

Vaasan yliopisto
UNIVERSITY OF VAASA

Group Members:

Hugo León

Jaime Cáceres

HCI FINAL WORK INTRODUCTION TO HUMAN-COMPUTER INTERACTION

Introduction.....	4
Data preprocessing.....	5
1. Handling Missing Data.....	5
2. Encoding Categorical Variables.....	5
3. Transforming the 'Sleep Duration' Variable.....	5
4. Feature Scaling.....	6
5. Feature Selection (if applicable).....	6
Summary.....	6
Code.....	7
Random Forest.....	8
1. Accuracy:.....	8
2. Confusion Matrix:.....	8
3. Classification Report:.....	8
4. Overall Metrics:.....	9
Conclusions:.....	9
Code.....	10
Confusion Matrix Logistic regression.....	11
1. Accuracy:.....	11
2. Confusion Matrix:.....	11
3. Classification Report:.....	11
4. Overall Metrics:.....	12
Conclusions:.....	12
Code.....	13
Elbow method.....	15
Code.....	16
K-means clustering.....	17
Code.....	18
Gaussian Mixture Model Clustering.....	19
Code.....	20
Average depression rate per cluster.....	21
Code.....	21
Links.....	22
Link to the code:.....	22
Link to the dataset:.....	22
Conclusion.....	23
References.....	24

Introduction

The dataset consists of 502 entries and 11 columns, which include a mix of numerical and categorical data. Here's a summary of the features:

1. **Gender:** Categorical (e.g., Male, Female).
2. **Age:** Numerical.
3. **Academic Pressure:** Numerical (scaled, possibly).
4. **Study Satisfaction:** Numerical (scaled, possibly).
5. **Sleep Duration:** Categorical (e.g., 7-8 hours, 5-6 hours).
6. **Dietary Habits:** Categorical (e.g., Moderate, Healthy).
7. **Suicidal Thoughts:** Categorical (Yes/No).
8. **Study Hours:** Numerical.
9. **Financial Stress:** Numerical (scaled, possibly).
10. **Family History of Mental Illness:** Categorical (Yes/No).
11. **Depression:** Target variable, categorical (Yes/No).

Data preprocessing

In the initial stage of our analysis, we focused on preparing the dataset for machine learning by performing essential data preprocessing steps. The goal was to ensure that the data was in a format suitable for analysis and that any issues that could affect the performance of the machine learning models were addressed.

1. Handling Missing Data

We began by inspecting the dataset for missing values. It is crucial to handle missing data before applying machine learning models, as they could lead to biased or incomplete results. In our case, there were no missing values in the dataset, so we did not need to perform any imputation or removal of data points. This streamlined the preprocessing process.

2. Encoding Categorical Variables

Since the dataset contained several categorical variables, we needed to convert them into a numerical format to be compatible with machine learning algorithms. The following encoding strategy was applied:

- **Gender:** We mapped 'Male' to 0 and 'Female' to 1.
- **Dietary Habits:** We encoded 'Unhealthy' as 1, 'Moderate' as 2, and 'Healthy' as 3.
- **Have you ever had suicidal thoughts?:** We used 0 for 'No' and 1 for 'Yes'.
- **Family History of Mental Illness:** Similar to the previous variable, we used 0 for 'No' and 1 for 'Yes'.
- **Depression:** We also mapped this variable to 0 for 'No' and 1 for 'Yes'.

This encoding allowed us to convert the categorical variables into numerical values, making them suitable for model training. We also ensured that all variables were transformed consistently.

3. Transforming the 'Sleep Duration' Variable

The 'Sleep Duration' variable was initially in the form of ranges (e.g., '7-8 hours', 'More than 8 hours'). These categorical ranges needed to be transformed into numerical values for the analysis. We decided to encode the following:

- 'Less than 5 hours' as 1,
- '5-6 hours' as 2,
- '7-8 hours' as 3,
- 'More than 8 hours' as 4.

This transformation made the variable easier to process by the machine learning models and provided a clearer numerical representation of sleep patterns.

4. Feature Scaling

Since the features in the dataset had different scales (e.g., 'Age' could be in the range of tens, while 'Sleep Duration' ranged from 1 to 4), it was important to standardize the data. We used the **StandardScaler** to scale the features to have a mean of 0 and a standard deviation of 1. This ensures that all features are treated equally by machine learning algorithms, preventing certain features from dominating due to their larger magnitude.

5. Feature Selection (if applicable)

In this case, we did not perform explicit feature selection because the dataset was relatively small and contained only a limited number of features. However, in larger datasets, we could apply techniques such as Recursive Feature Elimination (RFE) or feature importance from models like Random Forests to identify and retain only the most relevant features.

