

1: Linear Regression With R

This tutorial follows a data analysis problem typical of earth sciences, natural and water resources, and agriculture, proceeding from visualisation and exploration through univariate point estimation, bivariate correlation and regression analysis.

Contents Linear Regression Linear regression is used to predict the value of an outcome variable Y based on one or more input predictor variables X . The aim is to establish a linear relationship a mathematical formula between the predictor variable s and the response variable, so that, we can use this formula to estimate the value of the response Y , when only the predictors X s values are known. Introduction The aim of linear regression is to model a continuous variable Y as a mathematical function of one or more X variable s , so that we can use this regression model to predict the Y when only the X is known. This mathematical equation can be generalized as follows: Collectively, they are called regression coefficients. Example Problem For this analysis, we will use the cars dataset that comes with R by default. You can access this dataset simply by typing in cars in your R console. You will find that it consists of 50 observations rows and 2 variables columns "dist and speed. Lets print out the first six observations here.. The graphical analysis and correlation study below will help with this. Graphical Analysis The aim of this exercise is to build a simple regression model that we can use to predict Distance dist by establishing a statistically significant linear relationship with Speed speed. But before jumping in to the syntax, lets try to understand these variables graphically. Typically, for each of the independent variables predictors, the following plots are drawn to visualize the following behavior: Visualize the linear relationship between the predictor and response Box plot: To spot any outlier observations in the variable. To see the distribution of the predictor variable. Ideally, a close to normal distribution a bell shaped curve, without being skewed to the left or right is preferred. Let us see how to make each one of them. Scatter Plot Scatter plots can help visualize any linear relationships between the dependent response variable and independent predictor variables. Ideally, if you are having multiple predictor variables, a scatter plot is drawn for each one of them against the response, along with the line of best as seen below. This is a good thing, because, one of the underlying assumptions in linear regression is that the relationship between the response and predictor variables is linear and additive. BoxPlot " Check for outliers Generally, any datapoint that lies outside the 1. If we observe for every instance where speed increases, the distance also increases along with it, then there is a high positive correlation between them and therefore the correlation between them will be closer to 1. The opposite is true for an inverse relationship, in which case, the correlation between the variables will be close to 0. A value closer to 0 suggests a weak relationship between the variables. A low correlation The function used for building linear models is lm. The lm function takes in two main arguments, namely: The data is typically a data. But the most common convention is to write out the formula directly in place of the argument as written below. Is this enough to actually use this model? Before using a regression model, you have to ensure that it is statistically significant. How do you ensure this? Lets begin by printing the summary statistics for linearMod. Checking for statistical significance The summary statistics above tells us a number of things. The p-Values are very important because, We can consider a linear model to be statistically significant only when both these p-Values are less than the pre-determined statistical significance level, which is ideally 0. This is visually interpreted by the significance stars at the end of the row. Null and alternate hypothesis When there is a p-value, there is a null and alternative hypothesis associated with it. In Linear Regression, the Null Hypothesis is that the coefficients associated with the variables is equal to zero. The alternate hypothesis is that the coefficients are not equal to zero. A larger t-value indicates that it is less likely that the coefficient is not equal to zero purely by chance. So, higher the t-value, the better. What this means to us? In our case, linearMod, both these p-Values are well below the 0. It is absolutely important for the model to be statistically significant before we can go ahead and use it to predict or estimate the dependent variable, otherwise, the confidence in predicted values from that model reduces and may be construed as an event of chance. How to calculate the t Statistic and p-Values? When the model co-efficients and standard error are known, the formula for calculating t Statistic and p-Value

is as follows: What R-Squared tells us is the proportion of variation in the dependent response variable that has been explained by this model. Its a better practice to look at the AIC and prediction accuracy on validation sample when deciding on the efficacy of a model. Now thats about R-Squared. What about adjusted R-Squared? As you add more X variables to your model, the R-Squared value of the new bigger model will always be greater than that of the smaller subset. This is because, since all the variables in the original model is also present, their contribution to explain the dependent variable will be present in the super-set as well, therefore, whatever new variable we add can only add if not significantly to the variation that was already explained. It is here, the adjusted R-Squared value comes to help. Adj R-Squared penalizes total value for the number of terms read predictors in your model. Therefore when comparing nested models, it is a good practice to look at adj-R-squared value over R-squared. Therefore, by moving around the numerators and denominators, the relationship between R2 and Radj2 becomes: Both criteria depend on the maximized value of the likelihood function L for the estimated model. The AIC is defined as: The most common metrics to look at while selecting the model are:

2: Time Series Analysis Using ARIMA Model In R | DataScience+

Using R for Data Analysis and Graphics Introduction, Code and Commentary J H Maindonald Centre for Mathematics and Its Applications, Australian National University.

3: R regression models workshop notes

Regression Models for Data Science in R A companion book for the Coursera Regression Models class. Multivariable regression analysis The linear model.

Always theres a thud Ggratitude your day The Fantastic Art of Boris Vallejo Members-only collective bargaining Fire engineer oral exam study guide Pleurocarpous mosses W Tetzlaff and J D Steeves Importance of ayurveda in modern life Tissue Engineering II (Advances in Biochemical Engineering Biotechnology) Chocolate surprise Scandanavian airlines green engine case study Plantation agriculture and social control in northern Peru, 1875-1933 Two hours todarkness The Fail-Proof Enterprise Label a clock 2nd grade Against the Drimlith Alabama related laws to the Insurance Code. ISSE 2005 Securing Electronic Business Processes Catcher in the rye chapter 19 Spirit of Chinese philosophy H.R. 4570, to improve and strengthen the child support collection system Instructors resource guide to accompany Understanding organizational behavior Courage on the Causeway (Cover-to-Cover Novels: Adventure) Methods in personality assessment Physics of climate filetype Rakhaldas Bandyopadhyay The case of the Evanstons. Extreme and chronic poverty and malnutrition in India R. Radhakrishna, K. Hanumantha Rao, C. Ravi and B. Geometry of PDEs and mechanics Drinking water treatment plant design Raid/untold Sty Patto Politics in the age of Cobden The prison in heaven. Complications in Arterial Surgery Truth in Marketing 1995 dodge dakota service manual Training his mate samantha madisen Pass among the stars Continuing national emergency with respect to Nicaragua In the Shadow of Organization