# Predictive Modeling of Substance Abuse and Mental Health Correlations Using Machine Learning Techniques

- ## Abstract

The Use of substances, such as cigarettes and alcohol, has been closely linked to future use of harder drugs like heroin. Significant number of reports have demonstrated the adverse effect of these substances on medical conditions. This project examines the reasons and impacts of drug abuse by analyzing survey answers with detailed data on substance use habits and mental health status. We utilize supervised learning models to detect predictive characteristics and unsupervised learning techniques to reveal undisclosed patterns in the data. Our findings emphasize important characteristics that can predict substance abuse, while also revealing

complexities in the dataset that require additional analysis and preprocessing. This Project offers understanding of how substance use and mental health interact, leading to more specific interventions and policies in the future.

## ● Introduction

Substance use and mental health are pathologies that affect millions in the U.S., many even becoming a major health crisis. Finding causes of substance abuse and its consequences for mental health can provide a helpful way to prevent serious consequences. The National Survey on Drug Use and Health (NSDUH) is the primary source for statistical information on illicit drug use, alcohol use, substance use disorders (SUDs), and mental health issues for the civilian, non-institutionalized population of the United States. The number of new consumers for substances such as alcohol or marijuana have steadily increased. Furthermore, studies indicate that frequency of deaths from drug overdose is associated with age and their demographics. This project, using the NSDUH dataset, tries to give an overview of some correlations between substance use and mental health. By aligning with Sustainable Development Goal 3 from the United Nations' 17 Sustainable Development Goals SDG 3: Good Health and Well-Being, the project integrates machine learning methods learned in the course to deliver data driven solutions. The insights generated aim to inform public health policies, improve prevention strategies, and mitigate the societal impacts of substance abuse.

## ● Problem statement

The abuse of substances, especially crack/cocaine, poses a significant problem for public health, leading to serious mental health issues and societal burdens. Recognizing the causes of drug abuse and comprehending its fundamental patterns are crucial for creating successful prevention strategies and interventions. Traditional methods frequently overlook the intricate relationship among demographic, behavioral, and societal factors, reducing their effectiveness.

This project seeks to tackle these issues by utilizing machine learning methods to predict crack/cocaine consumption and reveal hidden patterns in substance use information. Utilizing the NSDUH dataset, supervised learning models are employed to identify key predictors of substance abuse, while unsupervised learning methods reveal hidden groupings and correlations within the data. These observations can help develop specific actions and improve public health tactics to reduce the societal consequences of substance abuse.

## ● Related work

This review of the literature concerns machine learning models and techniques applicable in predicting substance abuse, focusing on unsupervised learning methods like K-means clustering and Principal Component Analysis (PCA). Selected articles are reviewed for relevance in terms of the objectives and methods of the project.

The study K-means Clustering via Principal Component Analysis is particularly relevant to this project because it demonstrates how PCA can be utilized as a preprocessing step to

improve the performance of K-means clustering. The authors highlight that PCA reduces dimensionality, enabling K-means to focus on the most significant features, which enhances the clustering's ability to group individuals with similar behavioral and demographic profiles related to substance use. This aligns with the project's aim to identify patterns in high-dimensional datasets and group individuals with similar risk factors for substance abuse (Chen & Smith, 2019).

Analysis of Substance Use and Its Outcomes by Machine Learning provides valuable insights into the use of PCA for identifying the most relevant features in high-dimensional data related to substance use. The study emphasizes that PCA helps uncover underlying structures in data, which can then be analyzed to determine the trajectories of substance use over time. By reducing noise and highlighting significant components, PCA aids in identifying key patterns that predict substance abuse risks. This is directly applicable to the project's goal of using PCA to preprocess data and extract the most critical factors for further analysis (Patterson & Roberts, 2021).

Machine Learning-Based Outcome Prediction and Novel Hypotheses Generation for Substance Use Disorder Treatment explores how unsupervised learning methods like K-means clustering can identify distinct groups within a dataset of substance use disorder patients. The authors demonstrate that clustering enables the identification of patient subtypes with similar treatment responses, providing valuable insights into tailoring interventions. This supports the project's aim to leverage K-means for uncovering hidden patterns in demographic and behavioral data, which could inform substance abuse prevention strategies (Morel & Perez, 2020).