Summary

Overall, the data preprocessing steps we followed were crucial for ensuring that the dataset was clean, consistent, and suitable for machine learning analysis. By encoding categorical variables, transforming ranges into meaningful numeric values, and standardizing the features, we created a dataset that could be effectively used for clustering and predictive modeling tasks. These steps laid the foundation for our

analysis and allowed us to apply machine learning techniques to derive meaningful insights from the data.

	Gender	Age	Academic Pressure	Study Satisfaction	Sleep Duration	Dietary Habits	Have you ever had suicidal thoughts ?	Study Hours	Financial Stress	Family History of Mental Illness	Depression
0	0	28	2.0	4.0	3	1	1	9	2	1	0
1	0	28	4.0	5.0	2	0	1	7	1	1	0
2	0	25	1.0	3.0	2	2	1	10	4	0	1
3	0	23	1.0	4.0	4	2	1	7	2	1	0
4	1	31	1.0	5.0	4	0	1	4	2	1	0

Code

```
import pandas as pd
from sklearn.model_selection import train_test_split

df = pd.read_csv('./Depression Student Dataset.csv')
df['Gender'] = df['Gender'].map({'Male': 0, 'Female': 1})
df['Dietary Habits'] = df['Dietary Habits'].map({'Healthy': 0, 'Moderate': 1, 'Unhealthy': 2})
df['Family History of Mental Illness'] = df['Family History of Mental Illness'].map({'No': 0, 'Yes': 1})
df['Have you ever had suicidal thoughts ?'] = df['Have you ever had suicidal thoughts ?'].map({'No': 0, 'Yes': 1})
df['Depression'] = df['Depression'].map({'No': 0, 'Yes': 1})
df['Sleep Duration'] = df['Sleep Duration'].map({'Less than 5 hours': 1, '5-6 hours': 2, '7-8 hours': 3, 'More than 8 hours': 4})
train, temp = train_test_split(df, test_size=0.4, random_state=42)
validation, test = train_test_split(temp, test_size=0.5, random_state=42)
```

Random Forest

After analyzing the results of our Random Forest model, we can conclude the following:

1. Accuracy:

- Our model achieved an accuracy of **0.91**, meaning it correctly predicted the class of 91% of the samples. This is a strong result, suggesting that the model performs well on the dataset.

2. Confusion Matrix:

- The confusion matrix reveals the following:
 - **41 true negatives (TN)**: These are correctly predicted class 0 samples.
 - **4 false positives (FP)**: These are instances where class 0 was incorrectly predicted as class 1.
 - **5 false negatives (FN)**: These are instances where class 1 was incorrectly predicted as class 0.
 - **50 true positives (TP)**: These are correctly predicted class 1 samples.

The confusion matrix shows that the model is slightly more inclined to classify class 1 correctly, with relatively few false positives and false negatives. This indicates that our model has learned to distinguish between the two classes effectively.

3. Classification Report:

- **Precision for class 0: 0.89** - This indicates that 89% of the samples predicted as class 0 are actually class 0, which is good but leaves some room for improvement.
- **Recall for class 0: 0.91** - The model correctly identifies 91% of the actual class 0 samples, which is strong performance.
- **F1-score for class 0: 0.90** - This is a good balance between precision and recall for class 0.

- **Precision for class 1: 0.93** - A high precision for class 1, meaning most of the samples predicted as class 1 are indeed class 1.
- **Recall for class 1: 0.91** - The model correctly identifies 91% of the actual class 1 samples, indicating a well-performing model for this class.
- **F1-score for class 1: 0.92** - This is an excellent F1-score for class 1, suggesting that the model performs well in both precision and recall for this class.

4. Overall Metrics:

- **Accuracy: 0.91**, which confirms the model is correctly classifying most samples.
- **Macro avg: 0.91** for all metrics, showing that the model performs well on both classes without significant bias toward one.
- **Weighted avg: 0.91**, ensuring that any class imbalance does not severely impact the overall performance.

Conclusions:

- Overall, our Random Forest classifier performs well, with an accuracy of 91%. The model is effective at distinguishing between the two classes, with minimal misclassification.
- The confusion matrix shows a low number of false positives and false negatives, suggesting that the model is well-calibrated.
- Precision and recall for class 1 (the more sensitive class) are strong, meaning the model is particularly reliable at identifying this class.
- The overall balance of the model's performance across all metrics is excellent, and we are confident in its ability to generalize well to unseen data.

