

Assessing the Accuracy of the Model: Takeaways



by Dataquest Labs, Inc. - All rights reserved © 2021

Syntax

EVALUATING BIVARIATE RELATIONSHIPS

- Visualize distribution of residuals with a histogram:

```
library(ggplot2)
residuals_df <- data.frame(residuals = lm_fit$residuals)
ggplot(data = residuals_df,
       aes(x = residuals)) +
  geom_histogram()
```

- View linear model summary:

```
summary(lm_fit)
```

- Manually estimating the t-statistic:

```
(lm_fit$coefficients[[2]] - 0) / coef(summary(lm_fit))[, 2][[2]]
```

- Extract the p-value from a bi-variate linear model summary:

```
p_value <- coef(summary(lm_fit))[, 4][[2]]
```

- Manually estimate the residual sum of squares (RSS):

```
df <- df %>%
  mutate(residuals = resid(lm_fit)) %>%
  mutate(resid_squared = residuals^2)
RSS <- df %>%
  summarise(RSS = sum(resid_squared)) %>%
  pull()
```

- Extract RSS from model output:

```
RSS <- deviance(lm_fit)
```

- Manually estimate the residual standard error (RSE):

```
RSE <- sqrt(RSS / (nrow(df) - 2))
```

- Extract RSE from model output:

```
RSE <- sigma(lm_fit)
```

- Manually estimate total sum of squares (TSS):

```
TSS <- sum((df$response - mean(df$response))^2)
```

- Manually estimate r-squared:

```
r_squared <- 1 - RSS/TSS
```

- Extract r-squared value from linear model object:

```
r_squared <- summary(lm_fit)$r.squared
```

- Extract adjusted r-squared from linear model object:

```
adj_r_squared <- summary(lm_fit)$adj.r.squared
```

Equations

- Mathematical equation for hypothesis test:
 - Null hypothesis = $H_0 : \beta_1 = 0$
 - Alternative hypothesis = $H_a : \beta_1 \neq 0$
- Mathematical equation for the t-statistic:
 - $t = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)}$
- A 95% confidence interval for the intercept is *approximately* equal to:
 - $[\hat{\beta}_0 - 2 * SE(\hat{\beta}_0), \hat{\beta}_0 + 2 * SE(\hat{\beta}_0)]$
- Alternatively:
 - $\hat{\beta}_0 \pm 2 * SE(\hat{\beta}_0)$
- And the 95% confidence interval for the slope *approximately* equals:
 - $\hat{\beta}_1 \pm 2 * SE(\hat{\beta}_1)$
- Residual sum of squares (RSS):
 - $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- Residual standard error RSE is:
 - $RSE = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$
 - $RSE = \sqrt{\frac{1}{n-2} RSS}$
- Total sum of squares (TSS):
 - $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$
- R-squared:
 - $R^2 = \frac{\text{Variance explained by the model}}{\text{Total variance}} = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$

Concepts

- **Null hypothesis H_0 :** there is no relationship between predictor variable and the response variable.
- **Alternative hypothesis H_a :** there is a relationship between predictor variable and the response variable.
- **t-statistic:** is the number of standard deviations that $\hat{\beta}_1$ is from 0.
- **p-value:** is the probability of observing any value equal-to or larger than t if the null hypothesis is true. A smaller p-value is better.
- **Confidence interval:** a confidence interval of 95% means that there is a 95% probability that the true unknown value of the coefficient will fall within the specified range.

- **Residual standard error (RSE):** represents the average amount that our response variable measurements deviate from the true regression line. The RSE is an estimate of the standard deviation of ϵ .
- **R-squared (R^2):** a measure of the proportion of the variability in the response variable that can be explained by the predictor variable. The R^2 value falls between 0 and 1.

Resources

- [Dataquest blog post on linear regression for predictive modeling in R.](#)
- [Dataquest blog post on linear regression error metrics.](#)
- [Wikipedia entry on the null hypothesis.](#)
- [Wikipedia entry on the t-statistic.](#)
- [Wikipedia entry on the t-distribution.](#)
- [Wikipedia entry on the coefficient of determination \(r-squared\).](#)
- [Wikipedia entry on the total sum of squares.](#)
- [An Introduction to Statistical Learning with Applications in R by James et al.](#)