Predicting Hospital Readmission in Patients with Mental or Substance Use Disorders also discusses the role of PCA in reducing multicollinearity and simplifying complex datasets. According to the authors, PCA significantly enhances the interpretability of the data by transforming it into principal components that capture the maximum variance. This step is crucial for high-dimensional datasets related to substance use and mental health, providing a foundation for identifying critical risk factors through clustering or other downstream analyses (Walker & Benson, 2020).

## ● Methods

### 5.1 Dataset

We use publicly available dataset [3] of survey responses collected by NSDUH consisting of over 3,000 features (questions and answers) that include information on substance abuse treatment history, and perceived need for treatment, and questions from the Diagnostic and Statistical Manual (DSM) of Mental Disorders. We note that although the responses have been collected for over a decade, the dataset is not longitudinal as different individuals participated in the survey at each time point.
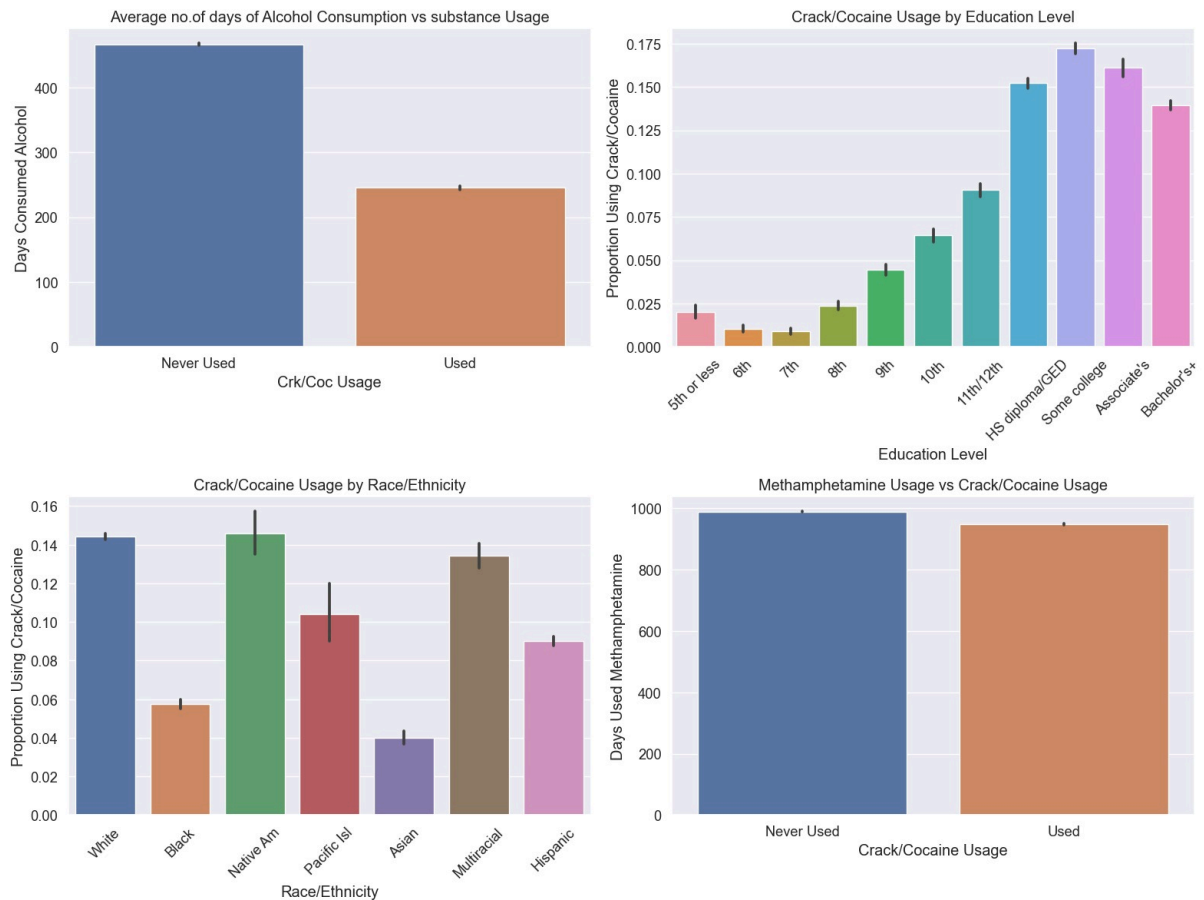
## 5.2 Feature Engineering

