

# BOOSTEDML

Articles on Statistics and Machine Learning for Healthcare

- › general
- › hypothesis test
- › linear and generalized linear models
- › Longitudinal Data Analysis
- › Missing Data
- › Neural Networks+Deep Learning
- › Numerical Linear Algebra
- › optimization
- › Survival
- › Time series
- › Uncategorized



[LINEAR AND GENERALIZED LINEAR MODELS](#)

# Linear Regression Summary(lm): Interpreting

POSTED ON [JUNE 1, 2019](#) BY [ALEX](#)**Contents** [ [hide](#) ]

- 1 [Introduction to Linear Regression Summary Printouts](#)
- 2 [Residual Summary Statistics](#)
- 3 [Coefficients](#)
  - 3.1 [Estimates](#)
  - 3.2 [Standard Error](#)
  - 3.3 [t-value](#)
  - 3.4 [Pr\(>|t|\)](#)
- 4 [Assessing Fit and Overall Significance](#)
  - 4.1 [Residual Standard Error](#)
  - 4.2 [Multiple and Adjusted](#)
  - 4.3 [F-Statistic and F-test](#)

## Introduction to Linear Regression Summary Printouts

In this post we describe how to interpret the summary of a linear regression model in `summary(lm)`. We discuss interpretation of the residual quantiles and summary statistics, standard errors and t statistics, along with the p-values of the latter, the residual standard error and the F-test. Let's first load the Boston housing dataset and fit a naive model. We will then discuss assumptions, which are described in other posts.

```

1 library(mlbench)
2 data(BostonHousing)
3 model<-lm(log(medv) ~ crim + rm + tax + lstat, data = BostonHousing)
4 summary(model)
5
6 Call:
7 lm(formula = log(medv) ~ crim + rm + tax + lstat, data = BostonHousing)
8
9 Residuals:
10      Min       1Q   Median       3Q      Max
11  -0.72730  -0.13031  -0.01628   0.11215   0.92987
12
13 Coefficients:

```

### Categories

- > [general](#)
- > [hypothesis test](#)
- > [linear and generalized linear models](#)
- > [Longitudinal Data Analysis](#)
- > [Missing Data](#)
- > [Neural Networks+Deep Learning](#)
- > [Numerical Linear Algebra](#)
- > [optimization](#)
- > [Survival](#)
- > [Time series](#)
- > [Uncategorized](#)



```

13 COEFFICIENTS:
14           Estimate Std. Error t value Pr(>|t|)
15 (Intercept)  2.646e+00  1.256e-01  21.056 < 2e-16 ***
16 crim        -8.432e-03  1.406e-03  -5.998 3.82e-09 ***
17 rm          1.428e-01  1.738e-02   8.219 1.77e-15 ***
18 tax         -2.562e-04  7.599e-05  -3.372 0.000804 ***
19 lstat       -2.954e-02  1.987e-03 -14.867 < 2e-16 ***
20 ---
21 Signif. codes:
22 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
23
24 Residual standard error: 0.2158 on 501 degrees of freedom
25 Multiple R-squared:  0.7236,    Adjusted R-squared:  0.72
26 F-statistic: 327.9 on 4 and 501 DF,  p-value: < 2.2e-16

```

## Residual Summary Statistics

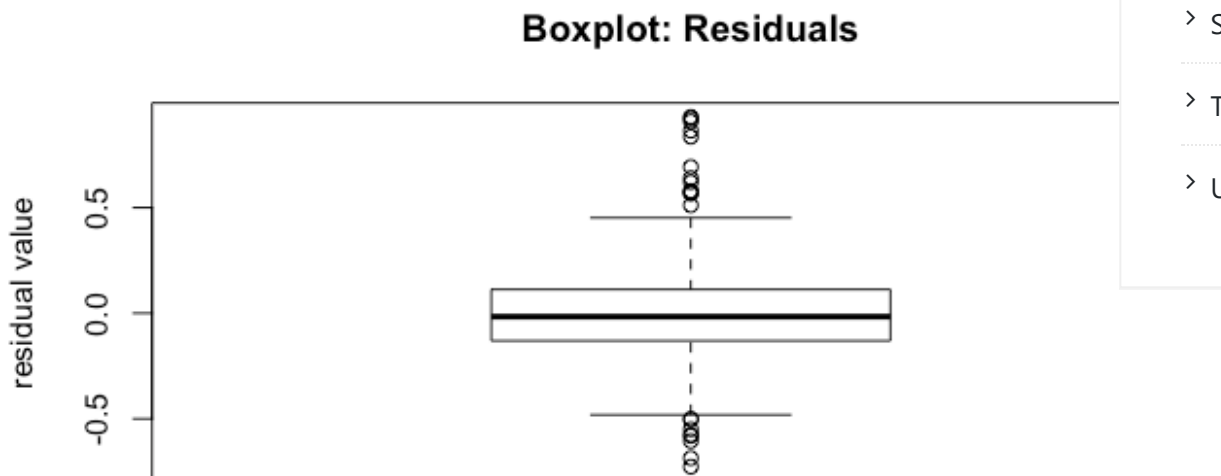
The first info printed by the linear regression summary after the formula is the residual statistics. One of the assumptions for hypothesis testing is that the errors follow a Gaussian distribution. As a consequence the residuals should as well. The residual summary statistics provide information about the symmetry of the residual distribution. The median should be close to the mean of the residuals is 0, and symmetric distributions have median=mean. Furthermore, the 1Q and 3Q should be close to each other in magnitude. They would be equal under a symmetric mean distribution. The max and min should also have similar magnitude. However, in practice, a large difference between the max and min may indicate an outlier rather than a symmetry violation.

