

# regression-1

September 13, 2023

```
[ ]: # Q1. Explain the difference between simple linear regression and multiple
      ↪ linear regression.

      # Provide an example of each
```

Simple Linear Regression is a statistical method used to model the relationship between two variables - one independent variable (predictor) and one dependent variable (response). The goal is to find the best-fitting straight line that represents the relationship between the variables.

Example of Simple Linear Regression:

Suppose we want to predict a student's final exam score (y) based on the number of hours they studied (x). We collect data from several students and find the following relationship:

Hours Studied (x) Exam Score (y) 2 60 3 70 4 75 5 85 6 90

Using simple linear regression, we can find the best-fitting line that represents this relationship and use it to predict exam scores for different levels of study hours.

Multiple Linear Regression: Multiple Linear Regression extends the concept of simple linear regression to include more than one independent variable. It models the relationship between a dependent variable and multiple independent variables.

Example of Multiple Linear Regression:

Example of Multiple Linear Regression:

Let's consider a real estate example. We want to predict the price of a house (y) based on its size in square feet (X1), the number of bedrooms (X2), and the age of the house in years (X3).

Size (x1) Bedrooms (x2) Age (x3) Price (y) 1500 3 10 250,000 2000 4 5 320,000 1800 3 12 280,000 2200 4 8 350,000 1600 2 15 230,000

Using multiple linear regression, we can find the best-fitting plane in a three-dimensional space that represents the relationship between these variables and use it to predict house prices based on their size, number of bedrooms, and age.

```
[ ]: # Q2. Discuss the assumptions of linear regression.

      # How can you check whether these assumptions hold in a given dataset?
```

Assumptions in Regression : Regression is a parametric approach. 'Parametric' means it makes assumptions about data for the purpose of analysis. Due to its parametric side, regression is restric-

tive in nature. It fails to deliver good results with data sets which doesn't fulfill its assumptions. Therefore, for a successful regression analysis, it's essential to validate these assumptions.

Violations of assumptions of linear regression can lead to biased or inefficient estimates, and it is important to assess and address these violations for accurate and reliable regression results. Assumptions of linear regression include:

Linearity: The relationship between the dependent and independent variables is linear. Independence: The observations are independent of each other. Homoscedasticity: The variance of the errors is constant across all levels of the independent variables. Normality: The errors follow a normal distribution. No multicollinearity: The independent variables are not highly correlated with each other. No endogeneity: There is no relationship between the errors and the independent variables.

Linearity: The relationship between the independent and dependent variables should be linear.

Independence of Errors: The errors (residuals) should be independent of each other, meaning that the error of one observation should not provide information about the error of another observation.

Normality of Residuals: The residuals should follow a normal distribution.

No Multicollinearity: In multiple linear regression, the independent variables should not be highly correlated with each other. High multicollinearity can lead to inflated standard errors and difficulty in interpreting the individual effects of each variable

No Outliers: Outliers can strongly influence the regression model, leading to biased estimates. It's important to identify and handle outliers appropriately.

Checking Assumptions:

Linearity: You can create scatter plots of each independent variable against the dependent variable to visually assess linearity.

Independence of Errors: This assumption can be assessed using residual plots. I

Homoscedasticity: Residual plots can also help with assessing homoscedasticity.

Normality of Residuals: A histogram or a normal probability plot of the residuals can provide insight into their distribution

No Multicollinearity: Calculate the correlation matrix between independent variables. If correlations are very high (close to 1 or -1), multicollinearity might be an issue.

No Outliers: Create a scatter plot of the standardized residuals or leverage values. Points that are far from the main cluster might be outliers

```
[ ]: # Q3. How do you interpret the slope and intercept in a linear regression model?  
    ↪  
  
    # Provide an example using a real-world scenario.
```

Intercept (b0): The intercept represents the estimated value of the dependent variable when the independent variable(s) are zero.

Slope (b1): The slope represents the change in the dependent variable for a one-unit change in the independent variable.

Example:

Let's consider a real-world scenario where we want to predict a person's monthly electricity bill (y) based on the number of kilowatt-hours (kWh) they consumed (x).

We use linear regression to model this relationship and obtain the following equation:  $y = 50 + 0.1x$

Here, y is the monthly bill, x is the number of kWh consumed, and 50 is the intercept, and 0.1 is the slope.

