

Interview Questions

Q. What is Normalization & Standardization and how is it helpful?

Normalization & Standardization are techniques used to scale numerical data. This can be useful to machine learning model process the data effectively.

Normalization = rescale of data fit range in specific way between 0 & 1

$$\text{norm}(x) \equiv \frac{x - \min(x)}{\max(x) - \min(x)}$$

The dataset has values ranging from 50 to 100. The normalization will scale these into the range [0, 1].

* The features is different unit height cm & weight kg

* The model sensitive is k-NN, k-mean, Neural Network.

Standardization :-

Standardization transforms data to have a mean of 0 & standard deviation of 1. This process converts the data to follow a standard normal distribution.

Formula

$$x_{\text{std}} = \frac{x - \bar{x}}{\sigma}$$

where \bar{x} is mean & σ is standard deviation

* \bar{x} = mean of feature

* σ = standard deviation of the feature

④ How is it helpful in Normalization & Standardization?

① Improve model performance:-

* speed up gradient descent

converge in optimization based on model (Ex:- logistic regression, Neural Network)

get rid of outliers & improve time

② Handle different Scale:-

* Avoid features with large

range dominating the model

③ Reduce Bias:-

- * Ensure all features contribute equally to the model's prediction.

④ Stability & Accuracy:-

- * Reduces numerical errors and improve the stability of computations.

⑤ Better distance calculation:-

- * Essential for distance-based models like kNN, k-means, SVM.

Q. What techniques can be used to address multicollinearity in multiple linear regression?

Multicollinearity occurs when independent variables in regression model are highly correlated with each other.

* Unstable co-efficients

* Reduced interpretability

* Inflated standard errors.

* Inflated standard errors.

① Detecting Multicollinearity

① Variance Inflation Factor (VIF):

- * Measures how much the variance of coefficient increases due to multicollinearity.
- * Rule of Thumb: If VIF > 5 multicollinearity present.

Formula:-

$$\text{VIF} = \frac{1}{1 - R^2}$$

* Correlation matrix:-

* Check the Pair wise

Pearson correlation co-efficients

* Correlations > 0.8 suggest

multicollinearity.

② Techniques to Handle Multicollinearity

A) Data driven Solutions:

- 1) Remove one of the correlated features:-

Drop the variable with the least predictive power.

2) Combine Correlated Features

* Create new features by combining correlated variables.

3) Principal Component Analysis (PCA)

* Reduces dimensionality while keeping most of the variance.

Want AIG to minimize *

B) Statistical Model Adjustments:-

4) Ridge Regression (L₂ Regularization)

* Add a penalty term to the loss function

co-efficients

formula:-

or other forms are good

$$\min \|y - x\beta\|^2 + \lambda \|\beta\|^2$$

5) LASSO Regression (L₁ Regularization)

* Shrinks some co-efficients to zero, effectively performing feature selection.

formula:-

$$\min \|y - x\beta\|^2 + \lambda \sum |\beta_i|$$

6) Elastic Net regression:
* A combination of Ridge & LASSO regularization.

(7) Diagnostic techniques:

c) model interpretation techniques:

7) Partial Least Squares (PLS):

* Similar to PCA but focuses on maximizing the explanation variance of the dependent variable.

8) Drop highly correlated Dummy variables:

* If using categorical variables, drop one dummy variable to prevent dummy variable trap.

(not discussed) 1) missing values

2) outliers - use standard deviation

3) multicollinearity - use correlation matrix

- tolerance

- VIF

$$191.21 + 118x_1 - 0.11x_2$$