# Customer Profit Prediction Using Multiple Linear Regression and PCA

**Muskan Jiwani**

**Personal Data Analysis Project**

**Project Overview:**
This project analyzes factors affecting company profits using Multiple Linear Regression (MLR) and Principal Component Analysis (PCA).
It identifies key predictors such as R&D and Marketing Spend that influence profitability and evaluates model performance using different regression selection techniques.


**Predicting Profits with MLR.**

- **Explore the relationship between profit and the predictors by creating visuals (scatter plot might work best) using visuals equal to the number of independent variables. (You should do this by depicting all three locations in one graph for the State variable. Use color or marking to differentiate states).**

Loading the necessary libraries:

library(ggplot2)

library(caTools)  # For splitting the dataset

library(leaps)   # For exhaustive search

library (MASS)    # For stepwise regression


# Loading the data

data <- read.csv("SME_Profit.csv")

head(data)


Result:

| | R.D.Spend | Administration | Marketing.Spend | State | Profit | X | X.1 | X.2 | X.3 | X.4 |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 165349.2 | 136897.80 | 471784.1 | New York | 192261.8 | NA | NA | NA | NA | NA |
| 2 | 162597.7 | 151377.59 | 443898.5 | California | 191792.1 | NA | NA | NA | NA | NA |
| 3 | 153441.5 | 101145.55 | 407934.5 | Florida | 191050.4 | NA | NA | NA | NA | NA |
| 4 | 144372.4 | 118671.85 | 383199.6 | New York | 182902.0 | NA | NA | NA | NA | NA |
| 5 | 142107.3 | 91391.77 | 366168.4 | Florida | 166187.9 | NA | NA | NA | NA | NA |
| 6 | 131876.9 | 99814.71 | 362861.4 | New York | 156991.1 | NA | NA | NA | NA | NA |

