

Let's go through the code step-by-step to understand what each part does and why it's essential.

Step 1: Import Necessary Libraries

```
```python
import pandas as pd
import numpy as np
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
import matplotlib.pyplot as plt
import seaborn as sns
```
```

****Explanation:****

This imports essential libraries:

- ****Pandas**** and ****Numpy**** are for data handling.
- ****Scikit-Learn**** provides tools for machine learning tasks like clustering (K-means), scaling (StandardScaler), dimensionality reduction (PCA), and model evaluation (Random Forest and accuracy score).
- ****Matplotlib**** and ****Seaborn**** are used for data visualization.

Step 2: Load the Dataset

```
` `` `python  
data = pd.read_csv('/content/cleaned_data.csv') # Adjust path as needed  
` `` `
```

****Explanation:****

The dataset is loaded from a file, `cleaned_data.csv`, using Pandas. Adjust the path as required depending on where you've stored your file (e.g., Google Drive if used in Colab).

Step 3: Explore Data

```
` `` `python  
print("Data Info:")  
print(data.info())  
print("\nFirst few rows:")  
print(data.head())  
` `` `
```

****Explanation:****

This step prints out:

1. ****Data Info****: Provides details on each column (data type, non-null counts) to help understand the dataset structure.


```
scaler = StandardScaler()

scaled_features = scaler.fit_transform(numerical_features)

...
```

****Explanation:****

To prepare the data for clustering, we select only numerical columns related to user activity, sleep, and health metrics.

- ****Standardization****: K-means clustering performs better on standardized data, where each feature has a mean of 0 and a standard deviation of 1. This is achieved using `StandardScaler`. Standardizing ensures that all features contribute equally to the distance calculations in clustering.

Step 6: Apply K-means Clustering and Visualize with PCA

```
```python

kmeans = KMeans(n_clusters=4, random_state=42)

clusters = kmeans.fit_predict(scaled_features)

data['Cluster'] = clusters

...
```

**\*\*Explanation:\*\***

K-means clustering is used to identify groups within the data. We set `n_clusters=4` to explore if clusters align with the four moods (`Neutral`, `Stressed`, `Happy`, `Tired`), although this may not be an exact match.

- **\*\*Adding Cluster Labels\*\***: The resulting cluster labels are added to the dataset for easier analysis and visualization.

---

#### #### PCA for Visualization

```
```python
```

```
pca = PCA(n_components=2)
```

```
pca_components = pca.fit_transform(scaled_features)
```

```
```
```

**\*\*Explanation:\*\***

Since the dataset has multiple numerical features, plotting all of them would be challenging. **\*\*PCA (Principal Component Analysis)\*\*** reduces data to two dimensions, helping us visualize the clusters in 2D.

- **\*\*n\_components=2\*\*** limits PCA to the top 2 principal components, which capture the highest variance in the data.

---

#### #### Plotting the Clusters

```
```python
```

```
plt.figure(figsize=(10, 6))
```

```

sns.scatterplot(x=pca_components[:, 0], y=pca_components[:, 1], hue=data['Cluster'],
palette="viridis", s=50)

plt.title("K-means Clustering Visualization (PCA-Reduced to 2D)")

plt.xlabel("PCA Component 1")

plt.ylabel("PCA Component 2")

plt.legend(title="Cluster")

plt.show()

` ``

```

****Explanation:****

This creates a 2D scatter plot showing the clusters using the PCA components. Each color represents a different cluster, which helps visualize the distinct groups formed by K-means. PCA helps us understand clustering results even in high-dimensional data.

Step 7: Feature Engineering and Accuracy Evaluation

****Feature Subsets Creation****

```

` `` python

feature_subsets = {

    "Activity & Heart": ['steps', 'active_minutes', 'heart_rate_avg'],

    "Calories & Distance": ['calories_burned', 'distance_km'],

    "Sleep & Activity": ['sleep_hours', 'active_minutes', 'steps']

}

` ``

```

****Explanation:****

This defines three feature subsets to test which combinations of variables work best for predicting `mood`. ****Feature engineering**** like this can help identify the most informative features for classification.

