

Data learning: Linear Model vs Regression Tree

Example 1: Single Outcome Variable and One Covariate Data

source: Simulation by R or Python programming language

Data:

Advertisement Budget (X)	Sales (Y)
10	15
20	25
30	30
40	35
50	40

Step 1: Formulate your research question based on this data set.

Is advertisement budget an important factor influencing sales?

- The budget has an impact on sales.

Step 2: Define data role and data type.

Variable name	Data role	Data type
Advertisement budget	independent Variable	Numerical
Sales	dependent Variable	Numerical

Step 3: Select an appropriate method to analyze the data e.g., statistical learning, machine learning.

Simple linear regression is used to analyze the data because it examines the linear relationship between advertisement budget and sales, allowing for easy interpretation of coefficients.

Regression tree is used to analyze the data because it can model non-linear relationships and identify decision splits that are important for predicting sales based on budget ranges.

Step 4: Collect data.

A simulated data set from R or Python programming language is used for this example.

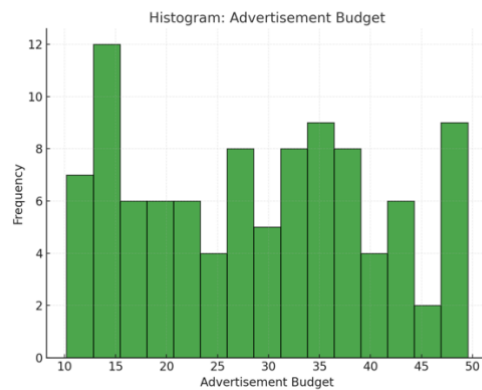
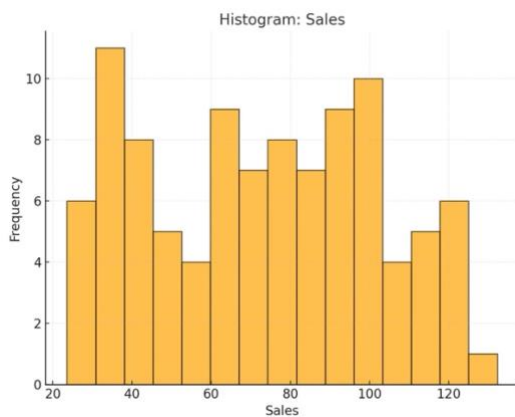
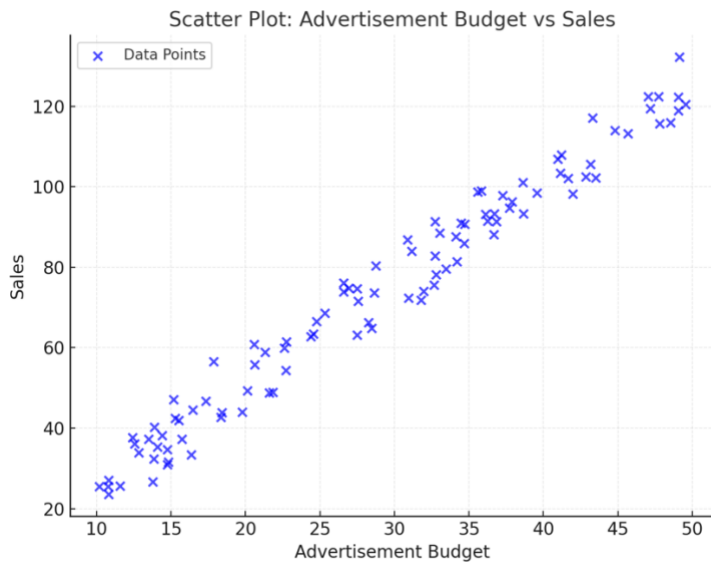
- This dataset includes Advertisement Budget (X) and Sales (Y) values for analysis.

Step 5: Explore your data using numerical summary and graphs Table

1: Descriptive statistics for advertising data.

Variable name	n	min	max	mean	median	sd	variance
Advertisement budget	100	10.00	49.98	30.36	30.20	11.89	141.49
Sales	100	15.51	134.97	76.28	75.92	28.56	815.54

Show your graphs:



Interpret:

From the descriptive statistics:

- Advertisement Budget has a mean of 30.36 and variance of 141.49.