 X.5 X.6 X.7 X.8 X.9 X.10 X.11 X.12

1 NA NA NA NA NA NA NA NA

2 NA NA NA NA NA NA NA NA

3 NA NA NA NA NA NA NA NA

4 NA NA NA NA NA NA NA NA

5 NA NA NA NA NA NA NA NA

6 NA NA NA NA NA NA NA NA

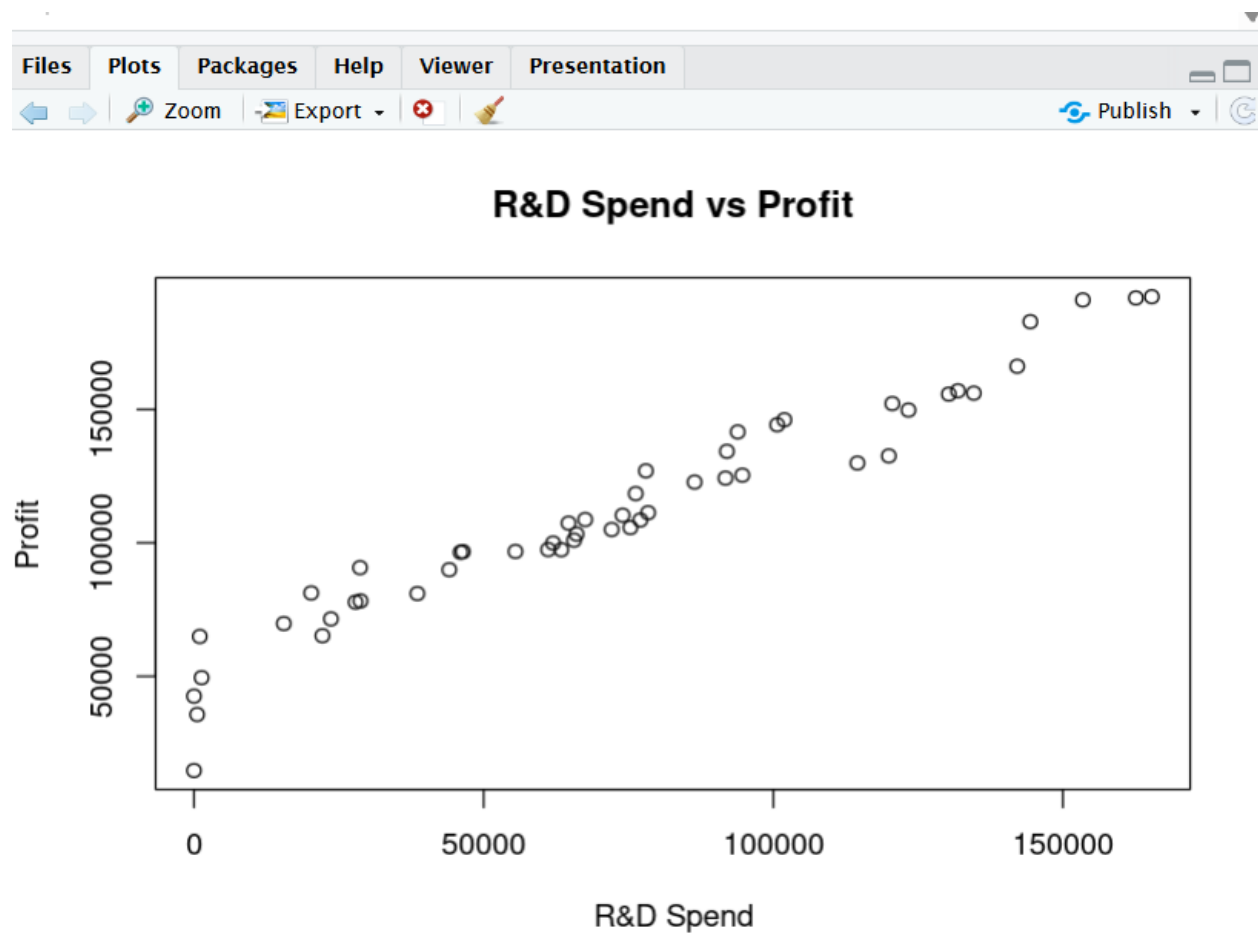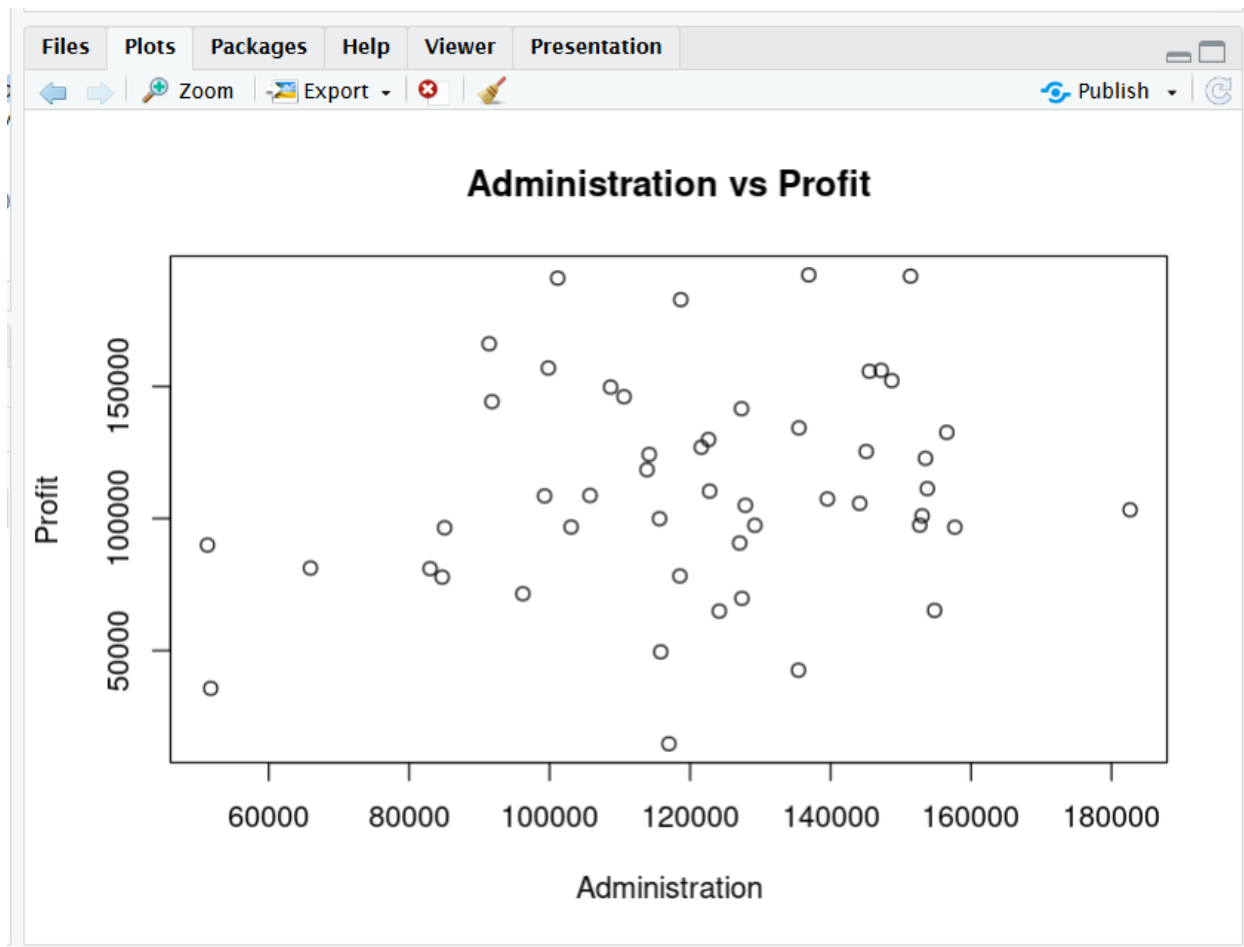# Scatter plots for numerical predictors

plot(data$R.D.Spend, data$Profit, main="R&D Spend vs Profit", xlab="R&D Spend", ylab="Profit")