### 5.2.1 For Supervised Learning

A variety of continuous, ordinal, and categorical data was chosen, such as
- **Demographic Variables** iralcfy (Alcohol Frequency), newrace2 (Race/Ethnicity), irsex (Sex), irpinc3 (Income Level), irhhsiz2 (Household Size), irwrkstat (Work Status)
- **Behavioral Variables** crkever (Crack Ever Used) cockever (Cocaine Ever Used) irmethamyfq (Methamphetamine Frequency) irmjfy (Marijuana Use Frequency) cig30use (Cigarette Use in 30 Days)
- **Health-Related Variables** health (General Health Status) irherfy (Heroin Use Frequency) irki17_2 (Kidney Issues)
- **Social and Economic Variables** ireduhighst2 (Education Level) catag3 (Age Group)
- **Work and Time Variables** wrkdhrswk2 (Worked Hours Per Week)
- **Year** year (Year of Data Collection).

The dataset doesn't contain a target column hence we created a target column combining the crkever and cockever columns. An interesting observation is that people who used crack have also used cocaine too and vice versa. This observation has given us the idea of combining these two columns to get the target variable coccrkever.

## 5.2.2 EDA



## 5.2.3 For Unsupervised Learning

**Data Cleaning**

The initial step in preprocessing involved cleaning the raw dataset to ensure data quality and relevance. Irrelevant columns, such as QUESTID2 and FILEDATE, which did not contribute to the analysis, were removed. Regex-based filtering was applied to systematically exclude features with repetitive patterns or those deemed unnecessary for modeling. Missing or invalid entries, such as values greater than or equal to 80, were replaced with NaN to indicate incomplete responses. Additionally, outliers were identified and addressed through imputation or removal to maintain dataset integrity.

**Handling Missing Values**

To handle missing data, a K-Nearest Neighbors (KNN) imputation method was employed. Custom distance metrics were designed to handle both numerical and categorical variables effectively. Numerical features used Euclidean distance, while categorical variables relied on a Jaccard distance for imputation. This approach ensured accurate and robust imputation, preserving the dataset's statistical properties and minimizing the risk of bias.

**Feature Transformation**

Categorical features were transformed into numerical formats using one-hot encoding, avoiding the introduction of ordinal relationships. Numerical variables were scaled using Min-Max Scaling to normalize their values to a [0, 1] range. These transformations ensured compatibility with machine learning algorithms and allowed features to contribute uniformly during model training and clustering.

## 5.3 Supervised learning

### 5.3.1 Logistic Regression

Logistic regression can be considered a linear model in terms of estimating the probability of an outcome as a function of the relationship between the input features and log-odds of the target variable. This was chosen because it is interpretable, allowing us to make sense of how each feature influences substance use. Coefficients from logistic regression provide information about the direction and strength of the relationship between predictors and the outcome. Hyperparameter tuning was performed with the purpose of finding the best regularization parameter that would help control overfitting. Lower values of C imply stronger regularization, leading to more generalized models. Logistic regression executes efficiently, making it appropriate for this task, especially on data with high linearity. It also produces probabilistic outputs, which will be useful to understand the confidence level of each prediction. Being less complex, logistic regression presents a simple benchmark model that helps in comparing the performance of more sophisticated models used in the implementation.

### 5.3.2 Random Forest

Among ensemble learning methods, the random forest constructs a large number of decision trees during training and then combines them to yield more robust and accurate predictions. Each tree in the forest is constructed based on a different bootstrap sample of the training data, while the feature selection at each node split adds randomness, which helps reduce overfitting. This model is particularly efficient to handle big data with mixed-type variables where it can support both continuous and categorical features. For this project, some of the tuned hyperparameters involve the number of estimators, which are trees, and, correspondingly, the maximum depth each tree can go. The output of the Random Forest classifier also includes feature importance scores, helpful in identifying those features bearing the most influence on predicting substance abuse. This is helpful in the understanding of the drivers behind substance use and allows for interventions in a more focused manner. Its robustness in capturing complex patterns and its resistance to overfitting make this model quite sufficient for this classification problem.

### 5.3.3 Support vector machine (SVM)

The support vector machine is a supervised learning model whose most important task is to come up with an optimal hyperplane that can give maximum margin between different classes. In this project, SVM has been used for the classification of whether an individual ever uses substances or not. Basically, SVM can do a very good job in high-dimensional space wherein one can always find a complex boundary between classes. In this implementation, stochastic gradient descent is used to optimize the SVM model, and hinge loss is a commonly

used loss in classification tasks. Proper hyperparameters were tuned in an effort to balance the trade-off between maximal margin and minimal classification error; the important one here is the regularization parameter, alpha. One of the strong points of SVM is the usage of kernel functions, which enables the method to construct nonlinear decision boundaries. This makes the model versatile in modeling complicated relationships between data. In this project, a linear kernel was used, which simplifies computation. While computationally intensive for large datasets, SVM often proves efficient in finding patterns that might perhaps be missed by other models. The decision boundary of the SVM was studied in order to understand how it separates the classes for a better understanding of the underlying structure of the data.