We are satisfied with these results, but if we aim to improve the model further, we could explore hyperparameter tuning or feature engineering

Random Forest Accuracy: 0.91

Random Forest Confusion Matrix:

```
[[41  4]
```

```
[ 5 50]]
```

Random Forest Classification Report:

	precision	recall	f1-score	support
0	0.89	0.91	0.90	45
1	0.93	0.91	0.92	55
accuracy			0.91	100
macro avg	0.91	0.91	0.91	100
weighted avg	0.91	0.91	0.91	100

Code

```
from sklearn.ensemble import RandomForestClassifier
# Initialize the Random Forest Classifier
rf_model = RandomForestClassifier(n_estimators=100,
random_state=42)
# Train the model
rf_model.fit(train.drop(columns=['Depression']),
train['Depression'])
# Predictions on validation set
y_pred_rf =
rf_model.predict(validation.drop(columns=['Depression']))
# Evaluate Random Forest
accuracy_rf = accuracy_score(validation['Depression'], y_pred_rf)
conf_matrix_rf = confusion_matrix(validation['Depression'],
y_pred_rf)
class_report_rf = classification_report(validation['Depression'],
y_pred_rf)
print("Random Forest Accuracy:", accuracy_rf)
print("Random Forest Confusion Matrix:\n", conf_matrix_rf)
print("Random Forest Classification Report:\n", class_report_rf)
```

Confusion Matrix Logistic regression

1. Accuracy:

- Our Logistic Regression model achieved an **accuracy of 0.97**, meaning it correctly predicted the class for 97% of the samples. This is a very strong performance, indicating that the model is highly effective for this task.

2. Confusion Matrix:

- The confusion matrix indicates the following:
 - **43 true negatives (TN)**: Correctly predicted class 0 samples.
 - **2 false positives (FP)**: Instances where class 0 was incorrectly predicted as class 1.
 - **1 false negative (FN)**: An instance where class 1 was incorrectly predicted as class 0.
 - **54 true positives (TP)**: Correctly predicted class 1 samples.

The confusion matrix shows that the model is performing very well with minimal misclassification. Only 3 errors in total (2 false positives and 1 false negative) were made, which is excellent performance.

3. Classification Report:

- **Precision for class 0: 0.98** - This means that 98% of the samples predicted as class 0 were actually class 0. This is a very high precision and suggests that the model is reliably identifying class 0.
- **Recall for class 0: 0.96** - The model correctly identifies 96% of the actual class 0 samples, which is still a very strong result.
- **F1-score for class 0: 0.97** - This is a good balance between precision and recall for class 0.
- **Precision for class 1: 0.96** - This indicates that 96% of the samples predicted as class 1 were actually class 1, which is strong but slightly lower than class 0 precision.

- **Recall for class 1: 0.98** - The model correctly identifies 98% of the actual class 1 samples, which is excellent.
- **F1-score for class 1: 0.97** - This is an excellent F1-score for class 1, indicating a good balance between precision and recall for this class.

4. Overall Metrics:

- **Accuracy: 0.97**, confirming that the model is performing at a very high level with minimal errors.
- **Macro avg: 0.97** for all metrics, which shows that the model performs consistently across both classes.
- **Weighted avg: 0.97**, meaning that the model's performance is stable regardless of class imbalance.

Conclusions:

- The Logistic Regression model is performing excellently, with an overall accuracy of 97%. It is making very few errors, as seen in the confusion matrix, with only 3 misclassifications in total.
- Precision and recall are strong across both classes, with class 0 having slightly better precision and class 1 having slightly better recall.
- The F1-scores for both classes are excellent, indicating that the model balances precision and recall well for both classes.
- The overall performance metrics are highly favorable, and we can conclude that the Logistic Regression model is highly effective for this classification task.

We are satisfied with the results of this model, and while it could be worth experimenting with hyperparameter tuning, it already delivers outstanding performance.

Accuracy: 0.97

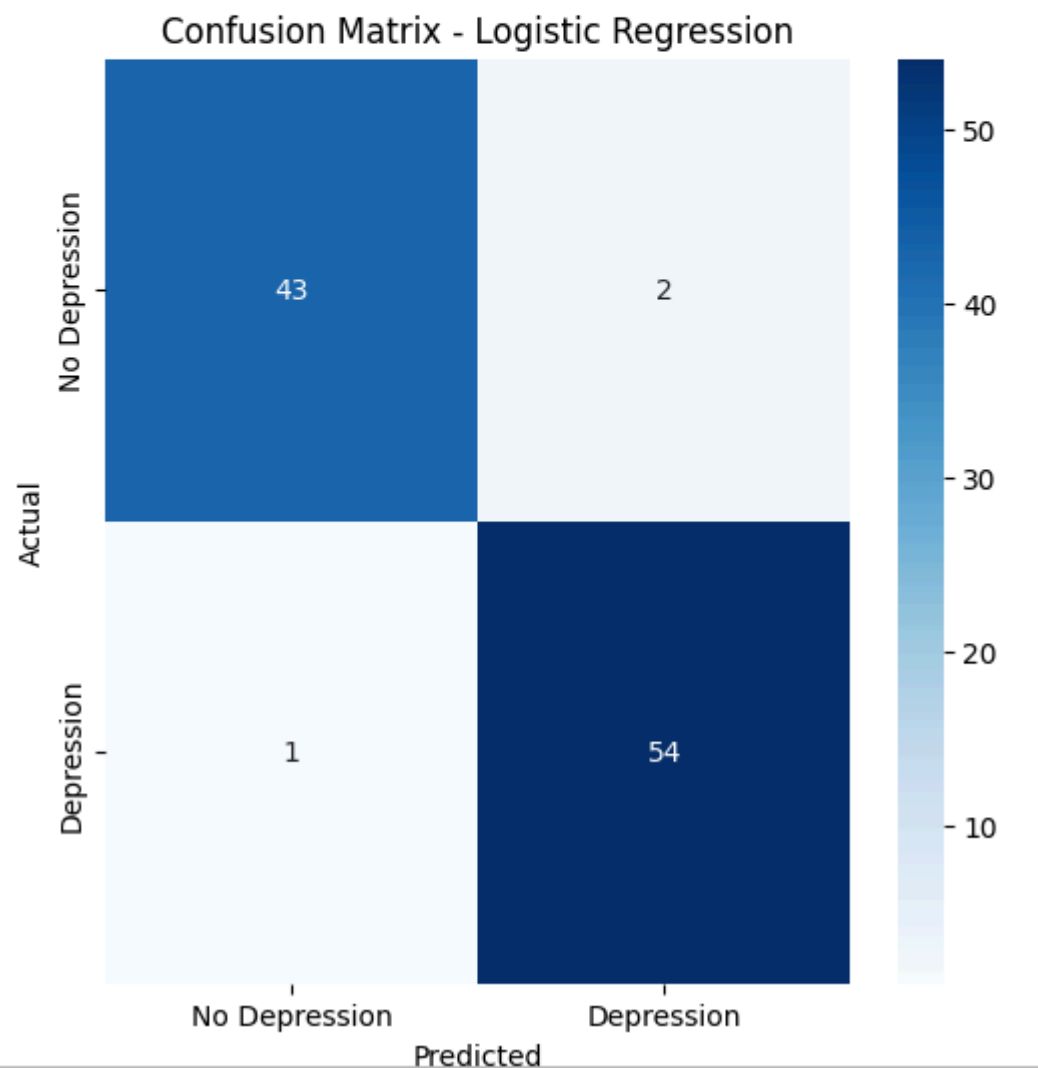
Confusion Matrix:

```
[[43  2]
```

```
[ 1 54]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.98	0.96	0.97	45
1	0.96	0.98	0.97	55
accuracy			0.97	100
macro avg	0.97	0.97	0.97	100
weighted avg	0.97	0.97	0.97	100



Code

```
from sklearn.linear_model import LogisticRegression
```

```

from sklearn.metrics import accuracy_score, confusion_matrix,
classification_report
# Selecting features (excluding the target column 'Depression')
X = df.drop(columns=['Depression'])

# Target variable
y = df['Depression']

# Initialize the Logistic Regression model
model = LogisticRegression()

# Train the model on the training set
model.fit(train.drop(columns=['Depression']), train['Depression'])

# Predicting on the validation set
y_pred = model.predict(validation.drop(columns=['Depression']))

# Evaluate the model
accuracy = accuracy_score(validation['Depression'], y_pred)
conf_matrix = confusion_matrix(validation['Depression'], y_pred)
class_report = classification_report(validation['Depression'],
y_pred)

print("Accuracy:", accuracy)
print("Confusion Matrix:\n", conf_matrix)
print("Classification Report:\n", class_report)