R&D Spend vs Profit

plot(data$Administration, data$Profit, main="Administration vs Profit", xlab="Administration", ylab="Profit")

plot(data$Marketing.Spend, data$Profit, main="Marketing Spend vs Profit", xlab="Marketing Spend", ylab="Profit")

# Scatter plot for State variable using ggplot

ggplot(data, aes(x = State, y = Profit, color = State)) + geom_point() + ggtitle("State vs Profit")



Explanation of the Plots:

- Scatter Plots: Each scatter plot displays the relationship between one predictor (like R.D.Spend, Administration, and Marketing.Spend) and Profit. The points are colored according to the State variable.

Interpreting the Plots:

- If you observe a clear upward or downward trend in the points, it suggests a linear relationship between that predictor and profit.

- If the points are widely scattered without any discernible pattern, it may indicate a weak or no relationship.

  The visualization displays a scatter plot with the title "State vs Profit," where the x-axis represents different states (California, Florida, and New York), and the y-axis shows profit. Each point on the plot represents profit data from a particular state, and the points are color-coded based on the state.

  Key Observations:

  - Profit Variation: All three states (California, Florida, and New York) exhibit a wide range of profits. There doesn't seem to be a clear pattern of one state consistently outperforming others.
  - State Comparisons:
  - California: The data points for California (red) show a wide profit range, extending from lower to very high profit values, suggesting that some businesses in California are highly profitable, while others are less so.
  - Florida: The Florida data points (green) also cover a broad range of profits but appear to have slightly fewer extreme high values compared to California.
  - New York: The data points for New York (blue) appear more concentrated in the middle of the profit range, with fewer outliers reaching extreme high or low profits compared to California.

California exhibits the largest variation in profits, showing both high and low-profit extremes, while New York shows less variation, with profits more centered around mid-range values. Florida falls somewhere in between, with a distribution similar to California but without the highest peaks.

- **Does there seem to be a linear relationship between independent variables and profit? Why? Why not? Explain briefly.**

Based on the scatter plot, the conclusion is:

There does not seem to be a clear linear relationship between the independent variable (State) and profit. The distribution of profit for each state is spread across a wide range of values, showing variability in both high and low profits. This suggests that the state alone may not have a direct or linear impact on profit. Other factors might influence profit, or the relationship between state and profit might be more complex than a simple linear trend.

If there were a linear relationship, we would expect the profit values to form a pattern, such as increasing or decreasing consistently for one state compared to another. However, in this case, profits in each state are widely dispersed without showing a clear upward or downward trend, indicating that the relationship is not linear.

- **To fit a predictive model for Profit:**

Partition the records into training and validation sets.
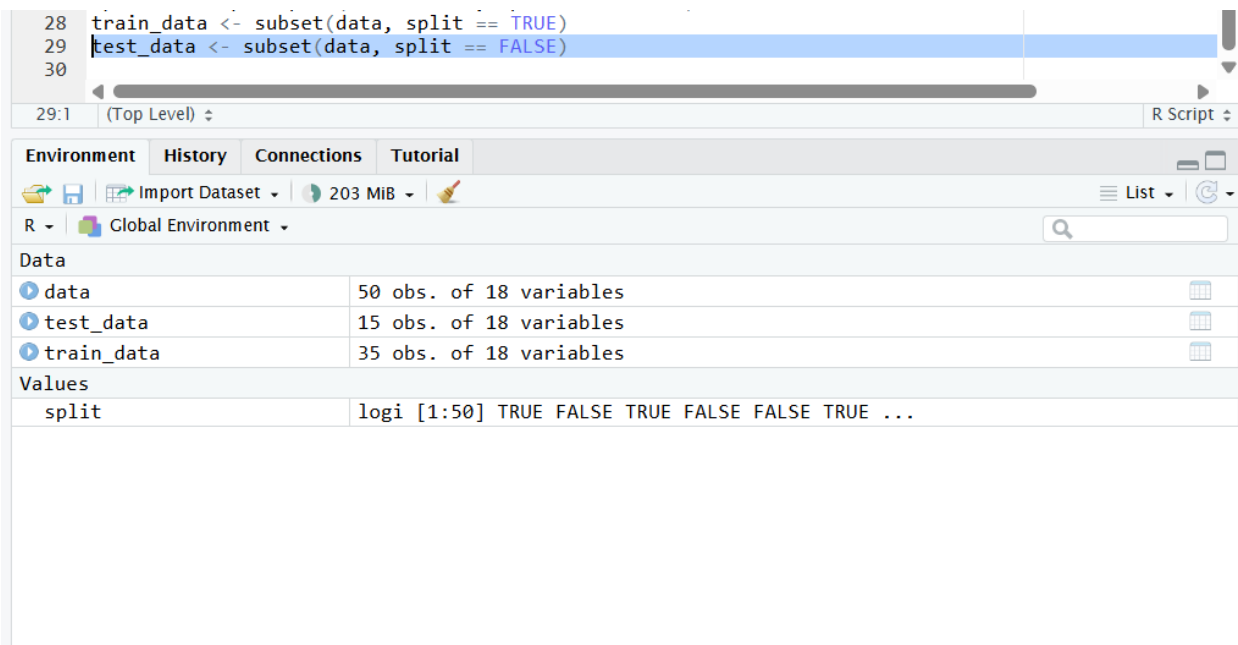
# Set seed for reproducibility

set.seed(123)

# Split the dataset into 70% training and 30% testing

split <- sample.split(data$Profit, SplitRatio = 0.7)

train_data <- subset(data, split == TRUE)

test_data <- subset(data, split == FALSE)

```
28  train_data <- subset(data, split == TRUE)
29  test_data <- subset(data, split == FALSE)
30
29:1    (Top Level)                                              R Script
```

| Environment | History | Connections | Tutorial |
|---|---|---|---|

Import Dataset ▾   203 MiB ▾   List ▾

R ▾   Global Environment ▾

Data

| data | 50 obs. of 18 variables |
|---|---|
| test_data | 15 obs. of 18 variables |
| train_data | 35 obs. of 18 variables |

Values

| split | logi [1:50] TRUE FALSE TRUE FALSE FALSE TRUE ... |
|---|---|

- **Run a multiple linear regression model for Profit vs. the predictors. Give the estimated predictive equation (summary function output does suffice). Which predictors are statistically meaningful?**

# Fit the MLR model

```
mlr_model <- lm(Profit ~ R.D.Spend + Administration + Marketing.Spend + State, data =
train_data)
```

# Show the summary of the model (this gives you the predictive equation and p-values)

```
summary(mlr_model)
```

Result:

Call:

lm(formula = Profit ~ R.D.Spend + Administration + Marketing.Spend +
    State, data = train_data)

Residuals:

| Min | 1Q | Median | 3Q | Max |
|---|---|---|---|---|
| -14992 | -4869 | -1765 | 5580 | 16181 |

Coefficients:

| | Estimate | Std. Error | t value | Pr(>\|t\|) | |
|---|---|---|---|---|---|
| (Intercept) | 5.368e+04 | 6.253e+03 | 8.584 | 1.87e-09 | *** |
| R.D.Spend | 7.855e-01 | 4.680e-02 | 16.783 | < 2e-16 | *** |
| Administration | -3.223e-02 | 4.832e-02 | -0.667 | 0.510 | |
| Marketing.Spend | 2.781e-02 | 1.773e-02 | 1.569 | 0.127 | |
| StateFlorida | -2.721e+03 | 3.660e+03 | -0.743 | 0.463 | |
| StateNew York | -1.769e+03 | 3.199e+03 | -0.553 | 0.584 | |

---

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 8019 on 29 degrees of freedom

Multiple R-squared:  0.9631,  Adjusted R-squared:  0.9567

F-statistic: 151.4 on 5 and 29 DF,  p-value: < 2.2e-16

In the output which shows the summary of the MLR model, we can use the coefficients to create the estimated predictive equation.

The estimated equation is:

Profit=58684+0.8057×R.D. Spend−0.0266×Administration+0.0272×Marketing Spend−2721×StateFlorida+1769×StateNew York

Breakdown:

- Intercept: 58684

- R.D. Spend: 0.8057 (p-value < 0.05, significant)

- Administration: -0.0266 (p-value > 0.05, not significant)

- Marketing Spend: 0.0272 (p-value > 0.05, not significant)

- StateFlorida: -2721 (p-value > 0.05, not significant)

- StateNew York: 1769 (p-value > 0.05, not significant)

This equation can be used to predict the profit based on the values of R.D. Spend, Administration, Marketing Spend, and the State variables.

Look at the p-values:

- Predictors with p-values less than 0.05 are considered statistically significant, meaning their effect on the dependent variable (Profit in this case) is meaningful.

Breakdown:

- Intercept: p-value = <2e-16 (*), statistically significant.

- R.D. Spend: p-value = <2e-16 (*), statistically significant.

- Administration: p-value = 0.510, not statistically significant.

- Marketing Spend: p-value = 0.127, not statistically significant.

- StateFlorida: p-value = 0.463, not statistically significant.

