# Internship Weekly Report - Week 6

**Name:** Sandeep Ravaji Patel

**Domain:** Data Science

**Week Number:** Week 6

**Objective:**

To understand clustering and dimensionality reduction techniques, specifically focusing on implementing K-Means clustering and performing PCA using Scikit-Learn.
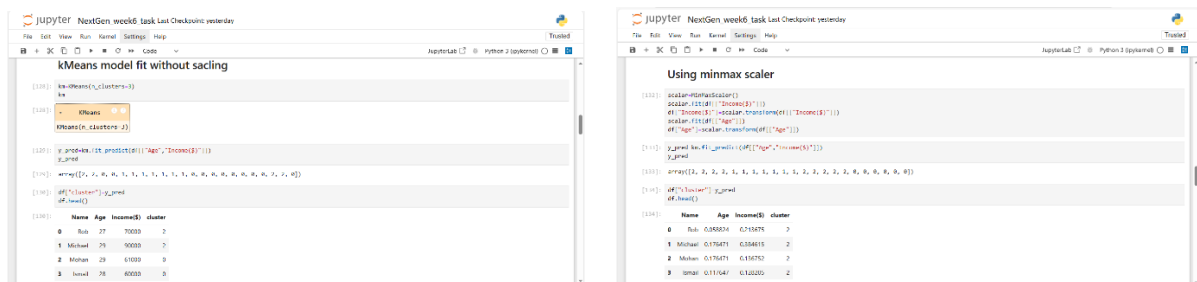
**Tasks Completed:**

- Learned about unsupervised learning and its applications.

- Understood the concept of clustering and the K-Means algorithm.

- Implemented K-Means clustering using Scikit-Learn's KMeans class.

- Learned about dimensionality reduction and Principal Component Analysis (PCA).

- Performed PCA using Scikit-Learn's PCA class to reduce the dimensionality of a dataset.

- Visualized clusters and the results of PCA.

- Evaluated the performance of K-Means clustering.

**Tools Used:**

- Scikit-Learn

- Pandas

- NumPy

- Matplotlib

- Jupyter Notebook

**Code Snippets / Design Screenshots**
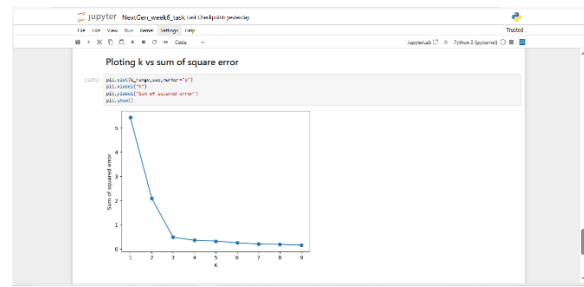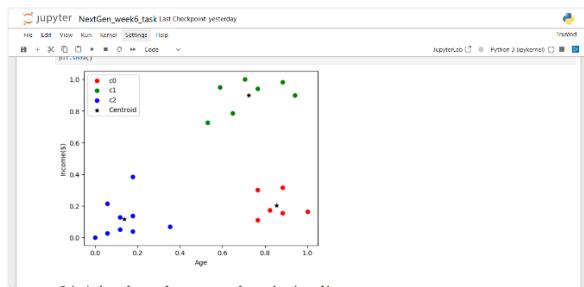
**Example 1: K-Means Clustering Implementation**

## Example 2: PCA Implementation and Results









## Challenges Faced

1. **Understanding the mathematical concepts behind PCA:**

   o Initially found it challenging to grasp the concepts of eigenvectors, eigenvalues, and variance explained.

   o **Resolution:** Reviewed linear algebra concepts and PCA tutorials, and experimented with different datasets to visualize the effects of PCA.

2. **Choosing the optimal number of clusters for K-Means:**

   o Was unsure about how to determine the best value for 'k' in K-Means clustering.

   o **Resolution:** Implemented the elbow method to identify the optimal 'k' by plotting the within-cluster sum of squares for different values of 'k' and selecting the 'k' at the "elbow" point.

## Learning Outcome

- Developed a solid understanding of unsupervised learning techniques, including clustering and dimensionality reduction.

- Implemented K-Means clustering and PCA using Scikit-Learn.

- Gained experience in applying these techniques to real-world datasets.

- Improved ability to analyze and interpret the results of clustering and dimensionality reduction.

- Learned how to visualize clusters and the effects of PCA.

**Resources**

- Scikit-Learn Documentation: https://scikit-learn.org/stable/

- Clustering Guide

- PCA Guide