****Training and Evaluating Each Subset****

```
` `` `python
```

```
accuracy_results = {}
```

```
for subset_name, features in feature_subsets.items():
```

```
    X_subset = X[features]
```

```
    # Split data into train and test sets
```

```
    X_train, X_test, y_train, y_test = train_test_split(X_subset, y, test_size=0.3,
random_state=42)
```

```
    # Train classifier
```

```
    clf = RandomForestClassifier(random_state=42)
```

```
    clf.fit(X_train, y_train)
```

```
    # Predict and calculate accuracy
```

```
    y_pred = clf.predict(X_test)
```

```
    accuracy = accuracy_score(y_test, y_pred)
```

```
    accuracy_results[subset_name] = accuracy
```

```
print("\nAccuracy Results for Feature Subsets:")  
  
print(accuracy_results)  
` ``
```

****Explanation:****

1. ****Splitting Data****: Each feature subset is split into training and testing sets using ``train_test_split``.
2. ****Random Forest Classifier****: We use a ****Random Forest Classifier**** (a robust, easy-to-tune classifier) to predict ``mood`` based on each feature subset.
3. ****Accuracy Calculation****: After training, we predict the ``mood`` for the test set and calculate the accuracy score for each subset. This helps identify which features are most predictive of mood.

Step 8: Apply PCA for Dimensionality Reduction

```
` `` `python  
  
pca_full = PCA(n_components=2) # Change n_components to desired number for analysis  
pca_data = pca_full.fit_transform(scaled_features)  
  
# Explained variance by each principal component  
print("\nExplained Variance Ratio by each PCA component:")  
print(pca_full.explained_variance_ratio_)  
` `` `
```

****Explanation:****

PCA is applied to reduce the dimensions of the full dataset, which simplifies data visualization and can improve model efficiency for high-dimensional data.

- **Explained Variance Ratio**: This shows how much variance each principal component captures, helping assess how well PCA represents the data. The more variance a component explains, the more useful it is for representing the dataset.

Summary

This code:

1. **Loads and explores the data** to understand its structure and content.
2. **Standardizes and clusters** the data using K-means, visualizing results with PCA for easier interpretation.
3. **Creates feature subsets** to analyze which features best predict the `mood` variable, evaluating accuracy with a Random Forest classifier.
4. **Applies PCA for dimensionality reduction**, making the dataset more manageable and improving visualization.

These steps provide a complete analysis pipeline for exploring, clustering, engineering features, and evaluating the dataset. Let me know if there's anything specific you'd like to dive deeper into!

Principal Component Analysis (PCA) is a technique used in machine learning and statistics for dimensionality reduction. It transforms a dataset with many correlated variables into a set of linearly uncorrelated variables called **principal components**. These components represent the data's main patterns, preserving most of the information while reducing its complexity.

1. What is PCA?

PCA finds the directions (principal components) in which the data varies the most and projects the data onto these directions. Here's how it works, in simplified steps:

1. **Calculate the Mean**: PCA starts by centering the data by subtracting the mean of each feature.
2. **Covariance Matrix**: It then calculates the covariance matrix to understand the relationships between variables.
3. **Eigenvalues and Eigenvectors**: PCA finds eigenvalues and eigenvectors of the covariance matrix, where eigenvectors represent directions of maximum variance.
4. **Projection**: The data is projected onto the top k eigenvectors (principal components) with the largest eigenvalues, preserving most variance while reducing dimensions.

2. Why Use PCA?

PCA is especially useful in the following scenarios:

- **Dimensionality Reduction**: In datasets with a large number of features, many of which may be redundant, PCA reduces the number of features by creating a smaller set of new variables. This makes data easier to visualize, understand, and manage.
- **Feature Extraction**: By capturing the most relevant patterns in the data, PCA can help reduce noise and emphasize underlying trends. This can be useful in complex datasets with highly correlated variables.

- **Improving Model Performance**: For machine learning models, fewer input features can mean faster training, improved generalization, and reduced risk of overfitting.
- **Data Visualization**: PCA can reduce a high-dimensional dataset to 2D or 3D, making it easier to visualize complex relationships and clusters in the data.

3. Advantages of Using PCA

- **Simplicity**: Reduces complexity by converting a large set of variables into a smaller, more manageable set, often with minimal loss of information.
- **Improved Efficiency**: Models trained on fewer features can often run faster and require less computational power, which is critical for large datasets.
- **Reduced Overfitting**: By focusing on the components with the most variance, PCA helps minimize noise and irrelevant details, making models more generalizable.
- **Enhanced Interpretability**: While PCA creates new components that may not have a straightforward interpretation, it can often clarify the main directions of variance in data, making the dataset's structure clearer.

Disadvantages to Be Aware Of

- **Loss of Interpretability**: The new features (principal components) are combinations of original features and may not be easy to interpret.
- **Sensitivity to Scaling**: PCA is sensitive to the scale of the data, so features need to be standardized.

- **Linear Assumption**: PCA assumes linear relationships between variables, so it may not capture nonlinear patterns well.

In summary, PCA is a powerful tool for reducing the dimensionality of complex datasets, enabling efficient analysis and visualization while maintaining most of the original data's structure and patterns.