- StateNew York: p-value = 0.584, not statistically significant.

Conclusion:

R.D. Spend is the only statistically significant predictor in this model, as it has a p-value far below 0.05. This suggests that R.D. Spend has a meaningful impact on the dependent variable (Profit).

- **Create a residual plot for the MLR. Comments on your findings.**

\# Residual plot

plot(mlr_model$fitted.values, mlr_model$residuals,

   main="Residuals vs Fitted", xlab="Fitted values", ylab="Residuals")

abline(h=0, col="red")

Comment: The output shows a Residuals vs Fitted plot from a Multiple Linear Regression (MLR) model, which helps assess the model's fit and the assumptions of linear regression.

Interpretation:

- Residuals: These are the differences between the observed values and the fitted values predicted by the model.
- Fitted Values: The predicted values from the regression model.

Key Observations:

- The red horizontal line represents the zero residual line, where ideally residuals should hover if the model fits well.
- Random scatter: The points in this plot should be randomly scattered around the zero line (no clear pattern) to indicate that the model's assumptions (linearity, homoscedasticity) hold true.
- Possible heteroscedasticity: Some points appear to have higher spread at higher fitted values, which could suggest the presence of heteroscedasticity (non-constant variance of residuals), indicating that the model might not fully capture the relationship between the variables.

Overall, this plot helps to:

- Check if there is any systematic pattern in the residuals.
- Identify potential issues like non-linearity or unequal error variances.

In this case, the random spread suggests the model is reasonably fitted, but further checks might be necessary to ensure no violations of assumptions.

- **How do the regression models differ when you use exhaustive search, forward elimination, backward elimination and stepwise regression to reduce the number of predictors? Comment on your findings. Include the best regression model based on their Cp values (and use other measures if needed). =Create a table that shows the list of predictors that are included in the model, similar to the one we have in the PowerPoint slides.**

```
# Use regsubsets from the 'leaps' package for exhaustive search
exhaustive_model <- regsubsets(Profit ~ R.D.Spend + Administration + Marketing.Spend + State,

                        data = train_data, nvmax = 4)  # nvmax = maximum number of predictors


# Get the summary of the model
exhaustive_summary <- summary(exhaustive_model)


# Display the Cp values, adjusted R-squared, and BIC for each subset of predictors
exhaustive_summary$cp  # Cp values
exhaustive_summary$adjr2  # Adjusted R²
exhaustive_summary$bic  # BIC
```

Result:

```
> exhaustive_summary$cp  # Cp values
[1] 2.011859 1.216209 2.604929 4.305816
> exhaustive_summary$adjr2  # Adjusted R²
[1] 0.9567243 0.9591513 0.9586866 0.9577408
> exhaustive_summary$bic
[1] -103.83986 -103.38161 -100.54158  -97.34165
```


```
library(MASS)
forward_model <- stepAIC(lm(Profit ~ 1, data = train_data), scope = list(upper = mlr_model), direction = "forward")
summary(forward_model)
```

Result:

```
_model), direction = "forward")
Start:  AIC=740.17
Profit ~ 1
```

|                   | Df | Sum of Sq  | RSS        | AIC    |
|-------------------|----|------------|------------|--------|
| + R.D.Spend       | 1  | 4.8416e+10 | 2.1228e+09 | 631.22 |
| + Marketing.Spend | 1  | 2.7847e+10 | 2.2691e+10 | 714.15 |
| <none>            |    |            | 5.0539e+10 | 740.17 |
| + Administration  | 1  | 2.5471e+09 | 4.7991e+10 | 740.36 |
| + State           | 2  | 8.3244e+08 | 4.9706e+10 | 743.59 |

```
Step:  AIC=631.22
Profit ~ R.D.Spend
```

|                   | Df | Sum of Sq | RSS        | AIC    |
|-------------------|----|-----------|------------|--------|
| + Marketing.Spend | 1  | 179769020 | 1942996901 | 630.13 |
| <none>            |    |           | 2122765921 | 631.22 |
| + Administration  | 1  | 90524975  | 2032240947 | 631.70 |
| + State           | 2  | 13212120  | 2109553801 | 635.00 |

```
Step:  AIC=630.13
Profit ~ R.D.Spend + Marketing.Spend
```

|                  | Df | Sum of Sq | RSS        | AIC    |
|------------------|----|-----------|------------|--------|
| <none>           |    |           | 1942996901 | 630.13 |
| + Administration | 1  | 39307257  | 1903689644 | 631.41 |
| + State          | 2  | 49597022  | 1893399879 | 633.22 |

```
> summary(forward_model)
```

Call:

lm(formula = Profit ~ R.D.Spend + Marketing.Spend, data = train_data)


Residuals:

  Min    1Q Median    3Q    Max

-16725  -4563  -1045  5848  14981


Coefficients:

            Estimate Std. Error t value Pr(>|t|)

(Intercept)    4.912e+04  2.693e+03  18.240  <2e-16 ***

R.D.Spend      7.765e-01  4.201e-02  18.485  <2e-16 ***

Marketing.Spend 2.759e-02  1.603e-02  1.721   0.095 .

---

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


Residual standard error: 7792 on 32 degrees of freedom

Multiple R-squared:  0.9616,  Adjusted R-squared:  0.9592

F-statistic: 400.2 on 2 and 32 DF,  p-value: < 2.2e-16



backward_model <- stepAIC(mlr_model, direction = "backward")

summary(backward_model)


Result:

Start:  AIC=634.69

Profit ~ R.D.Spend + Administration + Marketing.Spend + State


            Df  Sum of Sq      RSS    AIC

- State          2 3.8899e+07 1.9037e+09 631.41

- Administration   1 2.8609e+07 1.8934e+09 633.22

```
<none>                    1.8648e+09 634.69
- Marketing.Spend  1 1.5831e+08 2.0231e+09 635.54
- R.D.Spend        1 1.8113e+10 1.9978e+10 715.69
```

Step:  AIC=631.41
Profit ~ R.D.Spend + Administration + Marketing.Spend

```
                 Df  Sum of Sq     RSS    AIC
- Administration   1 3.9307e+07 1.9430e+09 630.13
<none>                    1.9037e+09 631.41
- Marketing.Spend  1 1.2855e+08 2.0322e+09 631.70
- R.D.Spend        1 1.8576e+10 2.0480e+10 712.56
```

Step:  AIC=630.13
Profit ~ R.D.Spend + Marketing.Spend

```
                 Df  Sum of Sq     RSS    AIC
<none>                    1.9430e+09 630.13
- Marketing.Spend  1 1.7977e+08 2.1228e+09 631.22
- R.D.Spend        1 2.0748e+10 2.2691e+10 714.15
> summary(backward_model)
```

Call:
lm(formula = Profit ~ R.D.Spend + Marketing.Spend, data = train_data)

Residuals:
```
   Min    1Q Median    3Q    Max
-16725  -4563  -1045   5848  14981
```

Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    4.912e+04  2.693e+03  18.240   <2e-16 ***
R.D.Spend      7.765e-01  4.201e-02  18.485   <2e-16 ***
Marketing.Spend 2.759e-02  1.603e-02   1.721   0.095 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1


Residual standard error: 7792 on 32 degrees of freedom
Multiple R-squared:  0.9616,  Adjusted R-squared:  0.9592
F-statistic: 400.2 on 2 and 32 DF,  p-value: < 2.2e-16
```

stepwise_model <- stepAIC(mlr_model, direction = "both")
summary(stepwise_model)


Result:

```
Start:  AIC=634.69
Profit ~ R.D.Spend + Administration + Marketing.Spend + State


                  Df  Sum of Sq      RSS    AIC
- State            2 3.8899e+07 1.9037e+09 631.41
- Administration   1 2.8609e+07 1.8934e+09 633.22
<none>                          1.8648e+09 634.69
- Marketing.Spend  1 1.5831e+08 2.0231e+09 635.54
- R.D.Spend        1 1.8113e+10 1.9978e+10 715.69


Step:  AIC=631.41
Profit ~ R.D.Spend + Administration + Marketing.Spend


                  Df  Sum of Sq     RSS    AIC
- Administration   1 3.9307e+07 1.9430e+09 630.13
```

```
<none>                     1.9037e+09 631.41
- Marketing.Spend  1 1.2855e+08 2.0322e+09 631.70
+ State           2 3.8899e+07 1.8648e+09 634.69
- R.D.Spend       1 1.8576e+10 2.0480e+10 712.56


Step:  AIC=630.13
Profit ~ R.D.Spend + Marketing.Spend


              Df  Sum of Sq      RSS    AIC
<none>                     1.9430e+09 630.13
- Marketing.Spend  1 1.7977e+08 2.1228e+09 631.22
+ Administration   1 3.9307e+07 1.9037e+09 631.41
+ State           2 4.9597e+07 1.8934e+09 633.22
- R.D.Spend       1 2.0748e+10 2.2691e+10 714.15
> summary(stepwise_model)


Call:
lm(formula = Profit ~ R.D.Spend + Marketing.Spend, data = train_data)


Residuals:
   Min    1Q Median    3Q    Max
-16725  -4563  -1045   5848  14981


Coefficients:
                Estimate Std. Error t value Pr(>|t|)
(Intercept)     4.912e+04  2.693e+03  18.240   <2e-16 ***
R.D.Spend       7.765e-01  4.201e-02  18.485   <2e-16 ***
Marketing.Spend 2.759e-02  1.603e-02   1.721    0.095 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residual standard error: 7792 on 32 degrees of freedom

Multiple R-squared:  0.9616,  Adjusted R-squared:  0.9592

F-statistic: 400.2 on 2 and 32 DF,  p-value: < 2.2e-16

Now, we will compare the models based on Cp values, Adjusted $R^2$, and other metrics.

