# Lecture 01:
# Performing Data Analysis using R-Language for Linear Regression

**Data Analysis 2 [125J5041]**

**Instructor: Dr. Abar**

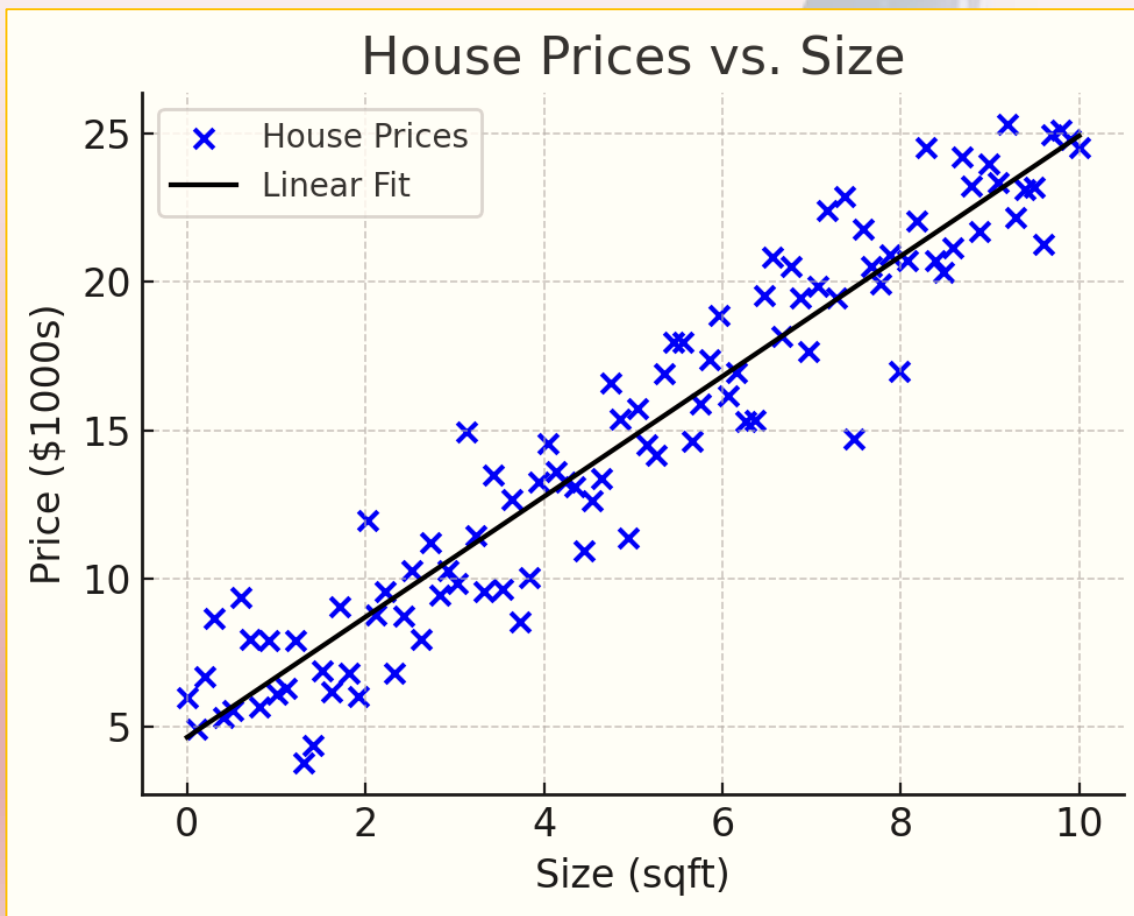**Spring Semester 2025**

- **Regression analysis** is a method for <u>modeling relationships between variables</u>.

- A variable we want to <u>infer or predict</u> is called the <u>dependent variable or outcome</u>; and variables we <u>use for prediction</u> are called <u>independent or explanatory variables or predictors</u>.

- **Linear regression** models the relationship between the independent variables and a dependent variable as a <u>straight/linear line</u>; e.g. since a person's height increases as age increases, age (independent variable) and height (dependent variable) have a linear relationship.

- It is used when the <u>dependent variable is continuous</u> and <u>can take any value within a range</u>. For instance, estimating house prices, measuring temperature, or predicting sales revenue, billing electricity consumption, etc.

- To understand the <u>performance of Linear Regression model</u>, carrying out <u>model evaluation</u> is necessary. <u>Error</u> (i.e. deviation in predicted findings) is typically measured using evaluation metrics like <u>Mean Squared Error (MSE)</u> or <u>Root Mean Squared Error (RMSE)</u>.

# Recap DA1 | Lecture 01: Linear Regression

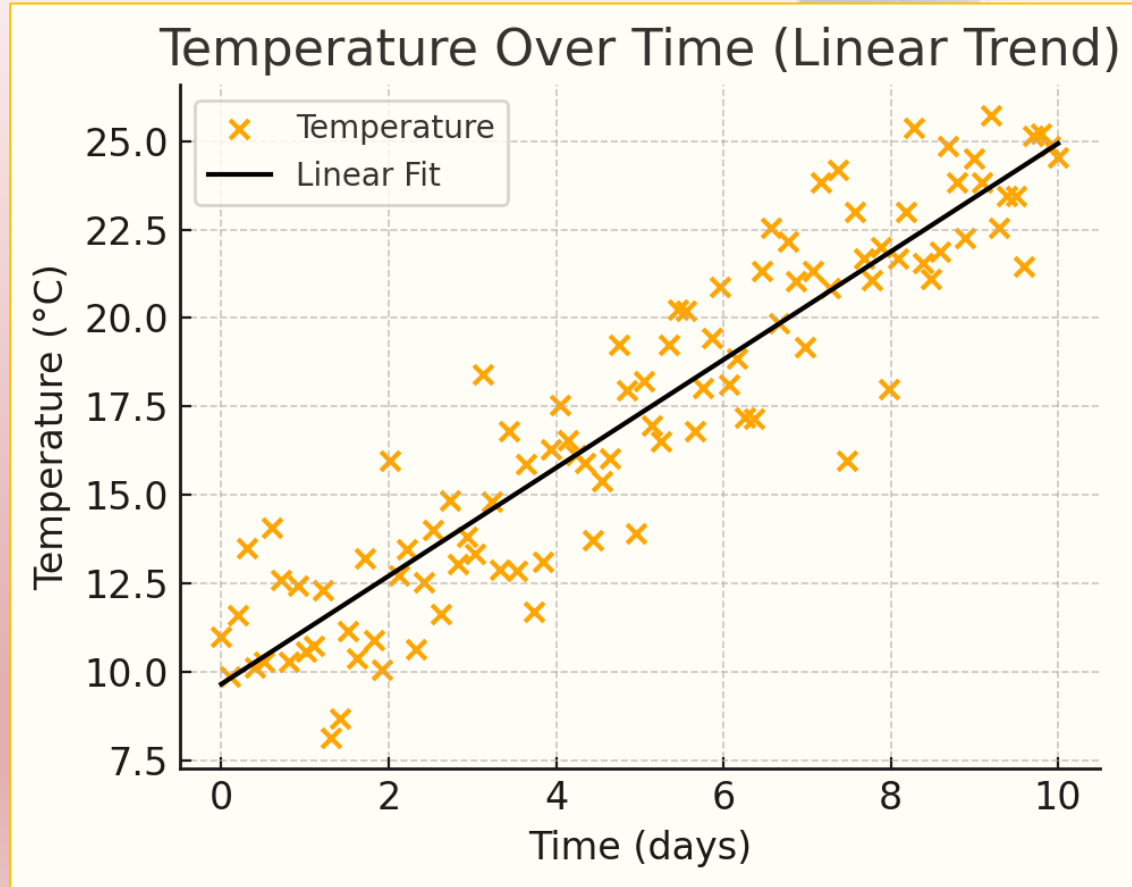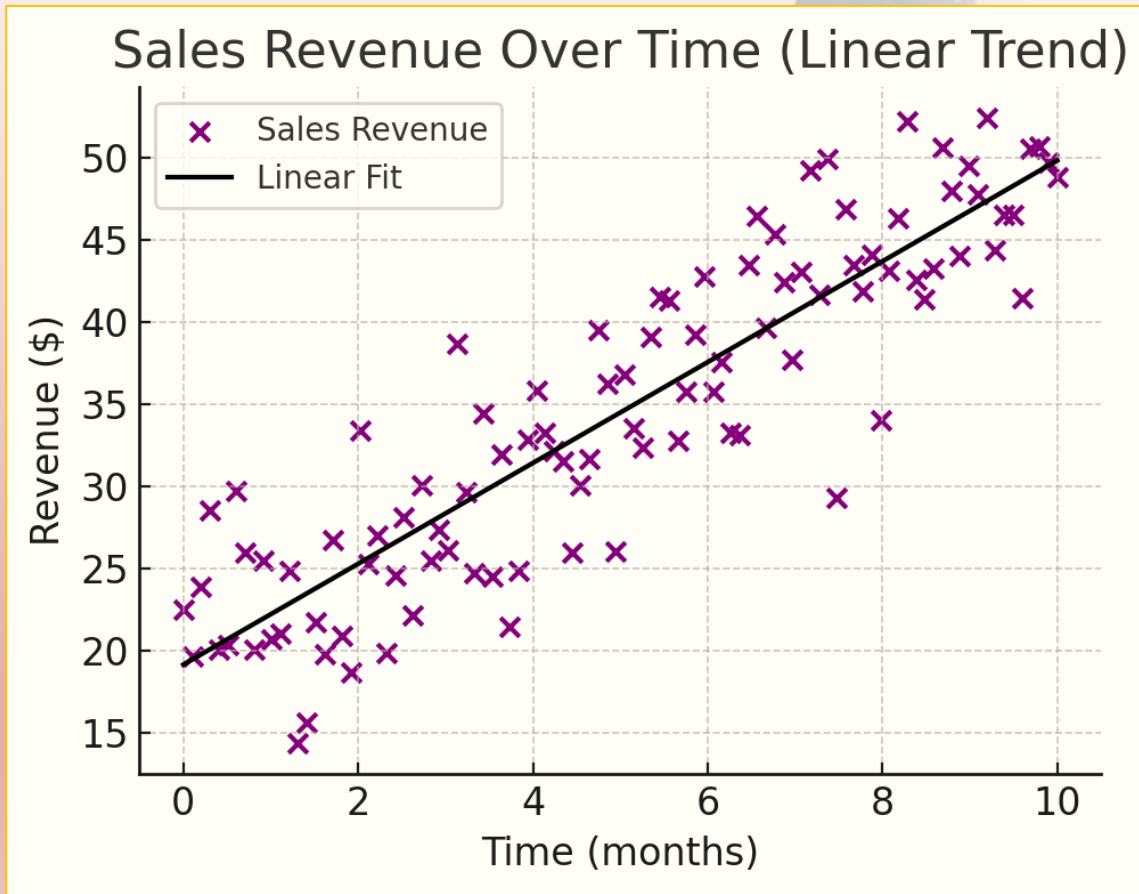**Why Linear regression for Machine Learning? | 4 Minutes Video**

Sales Revenue Over Time (Linear Trend)

Electricity Consumption Over Time (Linear Trend)

## Perform linear regression analysis in RStudio using a sample dataset

- **Objective:** Perform a linear regression analysis <u>to predict the sales of a product</u> based on its <u>advertising expenses</u> using <u>a sample dataset</u>.

- **Dataset:** You will use the "Advertising" dataset, which contains three columns: "TV," "Radio," and "Newspaper" advertising expenses in thousands of dollars, and "Sales" in thousands of units.

- **Instructions:**

1. Open RStudio

2. Click File → New Project → New Directory → New Project → Directory name: → Create Project

3. Create a new R Notebook (from File menu + select R Notebook) and form a new chunk.

4. Copy file **Advertising.csv** in your Directory or Folder created in Step: 1.

5. Load dataset into your R session using following command:
**advertising_data <- read.csv("Advertising.csv")** (Slide: 11)

6. Explore dataset using functions like **head(advertising_data)** for viewing first few rows of your dataset, **summary(advertising_data)** for viewing statistical summaries of your dataset, and **str(advertising_data)** to understand its structure and contents.

**kcg.edu**

## Perform linear regression analysis in RStudio using a sample dataset

7. Click on <u>Preview</u>, R Notebook in HTML format will appear inside RStudio. **(Slide: 12)**

8. Perform a simple linear regression analysis to predict "<u>Sales</u>" based on "<u>TV</u>" <u>advertising expenses</u>. For this, use **lm()** function to create linear regression model:

**model <- lm(Sales ~ TV, data = advertising_data)**

9. Print the summary of regression model using **summary(model).** (**Click on Preview again and refer to Slide: 14).**

10. Interpret results from the regression summary, paying attention to coefficients, p-values, and R-squared. **(see Slide: 17 Appendix)**

11. Create a scatter plot to visualize the relationship between "TV" advertising expenses and "Sales":

**plot(advertising_data$TV, advertising_data$Sales, xlab = "TV Advertising Expenses", ylab = "Sales", main = "Scatter Plot of TV vs. Sales")**

**abline(model, col = "red")**   [R function abline() is used to draw regression lines to a graph]

**(Click on Preview again and see Slide: 15)**

12. Calculate predicted values using your regression model:

**predicted_sales <- predict(model, newdata = data.frame(TV = advertising_data$TV))** [This line of code calculates predicted sales for each value of "TV" advertising expenses in your original dataset and stores the results in "predicted_sales" object]

**predicted_sales**

13. Calculate residual values (differences between actual and predicted values):

**residuals <- advertising_data$Sales - predicted_sales**

**residuals**

14. Calculate mean squared error (MSE) to measure model's accuracy:

**mse <- mean(residuals^2)**

**mse**

15. Discuss results, including interpretation of coefficients and model's accuracy. **(Next Page)**

## Perform linear regression analysis in RStudio using a sample dataset

**15. Discuss the results, including the interpretation of coefficients and the model's accuracy.**

Let's discuss the results of linear regression analysis performed in this exercise, focusing on the <mark>interpretation of coefficients</mark> and <mark>model's accuracy</mark>.

### 1. Coefficients:

<u>Linear regression model you created aims to predict "Sales" based on "TV" advertising expenses</u>. Coefficient for "TV" predictor variable can be found in the regression summary output.

Coefficient for "TV" represents the estimated change in "Sales" for a one-unit increase in "TV" advertising expenses, holding other variables constant. If the coefficient is positive, it implies that as "TV" advertising expenses increase, "Sales" is expected to increase. Conversely, a negative coefficient would imply a negative relationship.

<u>In your regression summary, if coefficient for "TV" is, for instance, 0.0475, this means that for each additional unit of money spent on TV advertising, the expected increase in sales is 0.0475 units, assuming all other factors remain constant.</u>

### 2. P-values:

<u>In regression summary, you will also observe p-value associated with the coefficient for "TV." This p-value is used to test the null hypothesis that coefficient is equal to zero (no effect).</u> <u>If the p-value is less than your chosen significance level (commonly 0.05), it indicates that the coefficient is statistically significant. In this case, a low p-value $2 \times 10^{-16}$ would suggest that "TV" advertising expenses have a statistically significant impact on "Sales".</u>

### 3. R-squared ($R^2$):

R-squared is a measure of the goodness of fit of your regression model. In your regression summary, you should see the R-squared value.

<u>An R-squared value closer to 1 indicates that the model explains a larger proportion of the variance in "Sales." For example, if R-squared is 0.75, it means that 75% of the variability (unpredictability/inconsistency/unevenness) in "Sales" is explained by the "TV" advertising expenses in your model.</u>

### Interpretation:

- If the coefficient for "TV" is positive and statistically significant (i.e. <mark>p-value < 0.05</mark>), you can conclude that there is a positive linear relationship between "TV" advertising expenses and "Sales". Therefore as you invest more in TV advertising, you can expect an increase in sales.
- <mark>R-squared</mark> value tells how well the model fits data. A higher R-squared suggests that "TV" advertising expenses explain a significant portion of the variance in "Sales."

Remember that linear regression assumes a linear relationship between the predictor and the response variable. Interpretation of the coefficients and model's accuracy is based on this assumption. If relationship is not truly linear, model may not be accurate in making predictions.
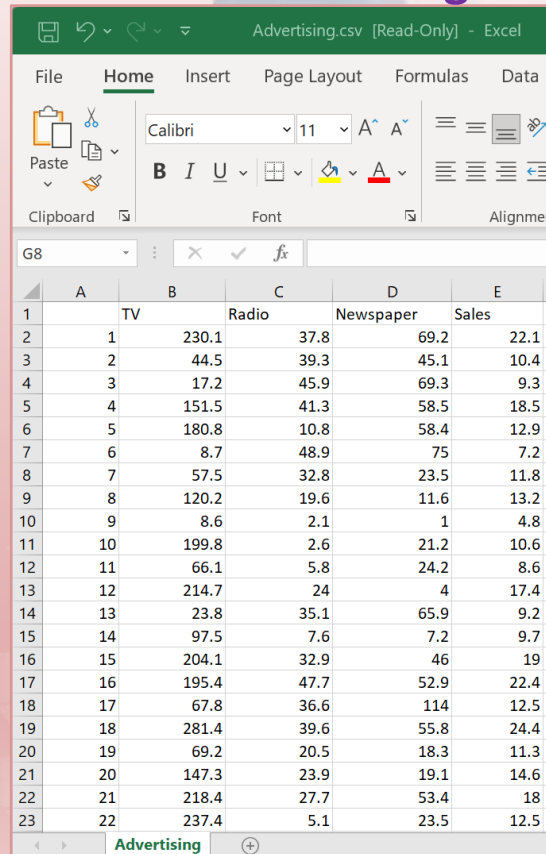
It's also essential to consider other factors, such as the residual plots and overall context of your analysis for assessing model's validity and whether it provides meaningful insights into the relationship between advertising expenses and sales.

# Exercise on Linear Regression [4/8]

## DATASET: Advertising.csv

Download this dataset from:

https://www.kaggle.com/datasets/bumba5341/advertisingcsv/

| | A | B TV | C Radio | D Newspaper | E Sales |
|---|---|---|---|---|---|
| 1 | | TV | Radio | Newspaper | Sales |
| 2 | 1 | 230.1 | 37.8 | 69.2 | 22.1 |
| 3 | 2 | 44.5 | 39.3 | 45.1 | 10.4 |
| 4 | 3 | 17.2 | 45.9 | 69.3 | 9.3 |
| 5 | 4 | 151.5 | 41.3 | 58.5 | 18.5 |
| 6 | 5 | 180.8 | 10.8 | 58.4 | 12.9 |
| 7 | 6 | 8.7 | 48.9 | 75 | 7.2 |
| 8 | 7 | 57.5 | 32.8 | 23.5 | 11.8 |
| 9 | 8 | 120.2 | 19.6 | 11.6 | 13.2 |
| 10 | 9 | 8.6 | 2.1 | 1 | 4.8 |
| 11 | 10 | 199.8 | 2.6 | 21.2 | 10.6 |
| 12 | 11 | 66.1 | 5.8 | 24.2 | 8.6 |
| 13 | 12 | 214.7 | 24 | 4 | 17.4 |
| 14 | 13 | 23.8 | 35.1 | 65.9 | 9.2 |
| 15 | 14 | 97.5 | 7.6 | 7.2 | 9.7 |
| 16 | 15 | 204.1 | 32.9 | 46 | 19 |
| 17 | 16 | 195.4 | 47.7 | 52.9 | 22.4 |
| 18 | 17 | 67.8 | 36.6 | 114 | 12.5 |
| 19 | 18 | 281.4 | 39.6 | 55.8 | 24.4 |
| 20 | 19 | 69.2 | 20.5 | 18.3 | 11.3 |
| 21 | 20 | 147.3 | 23.9 | 19.1 | 14.6 |
| 22 | 21 | 218.4 | 27.7 | 53.4 | 18 |
| 23 | 22 | 237.4 | 5.1 | 23.5 | 12.5 |

## Click on Preview; then R Notebook (HTML) will appear on the right-side pane

From Windows File Explorer, you can open R Notebook (HTML) in your browser; refresh every time you make changes by executing commands in RStudio

# Exercise on Linear Regression

On executing Step. 9 | RStudio IDE

**After Step. 11, RStudio IDE**

# Homework Activity: Further Practice

**To practice <span style="color:darkred">linear regression analysis</span> in <span style="color:blue">RStudio</span> and understand how to interpret the results of a regression model; repeat analysis for "<span style="color:darkred">Radio</span>" and "<span style="color:darkred">Newspaper</span>" advertising expenses to see how these relate to "<span style="color:darkred">Sales</span>".**

**Analyze the findings obtained in your R Notebook HTML format using the guidance provided in Appendix (Next Slide)**

## Step 10. Interpret results from the regression summary, paying attention to coefficients, p-values, and R-squared.

Interpreting the results from a regression summary in R, focusing on coefficients, p-values, and R-squared, is essential to understand the relationships within your dataset. Here's how to interpret these key elements:

### 1. Significance (importance) of Coefficients:
 - In a linear regression summary, you will see a table that includes coefficients for each predictor variable (independent variable) in your model. These coefficients represent change in dependent variable (in your case, "Sales") for a one-unit change in predictor variable, while holding other predictors constant.
 - Coefficients are typically labeled as "Estimate," and they represent the slope of regression line. For instance, if coefficient for "TV" is 0.0475, it means that for every additional unit of TV advertising expenses, Sales is expected to increase by 0.0475 units, assuming all other factors remain constant.

### 2. P-values: Validating Null/Alternative Hypothesis
 - P-values associated with each coefficient test the null hypothesis that the coefficient is equal to zero (no effect). A low p-value (usually less than 0.05) indicates that predictor variable is statistically significant in predicting the dependent variable.
 - If the p-value is less than your chosen significance level (commonly 0.05), you can reject the null hypothesis and conclude that predictor has a statistically significant effect on the dependent variable.

### 3. R-squared ($R^2$):
 - R-squared is a measure of how well the independent variables (predictors) explain the variability in dependent variable. It is a value between 0 and 1, where a higher value indicates a better fit of model to the data.
 - An R-squared value of 0 means that the model does not explain any of the variability in the dependent variable, while an R-squared value of 1 means that the model explains all the variability.
 - In the context of linear regression, R-squared measures the proportion of the variance in the dependent variable that is predictable from the independent variables. For example, if R-squared is 0.75, it means that 75% of the variance in Sales is explained by the predictor variables in your model.

In a nutshell, when interpreting the regression summary:

- Look at coefficients to understand the direction and strength of relationship between predictor variables and dependent variable.
- Examine p-values to determine whether predictor variables are statistically significant.
- Assess R-squared to gauge how well your model explains the variability in dependent variable.

A high R-squared, significant coefficients, and their directions can help you draw meaningful conclusions about the relationships within your data and predictive power of your model.

# Practice Class Activity: Linear Regression using Some Synthetic Data [1/2]

**This code generates synthetic data, fits a linear regression model, provides a summary of the model, plots data and regression line, creates a residual plot, calculates Mean Squared Error (MSE), and finally, offers a brief analysis of the outcome.**

**Step 1: Generate synthetic data**

```
set.seed(123)  # Set a seed for reproducibility
x <- 1:20
y <- 2 * x + rnorm(20, mean = 0, sd = 2)  # Generate synthetic data with some noise
```

**Step 2: Perform linear regression**

```
lm_model <- lm(y ~ x)
```

**Step 3: Create a summary of regression modeling**

```
summary(lm_model)
```

**Step 4: Plot data and regression line**

```
plot(x, y, main = "Linear Regression Example", xlab = "X", ylab = "Y")
abline(lm_model, col = "red")
```

**Step 5: Create a residual plot**

```
par(mfrow = c(1, 2))
plot(lm_model, which = 1)  # Residuals vs Fitted
plot(lm_model, which = 2)  # Normal Q-Q Plot
```

**Step 6: Calculate the Mean Squared Error (MSE)**

```
predicted <- predict(lm_model)
mse <- mean((y - predicted)^2)
cat("Mean Squared Error (MSE):", mse, "\n")
```

**Step 7: Analyze the outcome**

You can analyze the coefficient estimates, p-values, and R-squared from summary output.

Additionally, MSE provides a measure of the model's goodness of fit where lower values indicate better fit.