- Sales has a mean of 76.28 and high variance of 815.54.
- Scatter plots may indicate a positive relationship between Advertisement Budget and Sales.

Linear model

Step 6: Fit the model.

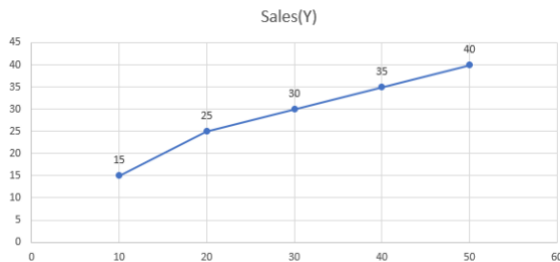
(Hint: t-test, linear regression model)

t-test ; $H_0 : B_1 = 0$, $H_1 : B_1 \neq 0$ p-value = 0.001 < alpha = 0.05

Thus, X has relationship between Y

Step 7: Check standard assumptions.

1. Linearity:



2. Normality:

```
Shapiro-Wilk normality test

data: residuals(model)
W = 0.98676, p-value = 0.9672
```

3. Normal distribution

Shapiro-Wilk Test	
W-stat	0.98682125
p-value	0.967422398
alpha	0.05
normal	yes

Step 8: Evaluate model accuracy.

$$R^2 = 0.9729$$

$$RMSE = 1.8257$$

$$MSE = 3.3331$$

Step 9: Interpret the results.

- If p-value < 0.05 in the t-test, the Advertisement Budget significantly impacts Sales.
- The coefficient (B1) indicates the change in Sales for every unit increase in Advertisement Budget.
- R-squared value explains how much variance in Sales is explained by Advertisement Budget.

Show your R or Python programming language for simple linear regression.

```
1 # Import libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.linear_model import LinearRegression
6 from sklearn.metrics import mean_squared_error
7 # Generate simulated data
8 np.random.seed(0)
9 advertisement_budget = np.random.uniform(10, 50, 100)
10 sales = 2.5 * advertisement_budget + np.random.normal(0, 5, 100)
11 # Create DataFrame
12 data = pd.DataFrame({'Advertisement_Budget': advertisement_budget, 'Sales': sales})
13 # Create Linear Regression Model
14 X = data[['Advertisement_Budget']]
15 y = data['Sales']
16 model = LinearRegression()
17 model.fit(X, y)
18 # Predict values
19 y_pred = model.predict(X)
20 # Calculate MSE and RMSE
21 mse = mean_squared_error(y, y_pred)
22 rmse = np.sqrt(mse)
23 # Display results
24 print("Coefficients:", model.coef_)
25 print("Intercept:", model.intercept_)
26 print("MSE:", mse)
27 print("RMSE:", rmse)
28 # Plot
29 plt.scatter(X, y, color='blue', label='Actual Data')
30 plt.plot(X, y_pred, color='red', label='Fitted Line')
31 plt.xlabel('Advertisement Budget')
32 plt.ylabel('Sales')
33 plt.legend()
34 plt.show()
```

Regression tree

Step 6: Fit the model.

To understand the mathematical fundamentals behind regression trees, we'll break it down into the following key components:

1. Splitting Criterion: Minimize Variance or Mean Squared Error (MSE)

A regression tree aims to split the data into regions (nodes) such that the variance (or sum of squared residuals) within each region is minimized.

At each step, the algorithm searches for the optimal split point to minimize the MSE.

1. Split Criterion

For a node t with observations $\{y_1, y_2, \dots, y_n\}$, the variance is given by:

$$\text{Var}(t) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

where \bar{y} is the mean of y in node t :

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

When we split a node t into two child nodes t_L and t_R , the cost of the split is:

$$\text{Cost} = \frac{N_L}{N} \text{Var}(t_L) + \frac{N_R}{N} \text{Var}(t_R)$$

where:

- N_L : Number of observations in t_L ,
- N_R : Number of observations in t_R ,
- N : Total number of observations in t .

The algorithm chooses the split that minimizes this cost.