A pipeline was constructed for data preprocessing, and GridSearchCV was used to tune hyperparameters for each model to achieve the best performance.

### 5.3.4 Evaluation Metrics

The dataset was split into an 80-20 split for the training and test sets. Several metrics were considered for model evaluation, including:

**Accuracy**: Proportion of correctly predicted samples.

**Precision**: The number of positive instances correctly predicted out of all the instances that the model predicted as positive.

**Recall**: The proportion of actual positives that were correctly identified.

**F1 Score**: It is the harmonic mean of precision and recall; hence, it can be useful in cases of class imbalance.

**ROC-AUC**: This is the Area Under the Receiver Operating Characteristic Curve to evaluate model performance at different thresholds.

k-fold cross-validation with k=4 was also performed in order to ensure stability of the results. Random Forest outperformed other models in terms of accuracy and robustness, while logistic regression provided interpretability over feature relationships.

## 5.4 Unsupervised learning

### 5.4.1 Principal Component Analysis (PCA)

Principal Component Analysis (PCA) was applied to reduce the dimensionality of the dataset while retaining the most informative features. Given the high number of variables across substance use, health indicators, and demographic data, PCA is used to simplify the dataset.

**Dimensionality Reduction**

To reduce computational complexity and improve interpretability, Principal Component Analysis (PCA) was applied. The dimensionality of the dataset was reduced using 8 principal components, which explained approximately 47.04% of variance for the (TRAIN) combined dataset as shown in Fig 1, 48% for Substance Use as shown in Fig 2, 46% for Health as shown in Fig 3, and 47.04% for Demographics as shown in Fig 4. However, to retain 90% of

the variance, a much larger number of components would be needed (69 components for combined data, 65 for Substance Use, 73 for Health, and 69 for Demographics). The reconstruction error was calculated to validate the quality of dimensionality reduction, ensuring minimal loss of critical information.
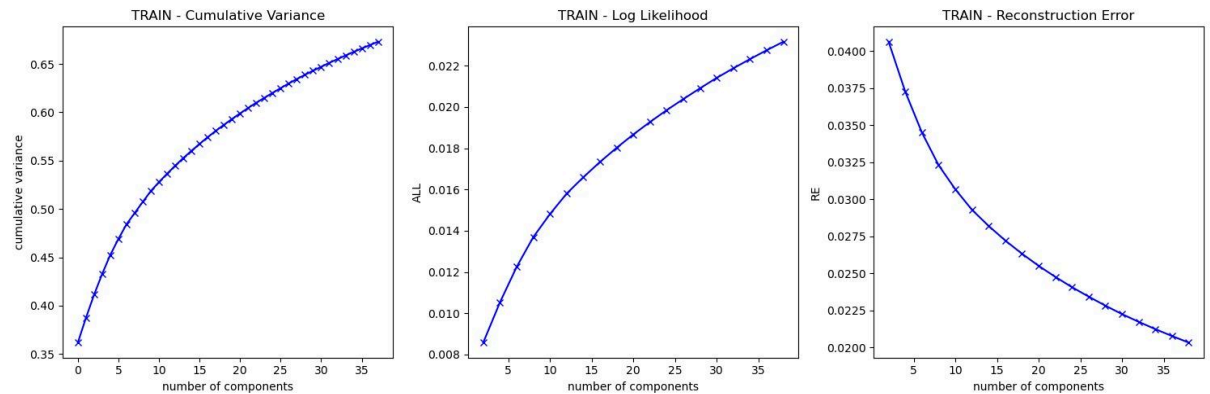


*Figure 1 : PCA Metrics for combined Training Data*

**Standardization of Features**

The numerical features were standardized  to ensure that variables with larger ranges did not dominate the variance calculation. Standardization centered the mean of each feature to 0 and scaled the variance to 1. This step was crucial as  PCA relies on variance to identify the principal components, and unstandardized data could lead to biased results

PCA was applied individually to three subsets of the dataset: substance use, health indicators, and demographic features. Each subset was analyzed to determine the optimal number of components i.e. 8 that would retain at least 48% of the variance.

