascinating Kings Gambit Wonder palacio Copeland refrigeration manual part 3 The heroic legend of arslan japanese novel The millionaire makeover naima simone US foreign policy in the 1990s Jane eyre chapter 20*