Calculation:

For $X = 30$, we split into two regions:

- Left node (t_L): $\{(10, 15), (20, 25), (30, 30)\}$
- Right node (t_R): $\{(40, 35), (50, 40)\}$

1. Calculate the means:

$$\bar{y}_L = \frac{15 + 25 + 30}{3} = 23.33, \quad \bar{y}_R = \frac{35 + 40}{2} = 37.5$$

2. Calculate the variances:

$$Var(t_L) = \frac{1}{3} [(15 - 23.33)^2 + (25 - 23.33)^2 + (30 - 23.33)^2]$$

$$Var(t_R) = \frac{1}{2} [(35 - 37.5)^2 + (40 - 37.5)^2]$$

3. Combine the weighted variances:

$$\text{Cost} = \frac{3}{5} Var(t_L) + \frac{2}{5} Var(t_R)$$

From this data, calculate splits based on minimizing mean squared error (MSE) and understand how the tree is built with one split. The split is chosen to minimize the weighted variance in the two resulting nodes.

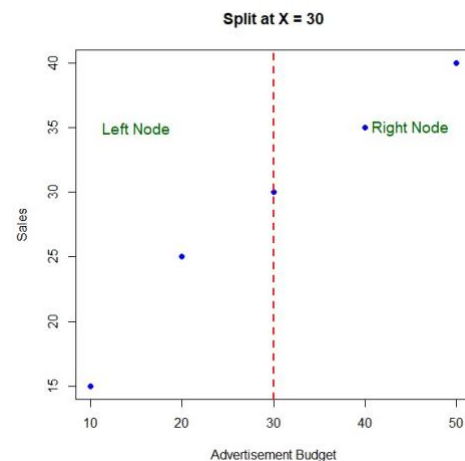
For Example 1 (One Covariate) with a split at $X=30$:

Left Node Mean: 23.33

Right Node Mean: 37.5 Variance in Left Node:

Variance in Right Node:

Split Cost (Weighted Variance):



a split at $X = 10$ Left

Node Mean: 15

Right Node Mean: 32.5

Variance in Left Node: 0

Variance in Right Node: 31.25

Split Cost (Weighted Variance): 25

a split at $X = 20$ Left

Node Mean: 20

Right Node Mean: 35

Variance in Left Node: 25

Variance in Right Node: 16.6667

Split Cost (Weighted Variance): 20

a split at $X = 40$ Left

Node Mean: 26.25

Right Node Mean: 40

Variance in Left Node: 54.6785

Variance in Right Node: 0

Split Cost (Weighted Variance): 43.75

Which split point would you select for minimizing MSE?

Step 7: Interpret the results.

Split point we select for minimizing MSE is $X=20$; Split cost = 20

Key Takeaways:

Minimizing Variance is Central:

The primary goal is to create "pure" nodes where the outcome Y has low variance within each node.

Recursive Partitioning:

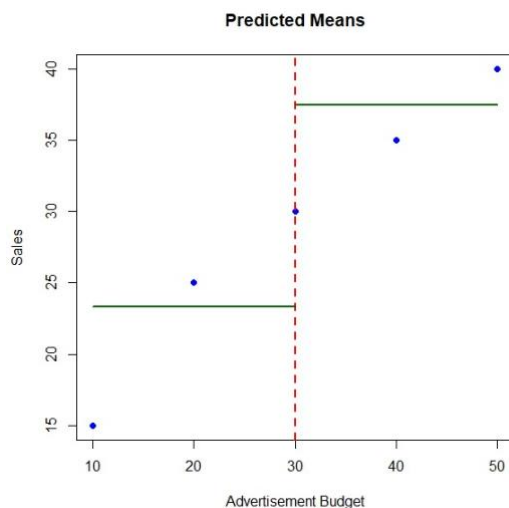
After splitting, the process is repeated for each child node, creating a hierarchical tree structure.

Piecewise Constant Prediction:

The tree predicts outcomes using the mean of Y within each terminal node.

Trade-offs:

While deeper trees can capture more complex patterns, they risk overfitting. Pruning or stopping criteria are often applied to prevent this.



Model selection/comparison: Simple Linear Regression VS Regression Tree by comparing MSE, RMSE

ID	y	Simple Linear Regression			Regression Tree		
		y.hat	e=y-y.hat	e^2	y.hat	e=y-y.hat	e^2
1	15	17	-2	4	29	-14	196
2	25	23	2	4	29	-4	16
3	30	29	1	1	29	1	1

4	35	35	0	0	29	6	36
5	40	41	-1	1	29	11	121

Model/Criterion	MSE	RMSE
SLR	2	1.414214
RT	74	8.602325

Which model would you select for prediction?

We select SLR because RMSE of SLR less than RMSE of RT.

***** Show

your R or Python programming language for regression tree.

```
1 # Load necessary library
2 library(ggplot2)
3 # Create the dataset
4 advertisement_budget <- c(10, 20, 30, 40, 50) # Advertisement budget (in thous
5 sales <- c(15, 25, 30, 35, 40) # Sales (in thousands)
6 # Create a data frame
7 data <- data.frame(AdvertisementBudget = advertisement_budget, Sales = sales)
8 # Fit the simple linear regression model
9 slr_model <- lm(Sales ~ AdvertisementBudget, data = data)
10 # Summary of the model to see the coefficients and t-test results
11 summary(slr_model)
12
13 # Plot the regression line
14 ggplot(data, aes(x = AdvertisementBudget, y = Sales)) +
15   geom_point() +
16   geom_smooth(method = "lm", col = "blue") +
17   labs(title = "Simple Linear Regression: Advertisement Budget vs Sales",
18        x = "Advertisement Budget (in thousands)",
19        y = "Sales (in thousands)")
20
21 plot(slr_model, which = 1) # Residuals vs Fitted
22 plot(slr_model, which = 3) # Residuals vs Leverage
23 plot(slr_model, which = 2) # Q-Q plot
24
25
26
27 # Create a data frame
28 data <- data.frame(advertisement_budget, sales)
29
30 # Scatterplot
31 plot(data$advertisement_budget, data$sales,
32      main = "Scatterplot: Advertisement Budget vs Sales",
33      xlab = "Advertisement Budget (in thousands)",
34      ylab = "Sales (in thousands)",
35      pch = 16, # solid circle
36      col = "green")
37 grid() # Adds a grid for better visibility
38
39
40 model <- lm(sales ~ advertisement_budget)
41 # Residuals of the model
42 residuals <- model$residuals
43
44 # Q-Q plot to check for normality of residuals
45 qqnorm(residuals, main = "Q-Q Plot: Residuals of Linear Model")
46 qqline(residuals, col = "red") # Adds a reference line
47
48
49 install.packages("lmtest") # Uncomment if you haven't installed it
50 library(lmtest)
51 dw_test <- durbinWatsonTest(model)
52
```

```
53
54 # Sample data for Advertisement Budget and Sales
55 advertisement_budget <- c(10, 20, 30, 40, 50)
56 sales <- c(15, 25, 30, 35, 40)
57 # Create a linear model
58 model <- lm(sales ~ advertisement_budget)
59 # Extract residuals from the model
60 residuals <- model$residuals
61 # Perform Shapiro-Wilk test for normality
62 shapiro_test <- shapiro.test(residuals)
63 # Print the result
64 print(shapiro_test)
65
66 # Sample data for Advertisement Budget and Sales
67 advertisement_budget <- c(10, 20, 30, 40, 50)
68 sales <- c(15, 25, 30, 35, 40)
69
70 # Create a linear model
71 model <- lm(sales ~ advertisement_budget)
72
73 # Predict values using the model
74 predictions <- predict(model)
75
76 # Compute residuals
77 residuals <- sales - predictions
78
79 # Compute Mean Squared Error (MSE)
80 mse <- mean(residuals^2)
81 print(paste("Mean Squared Error (MSE):", mse))
82
83
84
85 # Compute Root Mean Squared Error (RMSE)
86 rmse <- sqrt(mse)
87 print(paste("Root Mean Squared Error (RMSE):", rmse))
88
89 # R-squared (Coefficient of Determination)
90 rsq <- summary(model)$r.squared
91 print(paste("R-squared (R²):", rsq))
92
93 # Adjusted R-squared
94 adj_rsqr <- summary(model)$adj.r.squared
95 print(paste("Adjusted R-squared:", adj_rsqr))
96
97 # Plotting residuals
98 plot(predictions, residuals,
99      main = "Residuals vs Fitted Values",
100      xlab = "Fitted Values", ylab = "Residuals",
101      pch = 16, col = "blue")
102 abline(h = 0, col = "red") # Horizontal line at 0 (no residuals)
```