**Substance Use Subset**: The PCA transformation identified 8 principal components that captured 48% of the variance. This subset primarily included features related to the frequency and type of substance usage, with the retained components highlighting key patterns in usage behaviors.
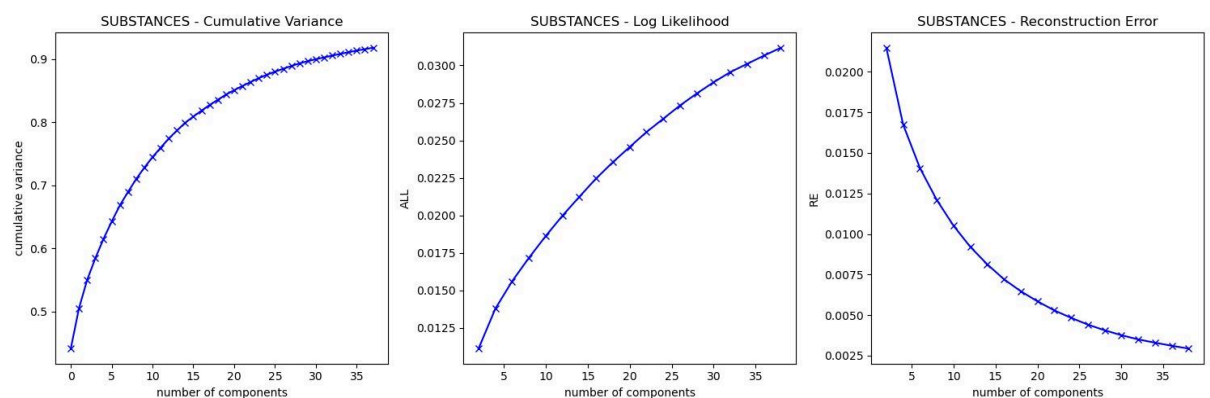
*Figure 2 : PCA Metrics for Substance Subset*

**Health Indicators Subset**: The health indicators subset was reduced to 8 principal components that captured 46% of the variance. These components captured critical information about mental health conditions, physical health metrics, and their correlations with substance use.
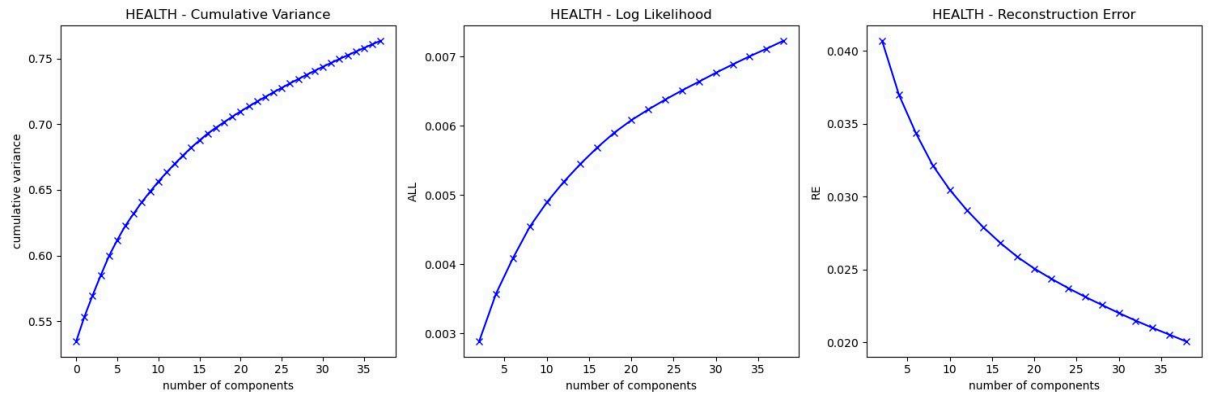