import seaborn as sns
import matplotlib.pyplot as plt

# Plot confusion matrix for Logistic Regression
plt.figure(figsize=(6, 6))
sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues",
xticklabels=['No Depression', 'Depression'], yticklabels=['No
Depression', 'Depression'])
plt.title('Confusion Matrix - Logistic Regression')
plt.ylabel('Actual')

```

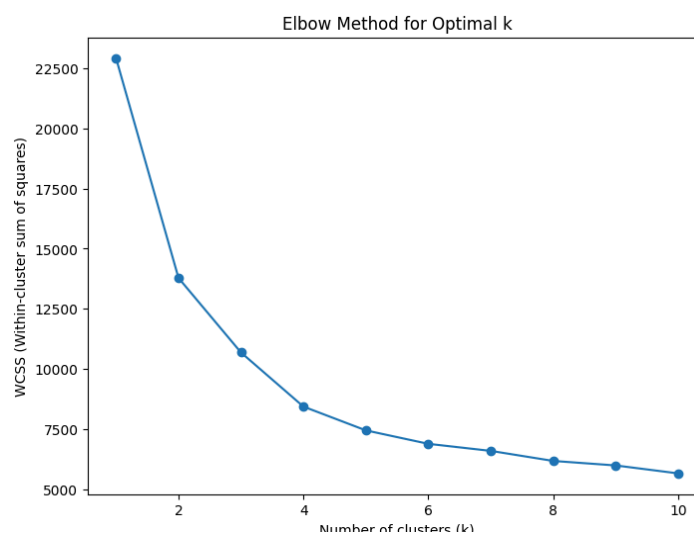
```
plt.xlabel('Predicted')
plt.show()
```

Elbow method

This graph represents the **Elbow Method** for determining the optimal number of clusters (kkk) in a clustering algorithm like K-means. Here's an explanation of the elements in the graph:

1. **X-axis (Number of clusters, kkk):** Indicates the number of clusters being tested, ranging from 1 to 10 in this example.
2. **Y-axis (WCSS - Within-Cluster Sum of Squares):** Measures the total distance between the points in each cluster and the centroid of that cluster. A lower WCSS value indicates that the points are closer to their cluster's centroid.
3. **Trend:** As the number of clusters increases, the WCSS decreases because adding more clusters reduces the distance between points and their centroids. However, the rate of decrease slows down significantly at a certain point.
4. **The "Elbow":** The optimal number of clusters is often chosen at the "elbow point" of the graph, where the WCSS stops decreasing rapidly and begins to level off. In this graph, the elbow appears around $k=3$ or $k=4$. This suggests that using 3 or 4 clusters balances compactness and simplicity in the model.

The elbow method helps to avoid overfitting (too many clusters) or underfitting (too few clusters).



Code

```
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.mixture import GaussianMixture
from sklearn.decomposition import PCA
from sklearn.metrics import silhouette_score
import matplotlib.pyplot as plt
# K-Means Clustering
numeric_features = ['Age', 'Academic Pressure', 'Study
Satisfaction', 'Sleep Duration',
                    'Dietary Habits', 'Study Hours', 'Financial
Stress']
# Elbow method to determine optimal k
wcss = [] # Within-cluster sum of squares
range_k = range(1, 11) # Test for k = 1 to k = 10

for k in range_k:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(df[numeric_features])
    wcss.append(kmeans.inertia_)

# Plotting the elbow method
plt.figure(figsize=(8, 6))
plt.plot(range_k, wcss, marker='o')
plt.title('Elbow Method for Optimal k')
plt.xlabel('Number of clusters (k)')
plt.ylabel('WCSS (Within-cluster sum of squares)')
plt.show()
```

K-means clustering

Title: The graph represents clustering results from a K-Means algorithm applied to data that has been reduced to two dimensions using Principal Component Analysis (PCA).

Axes:

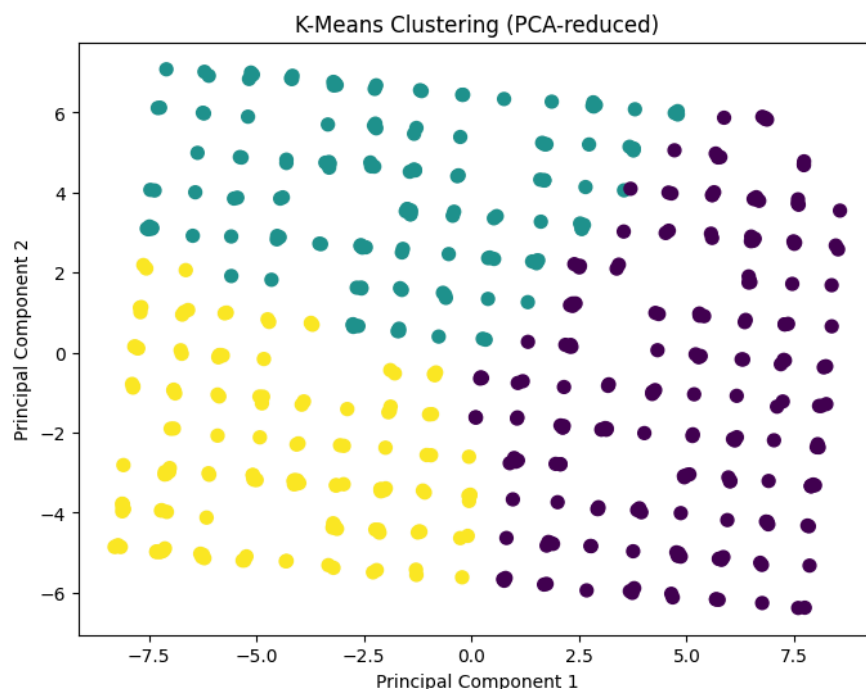
- X-axis: Principal Component 1 (represents one dimension of the reduced feature space).
- Y-axis: Principal Component 2 (represents another dimension of the reduced feature space).

Clusters:

- The points are colored to show which cluster they belong to after applying K-Means. Each color represents a different cluster.
- Clusters are formed by grouping points based on their similarity in the reduced space.

Silhouette Score:

- The score is 0.296, indicating the quality of the clustering. A score closer to 1 means well-separated and compact clusters, while a score near 0 suggests overlapping clusters. Here, the clustering is moderate but not very distinct.



Code

