

Assignment No. 3

Problem Statement: CORRELATION AND COVARIANCE:

- a. Find the correlation matrix on the iris dataset.
- b. Plot the correlation plot on the dataset and visualize giving an overview of relationships among data on iris dataset.

Objective: The objective is to calculate the correlation matrix of the Iris dataset and visualize the relationships among its variables through a correlation plot, providing insight into how the features of the dataset relate to each other.

Prerequisite :

1. A Python environment set up with libraries like pandas, xml.etree.ElementTree, and requests (for web access).
2. Internet connection (for reading datasets from the web).
3. Text editor and basic knowledge of file handling in Python.

Theory :

Correlation:

Correlation is a statistical measure that describes the strength and direction of the relationship between two variables. It quantifies how changes in one variable are associated with changes in another. The correlation coefficient, which ranges between -1 and 1, is used to represent this relationship:

- **Positive Correlation:** A correlation coefficient close to +1 indicates a strong positive relationship, meaning as one variable increases, the other also increases.
- **Negative Correlation:** A correlation coefficient close to -1 indicates a strong negative relationship, meaning as one variable increases, the other decreases.
- **No Correlation:** A coefficient near 0 suggests no significant relationship between the variables.

The formula for correlation (Pearson's correlation) is:

$$r = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}$$

Where r is the correlation coefficient, X_i and Y_i are individual data points, and \bar{X} and \bar{Y} are their respective means.

Correlation Matrix:

A correlation matrix is a table that shows the correlation coefficients between multiple variables. Each cell in the matrix contains the correlation coefficient between two variables. For example, in the Iris dataset, the matrix will show how the features like sepal length, sepal width, petal length, and petal width correlate with each other. The matrix helps identify patterns such as strong relationships between certain features.

Covariance:

Covariance measures how much two variables change together, giving a sense of the direction of their linear relationship. It indicates whether the variables tend to increase or decrease together:

- **Positive Covariance:** If the covariance is positive, it suggests that as one variable increases, the other tends to increase as well.
- **Negative Covariance:** A negative covariance suggests that as one variable increases, the other tends to decrease.
- **Zero Covariance:** No linear relationship exists if the covariance is zero.

Covariance is calculated using the formula:

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{n}$$

where X and Y are the variables, \bar{X} and \bar{Y} are their respective means, and n is the number of data points.

Correlation vs. Covariance:

- **Correlation** normalizes the covariance, making it unitless and easier to interpret. It provides both the strength and direction of the linear relationship, while covariance only gives the direction.
- Correlation is bounded between -1 and 1, whereas covariance can take any value, making correlation more useful for comparing the strength of relationships across different pairs of variables.

Visualization of Correlation:

To visualize the relationships between variables, a **correlation plot** (such as a heatmap) is used. This graphical representation of the correlation matrix provides an easy way to observe which variables have strong or weak correlations. In a heatmap, darker colors indicate stronger correlations, while lighter colors show weaker ones. Positive and negative correlations are often represented using different color schemes (e.g., blue for positive and red for negative correlations). Such plots help identify potential patterns in the dataset, revealing which features may be redundant or useful for predictive modeling.

In the case of the Iris dataset, visualizing the correlation can help understand how sepal length, sepal width, petal length, and petal width relate to each other, which is particularly valuable for feature selection and data reduction in machine learning.

References :

Database <https://www.kaggle.com/datasets>

Conclusion :

Understanding the correlation between features in a dataset is crucial for identifying relationships and dependencies. The correlation matrix and its visualization provide an overview of how variables interact, guiding decisions in data preprocessing, feature selection, and model development. This analysis is fundamental to improving model performance and eliminating redundant features.