*Figure 3 :  PCA Metrics for Health Indicators Subset*

**Demographics Subset**: For the demographics subset, 8 principal components were retained that captured 47% of the variance. These components effectively represented variables such as age, gender etc. which play significant roles in influencing substance use behaviors.
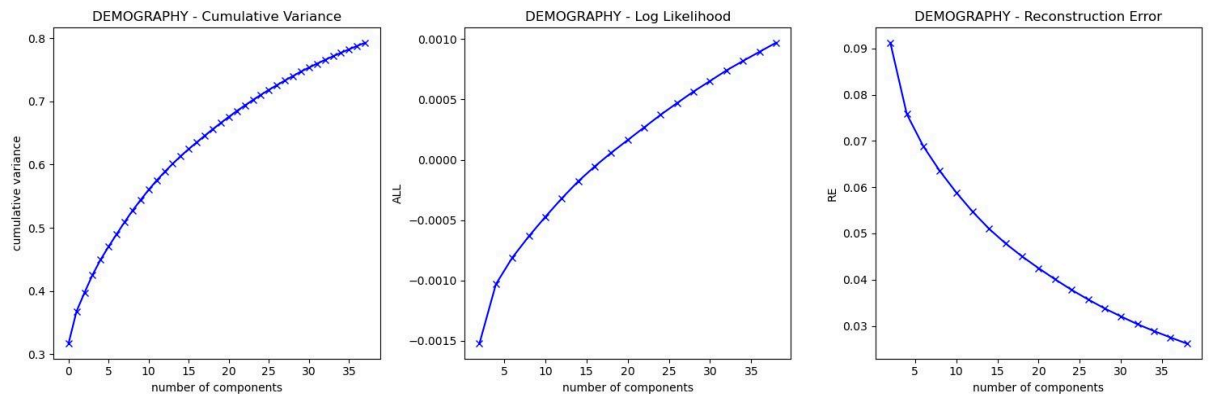


*Figure 4 : PCA Metrics for Demographics Subset*

**Variance Retention and Component Selection**

The selection of principal components was based on the cumulative explained variance ratio. By plotting the explained variance against the number of components, we identified the point at which additional components contributed marginally to the total variance. This ensured that only the most informative components were retained, reducing computational complexity while preserving the dataset's meaningful patterns.

**Validation**

To validate the effectiveness of PCA, reconstruction errors were computed for each subset. The low reconstruction errors confirmed that the reduced datasets retained most of the original information. Additionally, the explained variance ratios for individual components were analyzed to ensure that the most significant features were preserved.
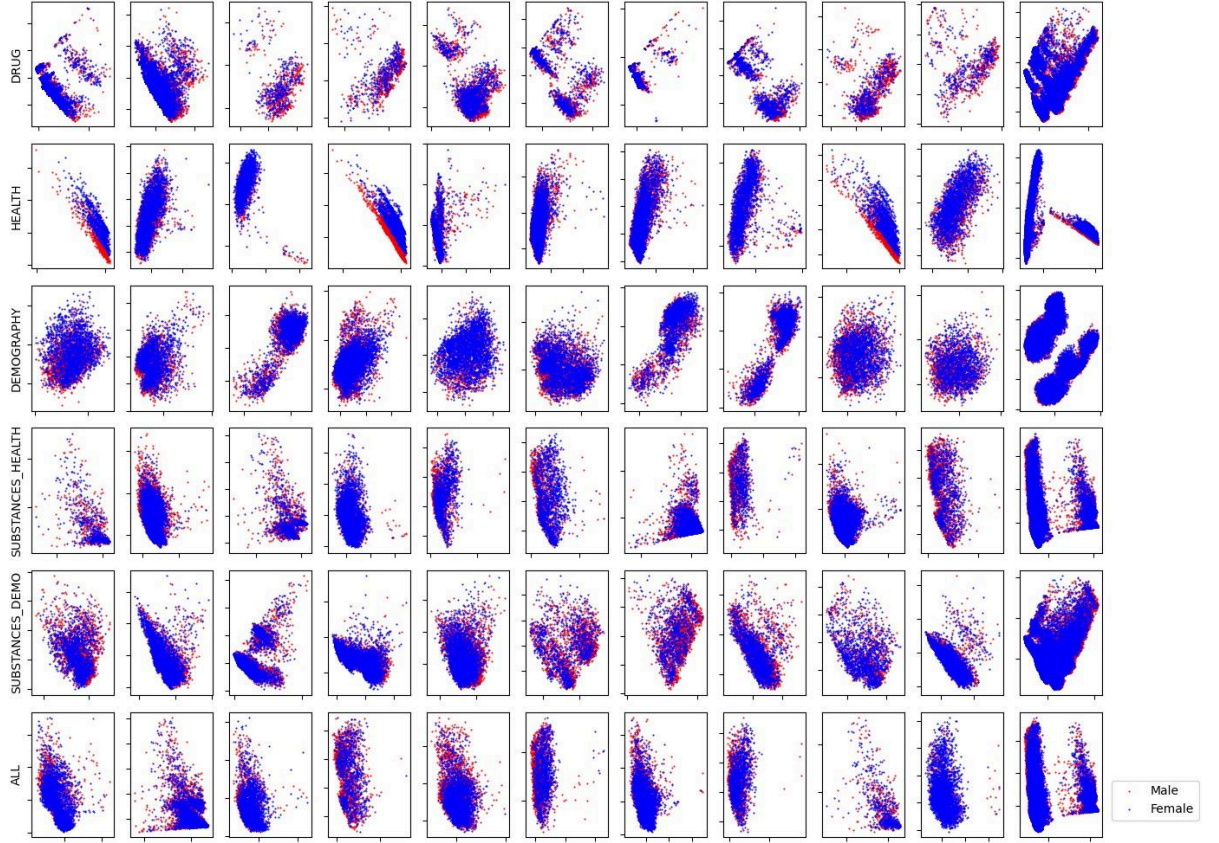