Intercept (): When a person consumes zero kWh (which is unlikely and not practical), the predicted monthly bill would be \$50. This intercept is the base cost that the person would need to pay even if they don't use any electricity.

Slope (): For each additional kWh consumed, the monthly bill increases by \$0.1.

For example, if a person consumed 300 kWh in a month, we can predict their monthly bill using the equation:  $y = 50 + 0.1 \times 300 = 50 + 30 = 80$

So, based on the model, we would predict that their monthly bill would be \$80.

```
[ ]: # Q4. Explain the concept of gradient descent. How is it used in machine learning?
```

Gradient Descent is an optimization algorithm used to minimize (or maximize) a function iteratively. It's a widely used technique in machine learning for training models, especially when dealing with models that have adjustable parameters (weights and biases) and an associated cost or loss function that needs to be minimized.

Usage in Machine Learning:

Gradient descent is fundamental in training machine learning models, particularly for optimizing the parameters of models like linear regression, logistic regression, neural networks, and more. In these models, the goal is to find the set of parameters that minimizes a cost or loss function, which quantifies how well the model's predictions match the actual data.

Gradient descent is a key optimization technique used in machine learning to iteratively update model parameters and minimize a cost function, enabling models to learn from data and make accurate predictions.

```
[ ]: # Q5. Describe the multiple linear regression model. How does it differ from simple linear regression?
```

Multiple Linear Regression is an extension of simple linear regression that allows for the modeling of relationships between a dependent variable and multiple independent variables

In simple linear regression, there is only one independent variable, whereas in multiple linear regression, there are two or more independent variables.

```
[ ]: # Q6. Explain the concept of multicollinearity in multiple linear regression.
# How can you detect and address this issue?
```

Multicollinearity is a phenomenon in multiple linear regression where two or more independent variables in the model are highly correlated with each other. In other words, there is a strong linear relationship between at least two of the independent variables. Multicollinearity can create problems in the regression analysis and affect the stability and interpretability of the model's coefficients.

Detecting Multicollinearity: Correlation Matrix: Calculate the correlation coefficients between all pairs of independent variables.

Addressing Multicollinearity:

Feature Selection: If multicollinearity is detected, consider removing one or more of the highly correlated variables from the model. However, be cautious not to remove variables that are theoretically important or necessary for the analysis.

```
[ ]: # Q7. Describe the polynomial regression model. How is it different from linear ↵  
      ↪ regression?
```

Polynomial Regression is a type of regression analysis that extends the idea of linear regression by introducing polynomial terms as predictors. In polynomial regression, the relationship between the independent variable(s) and the dependent variable is modeled as an nth-degree polynomial equation. This allows the model to capture more complex and nonlinear relationships between the variables.

A simple linear regression algorithm only works when the relationship between the data is linear. But suppose we have non-linear data, then linear regression will not be able to draw a best-fit line. Simple regression analysis fails in such conditions.

Polynomial regression is a form of Linear regression where only due to the Non-linear relationship between dependent and independent variables, we add some polynomial terms to linear regression to convert it into Polynomial regression.

```
[ ]: # Q8. What are the advantages and disadvantages of polynomial regression ↵  
      ↪ compared to linear regression? \  
  
# In what situations would you prefer to use polynomial regression?
```

Advantages of Polynomial Regression:

Flexibility: Polynomial regression can capture complex relationships and patterns that linear regression cannot. It can fit curves, bends, and nonlinear trends in data.

Better Fit: When the relationship between the variables is nonlinear, polynomial regression can provide a better fit to the data, leading to more accurate predictions.

Disadvantages of Polynomial Regression:

Overfitting: With higher-degree polynomials, there's a risk of overfitting the model to noise in the data, resulting in poor generalization to new, unseen data.

Complexity: Higher-degree polynomials introduce more parameters into the model, making it more complex and harder to interpret.

When to Prefer Polynomial Regression:

Polynomial regression is a valuable tool when the relationship between the dependent and independent variables is nonlinear and simple linear regression is not sufficient to capture the pattern.

Curved Relationships: When you visually observe that the scatter plot of data suggests a curved or nonlinear relationship.

[ ]: