**1. What is the difference between a dependent variable and an independent variable in a linear equation?**

- The dependent variable is the variable that is being predicted or explained by the independent variable(s). Its value depends on the values of the independent variable(s).

- The independent variable, on the other hand, is the variable that is believed to have an influence on the dependent variable. Its value is manipulated or controlled to observe its effect on the dependent variable.

**2. What is the concept of simple linear regression? Can you provide an example?**

- Simple linear regression is a statistical technique used to model the relationship between a single independent variable and a dependent variable. It assumes a linear relationship between the variables and aims to find the best-fit line that minimizes the sum of the squared differences between the observed and predicted values.

- Example: Suppose we want to study the relationship between the number of hours studied (independent variable) and the exam score (dependent variable) of a group of students. Using simple linear regression, we can estimate how the exam score changes for each additional hour studied.

**3. How would you define the slope in linear regression?**

- In linear regression, the slope represents the change in the dependent variable (y) per unit change in the independent variable (x). It indicates the rate at which the dependent variable is expected to change when the independent variable changes.

**4. Given two points (3, 2) and (2, 2) on a line, what is the slope of the line?**

- The slope of a line is calculated as the change in the y-coordinates divided by the change in the x-coordinates. In this case, both points have the same y-coordinate (2), so the change in y is 0. The change in x is (2 - 3) = -1. Therefore, the slope is 0 divided by -1, which is 0.

**5. What are the conditions for a positive slope in linear regression?**

- In linear regression, a positive slope indicates that there is a positive relationship between the independent and dependent variables. The conditions for a positive slope are when the points on the scatter plot generally increase as the independent variable increases, indicating a direct relationship.

**6. What are the conditions for a negative slope in linear regression?**

- In linear regression, a negative slope indicates that there is a negative relationship between the independent and dependent variables. The conditions for a negative slope are when the points on

the scatter plot generally decrease as the independent variable increases, indicating an inverse relationship.

### 7. What is multiple linear regression and how does it work?

- Multiple linear regression is a statistical technique used to model the relationship between multiple independent variables and a dependent variable. It extends the concept of simple linear regression to include multiple predictors.

- In multiple linear regression, the technique estimates the coefficients for each independent variable to find the best-fit line or hyperplane that minimizes the sum of the squared differences between the observed and predicted values of the dependent variable.

### 8. How would you define the sum of squares due to error in multiple linear regression?

- The sum of squares due to error, also known as the residual sum of squares (RSS) or the sum of squared residuals, is a measure of the overall discrepancy between the observed values of the dependent variable and the predicted values from the multiple linear regression model. It quantifies the unexplained variation in the dependent variable.

### 9. How would you define the sum of squares due to regression in multiple linear regression?

- The sum of squares due to regression, also known as the explained sum of squares (ESS), is a measure of the variation in the dependent variable that is accounted for by the multiple linear regression model. It quantifies the contribution of the independent variables in explaining the variation in the dependent variable.

### 10. What is multicollinearity in a regression equation?

- Multicollinearity refers to the high correlation or interdependence between two or more independent variables in a regression equation. It occurs when the independent variables are linearly related to each other, which can lead to issues in interpreting the individual effects of the variables and can affect the stability and reliability of the regression model.

### 11. What is heteroskedasticity, and what does it mean?

- Heteroskedasticity refers to the violation of the assumption of constant variance of errors in a regression model. It occurs when the variability of the errors (residuals) is not the same across all levels of the independent variables. In other words, the spread or dispersion of the residuals differs for different values of the independent variables, indicating a systematic pattern of variability.

**12. Can you describe the concept of ridge regression?**

- Ridge regression is a regularization technique used to mitigate multicollinearity and reduce the impact of highly correlated independent variables in a regression model. It adds a penalty term to the ordinary least squares estimation, which shrinks the coefficients of the independent variables. This penalty term helps to stabilize the estimates and can prevent overfitting by reducing the model's sensitivity to the data.

**13. Can you describe the concept of lasso regression?**

- Lasso (Least Absolute Shrinkage and Selection Operator) regression is another regularization technique used to select relevant independent variables and shrink the coefficients. It adds a penalty term to the ordinary least squares estimation, which includes the sum of the absolute values of the coefficients. Lasso regression encourages sparsity in the coefficient estimates, effectively setting some coefficients to exactly zero and performing feature selection.

**14. What is polynomial regression and how does it work?**

- Polynomial regression is a form of regression analysis where the relationship between the independent variable(s) and the dependent variable is modeled as an nth-degree polynomial. It allows for curved or nonlinear relationships to be captured by introducing polynomial terms of the independent variables. The regression model can then estimate the coefficients of these polynomial terms to fit the data.

**15. Can you describe the basis function?**

- In the context of regression analysis, a basis function is a mathematical function used to represent the relationship between the independent variable(s) and the dependent variable. The basis function transforms the input variables into a new set of variables, which are then used as inputs for the regression model. Common types of basis functions include polynomial functions, exponential functions, and sigmoid functions.

**16. How does logistic regression work?**

- Logistic regression is a statistical model used for binary classification problems, where the dependent variable is categorical and has two possible outcomes (e.g., 0 or 1). It estimates the probability of the dependent variable belonging to a particular category based on the values of the independent variables. Logistic regression uses a logistic (sigmoid) function to map the linear combination of the independent variables to the range of [0, 1], representing the probability of the event occurring.