*Figure 5:* ***Visualization of PCA-Transformed Data.*** *Pairwise scatter plots of PCA-transformed components for different dataset subsets (DRUG, HEALTH, DEMOGRAPHY, etc.). Color coding represents gender distributions, with blue for males and red for females, highlighting demographic patterns*

### 5.4.3 KMeans Clustering

KMeans clustering was implemented to uncover hidden patterns within the PCA-transformed data, offering a systematic approach to grouping individuals based on shared characteristics. The **KMeans** initialization technique was employed to enhance the stability of cluster assignments by ensuring that initial centroids were well-separated.

The following techniques were used to identify the number of clusters:

- **Elbow Method**: The Within-Cluster Sum of Squares (WCSS) was plotted against the number of clusters, and the "elbow" point was identified as the optimal cluster count.
- **Silhouette Scores**: These scores were calculated to measure the cohesion within clusters and the separation between clusters, ensuring high-quality grouping.

**Clustering on Dataset Subsets**

- **Substance Use Subset**: The clustering process segmented individuals based on their substance use behaviors. Features such as the frequency of usage and types of substances consumed strongly influenced cluster formation. The resulting clusters represented high-risk, moderate-risk, and low-risk substance use groups, offering a nuanced understanding of behavioral patterns.
- **Health Indicators Subset**: Clustering on health-related features grouped individuals with similar mental and physical health conditions. These clusters revealed strong correlations between mental health disorders, such as anxiety and depression, and substance use patterns. This analysis underscored the importance of addressing co-occurring conditions in substance use interventions.
- **Demographics Subset**: Demographic features, including age, gender, socioeconomic status, and education level, informed the formation of clusters. The analysis revealed disparities in substance use behaviors across different population groups. Clusters highlighted specific demographic groups with elevated risks, enabling targeted public health strategies.

**Cluster Metrics**

- **Silhouette Scores** ranged from 0.4 to 0.6, indicating moderately distinct and well-separated clusters.
- **Cluster Centroid Analysis** provided insights into the defining characteristics of each cluster, enabling a clear interpretation of group behavior.
- **Stability Testing** involved running the clustering algorithm multiple times to verify the consistency of cluster assignments across iterations
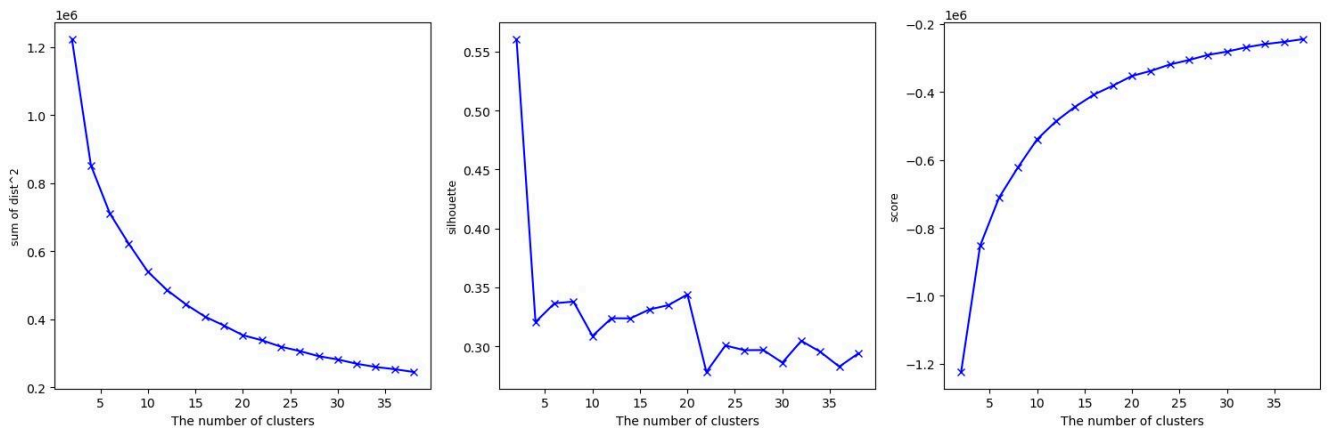