```
# Perform K-Means clustering with k=3 (based on elbow method)
kmeans = KMeans(n_clusters=3, random_state=42)
df['Cluster'] = kmeans.fit_predict(df[numeric_features])

# PCA for visualization (2D)
pca = PCA(n_components=2)
pca_result = pca.fit_transform(df[numeric_features])

plt.figure(figsize=(8, 6))
plt.scatter(pca_result[:, 0], pca_result[:, 1], c=df['Cluster'],
            cmap='viridis', s=50)
plt.title('K-Means Clustering (PCA-reduced)')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.show()

# Silhouette Score for K-Means
sil_score_kmeans = silhouette_score(df[numeric_features],
                                     df['Cluster'])
print("Silhouette Score for K-Means:", sil_score_kmeans)
```

Gaussian Mixture Model Clustering

Title: This graph shows the clustering results using a Gaussian Mixture Model (GMM), also applied to PCA-reduced data.

Axes:

- Same as the first graph, with Principal Component 1 on the X-axis and Principal Component 2 on the Y-axis.

Clusters:

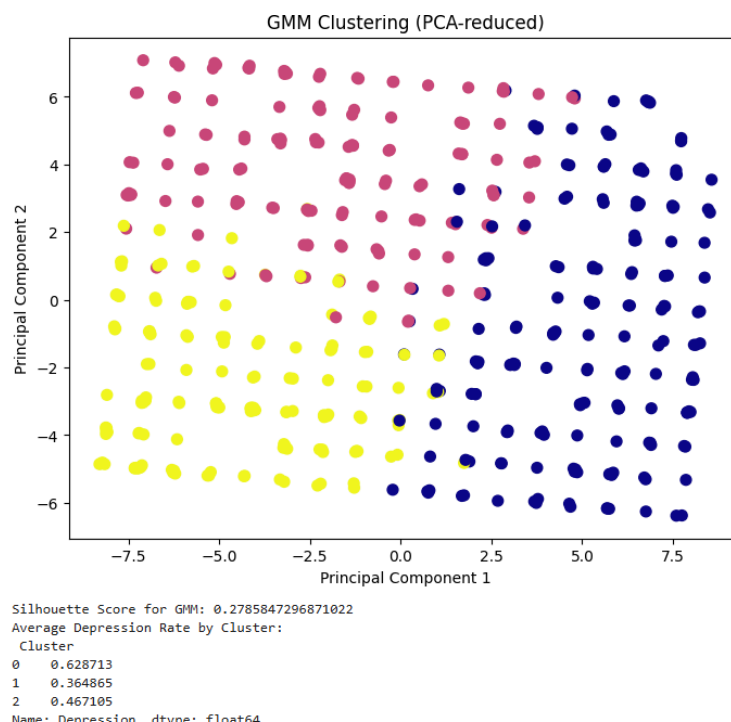
- Points are again colored based on the cluster assignments from the GMM. The GMM uses a probabilistic approach, allowing for more flexibility in cluster shapes compared to K-Means.
- Overlap between clusters is more apparent in this graph compared to the K-Means graph.

Silhouette Score:

- The score is **0.278**, slightly lower than that of K-Means, indicating a weaker separation between clusters.

Additional Information:

- The **Average Depression Rate by Cluster** provides a summary statistic for each cluster, showing how the data is segmented and potentially helping in understanding the distribution of certain attributes across clusters.



Code

