

# Performing linear regression analysis using `lm()` (R).

## Aim

To perform **linear regression analysis using the `lm()` function in R** in order to study the relationship between **number of playlists** and **song streams**, and to predict the number of streams for a given playlist count.

---

## Dataset

- **File Name:** Song.csv
  - **Source:** Local CSV file
  - **Variables Used:**
    - **Independent Variable (X):** Number of Playlists
    - **Dependent Variable (Y):** Streams (in millions)
- 

## Theory

Linear regression is a statistical method used to model the linear relationship between a dependent variable and an independent variable. In R, linear regression is performed using the `lm()` function.

## Mathematical Model:

**Name :- Priya Gupta**  
**Roll No :- So81**

# Sheth L.U.J. & Sir M.V. College

Mathematical Model:

$$Y = \beta_0 + \beta_1 X$$

Where:

- $Y$  = Streams (in millions)
  - $X$  = Number of Playlists
  - $\beta_0$  = Intercept
  - $\beta_1$  = Slope
- 

## Procedure

1. Loaded the dataset using `read.csv()` with proper file encoding.
  2. Converted playlist and stream columns into numeric format.
  3. Removed missing values using `na.omit()`.
  4. Converted stream values into millions for better interpretation.
  5. Applied linear regression using the `lm()` function.
  6. Obtained regression coefficients, R-squared value, and statistical significance from `summary()`.
  7. Predicted streams for a new playlist value.
  8. Visualized the regression using scatter plot and regression line.
- 

**Name :- Priya Gupta**

**Roll No :- So81**

# Sheth L.U.J. & Sir M.V. College

## Results

### Regression Coefficients

From the `summary(model)` output:

Parameter	Estimate
Intercept ( $\beta_0$ )	167.6
Slope ( $\beta_1$ )	0.0815

### Regression Equation:

$$\text{Streams (in millions)} = 167.6 + 0.0815 \times (\text{Number of Playlists})$$

```
Call:
lm(formula = streams_million ~ playlist, data = song)

Residuals:
    Min       1Q   Median       3Q      Max
-1337.0  -159.6   -74.2   137.9  1278.0

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 1.676e+02  4.650e+01   3.605 0.000495 ***
playlist     8.149e-02  4.383e-03  18.592 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 371.5 on 97 degrees of freedom
Multiple R-squared:  0.7809,    Adjusted R-squared:  0.7786
F-statistic: 345.7 on 1 and 97 DF,  p-value: < 2.2e-16
```

## Model Statistics

- **Multiple R-squared:** 0.7809
- **Adjusted R-squared:** 0.7786

**Name :- Priya Gupta**

**Roll No :- So81**

# Sheth L.U.J. & Sir M.V. College

- **F-statistic:** 345.7
- **p-value:**  $< 2.2e-16$

## Interpretation:

The model explains approximately **78.09%** of the variation in streams, indicating a **strong linear relationship** between playlists and streams. The very small p-value confirms that the regression model is statistically significant.