We can investigate this further with a boxplot of the residuals.

```

1 boxplot(model[['residuals']],main='Boxplot: Residuals',ylab='residual value')

```



### Categories

- > general
- > hypothesis test
- > linear and generalized linear models
- > Longitudinal Data Analysis
- > Missing Data
- > Neural Networks+Deep Learning
- > Numerical Linear Algebra
- > optimization
- > Survival
- > Time series
- > Uncategorized

We see that the median is close to 0. Further, the 25 and 75 percentile look approximately the same distance from 0, and the non-outlier min and max also look about the same distance from 0. All of this is good as it suggests correct model specification.

## Coefficients

The second thing printed by the linear regression summary call is information about the coefficients. This includes their estimates, standard errors, t statistics, and p-values.

## Estimates

The intercept tells us that when all the features are at 0, the expected response is the intercept. Note that for an arguably better interpretation, you should consider [centering your features](#). This changes the interpretation. Now, when features are at their mean values, the expected response is the intercept. For the other features, the estimates give us the expected change in the response due to a unit change in the feature.

## Standard Error

The standard error is the standard error of our estimate, which allows us to construct confidence intervals for the estimate of that particular feature. If  $s.e.(\hat{\beta}_i)$  is the standard error of  $\hat{\beta}_i$  is the estimated coefficient for feature  $i$ , then a 95% confidence interval is given by  $\hat{\beta}_i \pm 1.96 \cdot s.e.(\hat{\beta}_i)$ . Note that this requires two things for this confidence interval to be valid:

- your model assumptions hold
- you have enough data/samples to invoke the central limit theorem, as you need the distribution to be approximately Gaussian.

That is, assuming all model assumptions are satisfied, we can say that with 95% confidence (this is not probability) the true parameter  $\beta_i$  lies in  $[\hat{\beta}_i - 1.96 \cdot s.e.(\hat{\beta}_i), \hat{\beta}_i + 1.96 \cdot s.e.(\hat{\beta}_i)]$ . On this, we can construct confidence intervals

```
1 confint(model)
2              2.5 %      97.5 %
3 (Intercept)  2.3987332457  2.8924423620
4 crim        -0.0111943622 -0.0056703707
5 rm          0.1086963289  0.1769912871
6 tax         -0.0004055169 -0.0001069386
7 lstat       -0.0334396331 -0.0256328293
```

Here we can see that the entire confidence interval for number of rooms has a large effect size relative to the other covariates.

### Categories

- > general
- > hypothesis test
- > linear and generalized linear models
- > Longitudinal Data Analysis
- > Missing Data
- > Neural Networks+Deep Learning
- > Numerical Linear Algebra
- > optimization
- > Survival
- > Time series
- > Uncategorized

## t-value

The [t-statistic](#) is

$$\frac{\hat{\beta}_i}{s.e.(\hat{\beta}_i)}$$

which tells us about how far our estimated parameter is from a hypothesized 0 value, standard deviation of the estimate. Assuming that  $\hat{\beta}_i$  is Gaussian, under the null hypothesis  $\beta_i = 0$ , this will be t distributed with  $n - p - 1$  degrees of freedom, where  $n$  is the observations and  $p$  is the number of parameters we need to estimate.