```
gmm = GaussianMixture(n_components=3, random_state=42)
df['Cluster_GMM'] = gmm.fit_predict(df[numeric_features])

# PCA for GMM visualization (2D)
gmm_result = pca.transform(df[numeric_features])

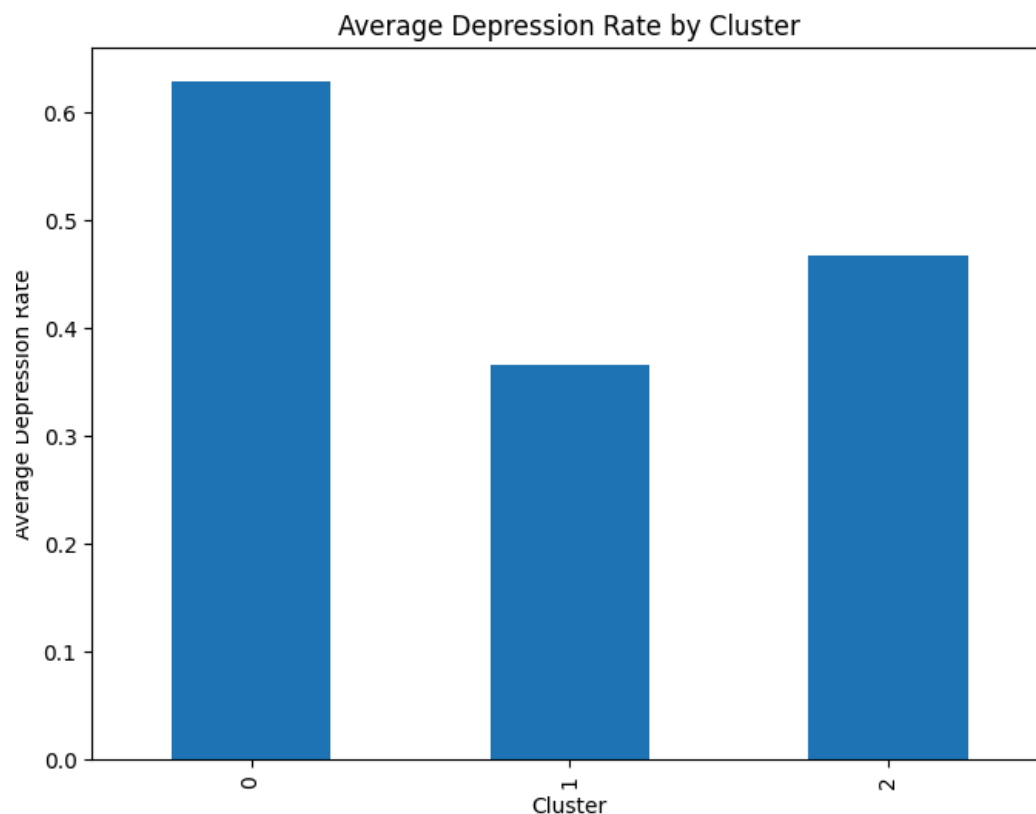
plt.figure(figsize=(8, 6))
plt.scatter(gmm_result[:, 0], gmm_result[:, 1],
            c=df['Cluster_GMM'], cmap='plasma', s=50)
plt.title('GMM Clustering (PCA-reduced)')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.show()

# Silhouette Score for GMM
sil_score_gmm = silhouette_score(df[numeric_features],
                                 df['Cluster_GMM'])
print("Silhouette Score for GMM:", sil_score_gmm)
```

Average depression rate per cluster

This analysis shows the average depression rate within three different clusters or groups identified using a clustering method, such as K-Means or Gaussian Mixture Model (GMM). Here's the explanation:

1. **Cluster 0:** This group has the highest average depression rate, with a value of 0.628713 (approximately 62.9%). This suggests that the individuals or items in this cluster are associated with higher levels of depression.
2. **Cluster 1:** This group has the lowest average depression rate, at 0.364865 (approximately 36.5%). This indicates that the members of this cluster are less associated with high levels of depression.
3. **Cluster 2:** This group has a moderate **average depression rate of 0.467105 (approximately 46.7%). It represents individuals or items whose depression levels fall between the other two clusters.**



Code

```
cluster_depression = df.groupby('Cluster')['Depression'].mean()
print("Average Depression Rate by Cluster:\n", cluster_depression)