---

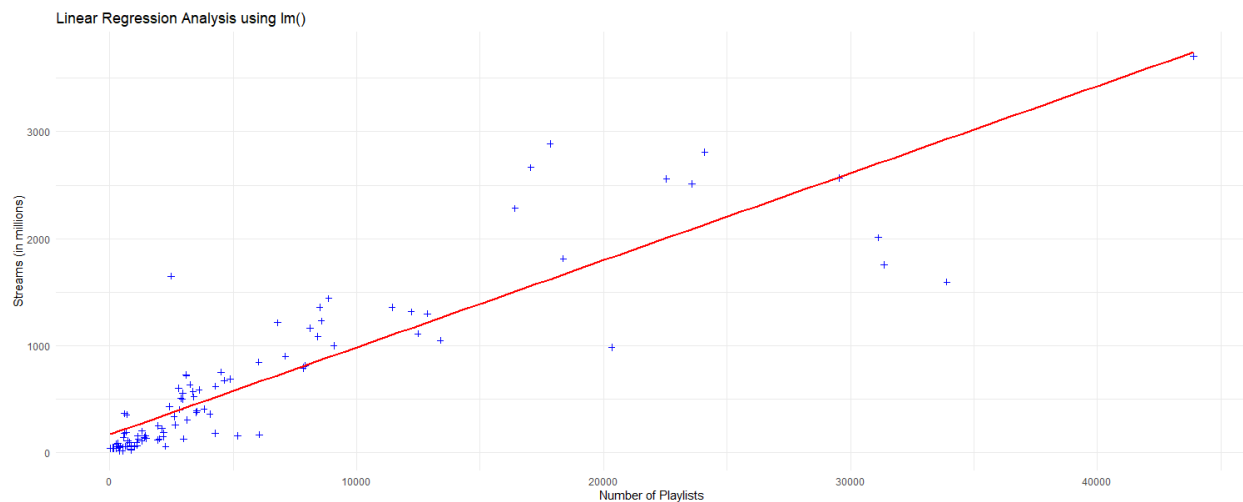
## Prediction Result

For **2000** playlists:

```
+ predicted_value, "\n")
Predicted streams (in millions) for 2000 playlists: 330.6156
>
```

---

## Graph



## Graph Interpretation:

**Name :- Priya Gupta**

**Roll No :- So81**

## Sheth L.U.J. & Sir M.V. College

- Blue points represent actual observations.
  - The red line represents the best-fit regression line generated using `lm()`.
  - The upward trend shows a positive relationship between playlists and streams.
- 

### Conclusion

Linear regression analysis was successfully performed using the `lm()` function in R. The results show that the number of playlists has a significant positive impact on the number of song streams. The high R-squared value indicates that the model fits the data well and can be effectively used for prediction.

---

### Screenshots

**Name :- Priya Gupta**  
**Roll No :- So81**

# Sheth L.U.J. & Sir M.V. College

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)

Source
Console Background Jobs
R - R 4.5.2 - ~/
> # =====
> # Linear Regression using lm()
> # =====
>
> # Load required library
> library(ggplot2)
>
> # Read CSV file
> song <- read.csv("c:/users/itlab/downloads/song.csv",
+ fileEncoding = "ISO-8859-1")
>
> # Convert required columns to numeric
> # Column 7 = Number of Playlists
> # Column 9 = Streams
> song$playlist <- as.numeric(song[, 7])
> song$streams <- as.numeric(song[, 9])
>
> # Remove missing values
> song <- na.omit(song)
>
> # Convert streams into millions
> song$streams_million <- song$streams / 1e6
>
> # =====
> # Linear Regression Model using lm()
> # =====
>
> model <- lm(streams_million ~ playlist, data = song)
>
> # Display model summary
> summary(model)

Call:
lm(formula = streams_million ~ playlist, data = song)

Residuals:
    Min       1Q   Median       3Q      Max
-1337.0  -159.6   -74.2   137.9  1278.0

Coefficients:
(Intercept) 1.676e+02  4.650e+01  3.605 0.000495 ***
playlist     8.149e-02  4.383e-03 18.592 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 371.5 on 97 degrees of freedom
Multiple R-squared:  0.7809, Adjusted R-squared:  0.7786
F-statistic: 345.7 on 1 and 97 DF, p-value: < 2.2e-16

>
> # =====
> # Prediction
> # =====
>
> newdata <- data.frame(playlist = 2000)
> predicted_value <- predict(model, newdata)
>
> cat("Predicted streams (in millions) for 2000 playlists:",
+ predicted_value, "\n")
Predicted streams (in millions) for 2000 playlists: 330.6156
>
> # =====
> # Plot
> # =====
>
> ggplot(song, aes(x = playlist, y = streams_million)) +
+   geom_point(color = "blue", shape = 3) +
+   geom_smooth(method = "lm", se = FALSE, color = "red") +
+   labs(
+     title = "Linear Regression Analysis using lm()",
+     x = "Number of Playlists",
+     y = "Streams (in millions)"
+   ) +
+   theme_minimal()
> geom_smooth() using formula = 'y ~ x'
>
```

**Name :- Priya Gupta**  
**Roll No :- So81**

# **Sheth L.U.J. & Sir M.V. College**

**Name :- Priya Gupta**

**Roll No :- So81**