## Pr(>|t|)

This is the [p-value](#) for the individual coefficient. Under the t distribution with  $n - p - 1$  degrees of freedom, this tells us the probability of observing a value at least as extreme as our  $\hat{\beta}_i$ . If the probability is sufficiently low, we can reject the null hypothesis that this coefficient is zero. Note that when we care about looking at *all* of the coefficients, we are actually doing multiple hypothesis tests, and need to correct for that. In this case we are making five hypothesis tests, one for each feature and one for the coefficient. Instead of using the standard p-value of 0.05, we use the Bonferroni correction and divide by the number of hypothesis tests, and thus use a value threshold of 0.01.

## Assessing Fit and Overall Significance

The linear regression summary printout then gives the residual standard error, the  $R^2$  statistic and test. These tell us about how good a fit the model is and whether any of the coefficients are significant.

## Residual Standard Error

The residual standard error is given by  $\hat{\sigma} = \sqrt{\frac{\sum \hat{\epsilon}_i^2}{n-p}}$ . It gives the standard deviation of the residuals, and tells us about how large the prediction error is *in-sample* or *on the training data*. We'd like this to be significantly different from the variability in the marginal response, otherwise it's not clear that the model explains much.

## Multiple and Adjusted $R^2$

Intuitively  $R^2$  tells us what proportion of the variance is explained by our model, and is given by

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{\sum_i \hat{\epsilon}_i^2}{\sum_i y_i^2}$$

(1)

### Categories

- > general
- > hypothesis test
- > linear and generalized linear models
- > Longitudinal Data Analysis
- > Missing Data
- > Neural Networks+Deep Learning
- > Numerical Linear Algebra
- > optimization
- > Survival
- > Time series
- > Uncategorized



(2)

$$\sum_i (y_i - \bar{y})^2$$

both  $R^2$  and the residual standard standard deviation tells us about how well our model fits the data. The adjusted  $R^2$  deals with an increase in  $R^2$  spuriously due to adding features fitting noise in the data. It is given by

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

thus as the number of features  $p$  increases, the required  $R^2$  needed will increase as  $n$  increases to maintain the same adjusted  $R^2$ .

## F-Statistic and F-test

In addition to looking at whether individual features have a significant effect, we may want to test whether *at least one* feature has a significant effect. That is, we would like to test the null hypothesis

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$$

that all coefficients are 0 against the alternative hypothesis

$$H_1 : \exists i : 1 \leq i \leq p - 1 : \beta_i \neq 0$$

Under the null hypothesis the F statistic will be F distributed with  $(p - 1, n - p)$  degrees of freedom. The probability of our observed data under the null hypothesis is then the p-value. If we use the F-test alone without looking at the t-tests, then we do not need a Bonferroni correction while if we do look at the t-tests, we need one.

« Previous

## Related Posts

[Linear Regression: Comparing Models Between Two Groups with linearHypothesis](#)

September 21, 2019

[Why You Should Center Your Features in Linear Regression](#)

August 31, 2019

## Categories

- > general
- > hypothesis test
- > linear and generalized linear models
- > Longitudinal Data Analysis
- > Missing Data
- > Neural Networks+Deep Learning
- > Numerical Linear Algebra
- > optimization
- > Survival
- > Time series
- > Uncategorized



## Leave a Reply

Your email address will not be published. Required fields are marked<sup>\*</sup>

Comment<sup>\*</sup>

Name<sup>\*</sup>

Email<sup>\*</sup>

Website

Post Comment

This site uses Akismet to reduce spam. [Learn how your comment data is processed.](#)

### Categories

- › general
- › hypothesis test
- › linear and generalized linear models
- › Longitudinal Data Analysis
- › Missing Data
- › Neural Networks+Deep Learning
- › Numerical Linear Algebra
- › optimization
- › Survival
- › Time series
- › Uncategorized



## Categories

- › general
- › hypothesis test
- › linear and generalized linear models
- › Longitudinal Data Analysis
- › Missing Data
- › Neural Networks+Deep Learning
- › Numerical Linear Algebra
- › optimization
- › Survival
- › Time series
- › Uncategorized





## Categories

- › general
- › hypothesis test
- › linear and generalized linear models
- › Longitudinal Data Analysis
- › Missing Data
- › Neural Networks+Deep Learning
- › Numerical Linear Algebra
- › optimization
- › Survival
- › Time series
- › Uncategorized



Categories

- > general
- > hypothesis test
- > linear and generalized linear models
- > Longitudinal Data Analysis
- > Missing Data
- > Neural Networks+Deep Learning
- > Numerical Linear Algebra
- > optimization
- > Survival
- > Time series
- > Uncategorized

© All Right Reserved