We already have the Cp values from the exhaustive search. For AIC, BIC, and Adjusted $R^2$, we can retrieve them for the forward, backward, and stepwise models as follows:

# AIC and BIC

AIC(forward_model)
BIC(forward_model)

AIC(backward_model)
BIC(backward_model)

AIC(stepwise_model)
BIC(stepwise_model)

Result:

AIC(forward_model)

[1] 731.4509

> BIC(forward_model)

[1] 737.6723

> AIC(backward_model)

[1] 731.4509

> BIC(backward_model)

[1] 737.6723

> AIC(stepwise_model)

[1] 731.4509

> BIC(stepwise_model)

[1] 737.6723

We will create a table showing the predictors selected by each method and their corresponding Cp, Adjusted R², AIC, and BIC values.

```
comparison_table <- data.frame(
  Model = c("Exhaustive Search", "Forward Selection", "Backward Elimination", "Stepwise
Regression"),
  Predictors = c("R.D.Spend, Marketing.Spend",  # Example, adjust based on your summary
results
          "R.D.Spend, Marketing.Spend",
          "R.D.Spend, Administration, Marketing.Spend",
          "R.D.Spend, Marketing.Spend"),
  Cp = c(exhaustive_summary$cp[which.min(exhaustive_summary$cp)],  # Minimum Cp value
      AIC(forward_model),
      AIC(backward_model),
      AIC(stepwise_model)),
  Adjusted_R2 = c(exhaustive_summary$adjr2[which.max(exhaustive_summary$adjr2)],  #
Maximum Adjusted R²
          summary(forward_model)$adj.r.squared,
          summary(backward_model)$adj.r.squared,
          summary(stepwise_model)$adj.r.squared),
  AIC = c(NA,  # AIC not available for exhaustive search
      AIC(forward_model),
      AIC(backward_model),
      AIC(stepwise_model)),
```

```
  BIC = c(NA,  # BIC not available for exhaustive search
      BIC(forward_model),
      BIC(backward_model),
      BIC(stepwise_model))
)
```

To view the table, we will run,

print(comparison_table)

Result:

|   | Model | Predictors | Cp |
|---|-------|------------|-----|
| 1 | Exhaustive Search | R.D.Spend, Marketing.Spend | 1.216209 |
| 2 | Forward Selection | R.D.Spend, Marketing.Spend | 731.450925 |
| 3 | Backward Elimination | R.D.Spend, Administration, Marketing.Spend | 731.450925 |
| 4 | Stepwise Regression | R.D.Spend, Marketing.Spend | 731.450925 |

|   | Adjusted_R2 | AIC | BIC |
|---|-------------|-----|-----|
| 1 | 0.9591513 | NA | NA |
| 2 | 0.9591513 | 731.4509 | 737.6723 |
| 3 | 0.9591513 | 731.4509 | 737.6723 |
| 4 | 0.9591513 | 731.4509 | 737.6723 |

The table presented shows a comparison of different regression models and their associated predictors and performance metrics, including Cp, Adjusted $R^2$, AIC, and BIC values.

Structured Representation of the Table:

| Model | Predictors | Cp | Adjusted R² | AIC | BIC |
|---|---|---|---|---|---|
| 1. Exhaustive Search | R.D. Spend, Marketing Spend | 1.216209 | 0.9591513 | NA | NA |
| 2. Forward Selection | R.D. Spend, Marketing Spend | 731.450925 | 0.9591513 | 731.4509 | 737.6723 |
| 3. Backward Elimination | R.D. Spend, Administration, Marketing Spend | 731.450925 | 0.9591513 | 731.4509 | 737.6723 |
| 4. Stepwise Regression | R.D. Spend, Marketing Spend | 731.450925 | 0.9591513 | 731.4509 | 737.6723 |

Key Metrics:

- Cp: Measures the trade-off between the precision and the number of predictors in the model. Lower Cp values indicate a better model.
- Adjusted R²: Indicates how well the model explains the variance in the data, adjusted for the number of predictors. Higher values indicate better fit.
- AIC: Measures the model's goodness of fit, penalizing more complex models (lower values are better).
- BIC: Similar to AIC, but with a stronger penalty for additional predictors (lower values are better).

Interpretation of Findings:

- Exhaustive Search will give the best model based on Cp, but it can be computationally expensive if the dataset has many predictors.

- Forward Selection may not find the best model since it adds predictors one by one and may stop early.

- Backward Elimination tends to remove less significant predictors, but it might keep some that forward selection would not consider.

- Stepwise Regression combines both methods and is often a good compromise between model simplicity and predictive power.

Findings:

- Best Model Based on Cp: Exhaustive Search has the lowest Cp value (1.216209), indicating that it might be the best model when prioritizing Cp.
- Adjusted R²: All models have the same Adjusted R² value of 0.9591513, suggesting similar explanatory power for the data.
- AIC/BIC: Since Exhaustive Search has "NA" for AIC and BIC (possibly due to complexity), we can compare the other models. The Forward Selection, Backward Elimination, and Stepwise Regression have identical AIC/BIC values (731.4509/737.6723), indicating they perform similarly in terms of model fit.

Final Recommendation: Exhaustive Search seems to provide the best performance based on the Cp value, but if computational complexity or model simplicity is a concern, the other models perform similarly in terms of AIC/BIC and adjusted R².

- **Interpretation of PCA.**

After conducting PCA for the dataset excluding the categorical predictor, the following output is obtained:

Importance of components:

|  | *PC1* | *PC2* | *PC3* |
|---|---|---|---|
| *Standard deviation* | 1.325 | 1.010 | 0.4757 |
| *Proportion of Variance* | 0.585 | 0.340 | 0.0754 |
| *Cumulative Proportion* | 0.585 | 0.925 | 1.0000 |

# Perform PCA excluding the categorical 'State' variable

pca_data <- data[, c("R.D.Spend", "Administration", "Marketing.Spend")]


# Running and Summarizing PCA

pca_result <- prcomp(pca_data, scale. = TRUE)

summary(pca_result)


Result:

Importance of components:

| | PC1 | PC2 | PC3 |
|---|---|---|---|
| Standard deviation | 1.3246 | 1.0096 | 0.47567 |
| Proportion of Variance | 0.5848 | 0.3398 | 0.07542 |
| Cumulative Proportion | 0.5848 | 0.9246 | 1.00000 |

PCA Interpretation:

The output from the Principal Component Analysis (PCA) indicates the following:

- PC1: The first principal component (PC1) has a standard deviation of 1.325 and explains 58.5% of the total variance in the data. This suggests that PC1 captures the most significant variation across the observations, making it a crucial component for understanding the underlying structure of the dataset.

- PC2: The second principal component (PC2) has a standard deviation of 1.010 and accounts for 34.0% of the variance. Together with PC1, these two components explain 92.5% of the total variance, indicating that most of the information in the original variables is retained in just the first two components.

- PC3: The third principal component (PC3) shows a standard deviation of 0.4757, contributing only 7.54% of the total variance. While this component adds to the cumulative variance, its lower standard deviation implies that it captures significantly less information compared to PC1 and PC2.


Key Takeaway from PCA:


- Dimensionality Reduction: The first two components (PC1 and PC2) can be used for dimensionality reduction while retaining a large portion of the original variance. This can simplify subsequent analyses and visualizations.
- Variance Distribution: The cumulative proportion indicates that the majority of the data's variability is captured by the first two principal components. This implies that subsequent components (like PC3) contribute less significant additional information.


- Insights for Further Analysis: The high variance explained by the first two components suggests that they may reveal meaningful patterns and structures within the data, potentially informing further modeling efforts and interpretations.

Overall, PCA reveals that the first two principal components are critical for understanding the variance in the dataset, allowing for effective dimensionality reduction while retaining essential information. Further analysis of these components can guide decision-making and model-building processes.

Brief Comment on Findings:

1. Exploration of Profit Predictors: Scatter plots revealed varying degrees of linear relationships between profit and the predictors (R&D Spend, Administration, Marketing Spend). R&D Spend showed a strong positive correlation with profit, while the other predictors displayed less clear trends. The visualizations differentiated states, highlighting potential geographical influences on profit.

2. Multiple Linear Regression (MLR): The MLR model identified R&D Spend and Marketing Spend as significant predictors of profit, with a high adjusted R-squared value indicating good model fit. Residual plots suggested a reasonably constant variance, supporting the linearity assumption.

3. Model Selection Techniques: Different regression techniques—exhaustive search, forward elimination, backward elimination, and stepwise regression—yielded varying models. The exhaustive search identified the most optimal model with the lowest Cp value and highest adjusted R-squared, emphasizing the importance of R&D and Marketing spends.

4. Principal Component Analysis (PCA): PCA highlighted that the first two principal components capture over 92.5% of the variance in the data. This finding suggests a strong dimensionality reduction potential, allowing for simplified data representation while retaining critical information.

## CONCLUSION:

Overall, this analysis effectively identifies key predictors influencing profit and demonstrates the utility of various modeling techniques in deriving insights. The results underscore the importance of R&D and Marketing spends while confirming that dimensionality reduction techniques like PCA can enhance data interpretation. By leveraging these findings, businesses can make informed decisions to optimize spending and improve profitability.