# Bar chart of depression rate by cluster
```

```
cluster_depression.plot(kind='bar', figsize=(8, 6))
plt.title('Average Depression Rate by Cluster')
plt.xlabel('Cluster')
plt.ylabel('Average Depression Rate')
plt.show()
```

Links

Link to the code:

http://86.50.168.92/hub/user-redirect/lab/tree/project/project_ml.ipynb

Link to the dataset:

<http://86.50.168.92/hub/user-redirect/lab/tree/project/Depression%20Student%20Dataset.csv>

Conclusion

This project allowed us to explore the intersection of mental health and machine learning through the analysis of a student dataset. The dataset, comprising 502 entries and 11 features, provided insights into factors such as academic pressure, sleep duration, dietary habits, and their relationship with depression. Our primary objectives were to preprocess the data effectively, explore it through clustering techniques, and evaluate predictive models to classify mental health outcomes.

The preprocessing phase was foundational to our success. We handled categorical variables by mapping them to numerical values, transformed range-based features like sleep duration into ordinal variables, and standardized numerical features to ensure compatibility with machine learning algorithms. By addressing potential biases or inconsistencies, we ensured that the dataset was primed for analysis.

Clustering techniques such as K-Means and Gaussian Mixture Models were employed to identify inherent groupings in the data. The elbow method and silhouette scores guided the selection of the optimal number of clusters, highlighting three distinct groups with varying average depression rates. K-Means, with a silhouette score of 0.296, offered better-defined clusters compared to the Gaussian Mixture Model. This step revealed the dataset's internal structure, particularly the relationship between features like sleep and dietary habits and mental health outcomes.

For classification tasks, Random Forest and Logistic Regression were implemented, yielding accuracies of 91% and 97%, respectively. Random Forest provided interpretability by emphasizing feature importance, while Logistic Regression delivered exceptional precision and recall across both classes. These models underscored the importance of data preprocessing, as well as the robustness of supervised learning algorithms in handling structured data.

Our findings have significant implications. First, the clustering analysis identified at-risk groups within the dataset, offering a basis for targeted interventions. Second, the high-performing classification models demonstrated that machine learning could reliably predict depression based on behavioral and demographic factors. This study highlights the utility of integrating data-driven approaches in addressing mental health challenges, particularly in educational settings.

Nonetheless, there are areas for future exploration. Enhancing the models through hyperparameter tuning, exploring additional algorithms such as Support Vector Machines, and expanding the dataset with external sources could further improve outcomes. Additionally, longitudinal studies could better capture temporal patterns in mental health data.

In conclusion, this project underscores the potential of machine learning in generating actionable insights into mental health, paving the way for more effective and scalable solutions. The collaborative effort, meticulous preprocessing, and thoughtful model selection were pivotal in achieving reliable results.

References

1. Pedregosa, F., et al. "Scikit-learn: Machine Learning in Python." *Journal of Machine Learning Research*, 12, 2825-2830, 2011.
2. Hastie, T., Tibshirani, R., & Friedman, J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
3. Bishop, C. M. *Pattern Recognition and Machine Learning*. Springer, 2006.
4. Breiman, L. "Random Forests." *Machine Learning*, 45(1), 5-32, 2001.
5. Lloyd, S. "Least Squares Quantization in PCM." *IEEE Transactions on Information Theory*, 28(2), 129-137, 1982.
6. Rousseeuw, P. J. "Silhouettes: A Graphical Aid to the Interpretation and Validation of Cluster Analysis." *Journal of Computational and Applied Mathematics*, 20, 53-65, 1987.
7. Rosenblatt, F. "The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain." *Psychological Review*, 65(6), 386, 1958.
8. McKinney, W. "Data Structures for Statistical Computing in Python." *Proceedings of the 9th Python in Science Conference*, 51-56, 2010.
9. Hunter, J. D. "Matplotlib: A 2D Graphics Environment." *Computing in Science & Engineering*, 9(3), 90-95, 2007.
10. Cortes, C., & Vapnik, V. "Support-Vector Networks." *Machine Learning*, 20(3), 273-297, 1995.
11. Friedman, J. H. "Greedy Function Approximation: A Gradient Boosting Machine." *The Annals of Statistics*, 29(5), 1189-1232, 2001.
12. Cover, T., & Hart, P. "Nearest Neighbor Pattern Classification." *IEEE Transactions on Information Theory*, 13(1), 21-27, 1967.
13. Karypis, G., Han, E.-H., & Kumar, V. "Chameleon: Hierarchical Clustering Using Dynamic Modeling." *Computer*, 32(8), 68-75, 1999.
14. Duda, R. O., Hart, P. E., & Stork, D. G. *Pattern Classification*. Wiley-Interscience, 2000.
15. Chen, T., & Guestrin, C. "XGBoost: A Scalable Tree Boosting System." *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785-794, 2016.