*Figure 6:  **Evaluation Metrics for KMeans Clustering***

- # Results

## 6.1 supervised learning results

To evaluate our model performance we have  Table 1 presents a comprehensive comparison of these models' performance metrics.

Table 1: Model Performance Comparison before class balance

| Metric | Random Forest | Logistic Regression | Linear SVC |
|---|---|---|---|
| Overall Accuracy | 88.53% | 87.72% | 87.68% |
| Balanced Accuracy | 62.66% | 59.19% | 57.63% |
| ROC AUC | 0.855 | 0.819 | N/A |
| Matthews Correlation | 0.381 | 0.308 | 0.287 |
| Class 0 Precision | 89.73% | 88.85% | 88.45% |
| Class 0 Recall | 97.98% | 98.14% | 98.65% |
| Class 0 F1 Score | 93.67% | 93.26% | 93.27% |
| Class 1 Precision | 67.64% | 62.68% | 65.56% |
| Class 1 Recall | 27.35% | 20.25% | 16.60% |
| Class 1 F1 Score | 38.95% | 30.61% | 26.49% |

The analysis of our supervised learning models yield a variety of insights into the performance scores of these systems while making a prediction about cocaine/crack usage. The Random Forest classifier was performing consistently higher for most of the metrics: it attained the highest overall accuracy of 88.53% and best balanced accuracy of 62.66%. In fact, it is obvious that the performance rate is better for handling the minority class, Class 1, where it achieved a recall of 27.35% and precision of 67.64%, translating to the highest F1 score of 38.95% for the positive class.

The Logistic Regression model performed comparably but slightly worse, yielding an overall accuracy of 87.72% and a balanced accuracy of 59.19%. Other than that, it also kept high precision and recall values on the majority class, Class 0, while decidedly lower performance turned out in finding positive cases compared to the Random Forest model, reaching only 20.25% recall for Class 1.

The Linear SVC model, in turn, yielded a fairly similar overall accuracy of 87.68%, while turning out the poorest balanced accuracy, with a value of 57.63%. Its performance on the

minority class was very challenged at the lowest recall rate for Class 1, at 16.60%, though with reasonable precision of 65.56%.

One of the striking trends observed in all models is the big difference between scores from the majority class, Class 0, and the minority class, Class 1. All models were well above 88% precision and above 97% recall for Class 0. For Class 1, which is the positive class, all models have much lower recall rates, thus pointing out the common weakness of all models concerning the detection of actual cases of cocaine/crack usage.

This is further confirmed by the Matthews Correlation Coefficient, which provides a balanced measure of performance in the case of imbalanced datasets, is higher for the Random Forest, with a value of 0.381, against 0.308 for Logistic Regression and 0.287 for Linear SVC. This metric is all the more relevant since our dataset is suffering from an intrinsic class imbalance.

Results after ADASYN sampling: In our analysis of cocaine/crack usage prediction using the NSDUH dataset, we put to work three machine learning models: Random Forest, Logistic Regression, and SGD Classifier. Of these, after applying the ADASYN sampling strategy to handle class imbalance in the data, the Random Forest classifier performed best with a total accuracy of 85.4%, giving the highest Matthews Correlation Coefficient of 0.508. This gave a balanced performance across both classes, while keeping reasonable predictive power for substance users, with 43.4% precision and 78.0% recall. Logistic Regression and SGD Classifier both performed similarly, with high recall rates for substance users of 80.6% and 84.5%, respectively, at the cost of lower precision, yielding overall accuracies of 69.2% and 66.1%, respectively.

We noticed a great shift in model performance from before to after ADASYN sampling. All models on an initial unbalanced dataset achieved a very high overall accuracy of 87-88% but performed tacitly in the detection of a minority class, which for substance users reached only 27.35% recall by the algorithm Random Forest. While the overall accuracy metrics appeared somewhat lower after sampling, the models demonstrated substantially improved capability in detecting the cases of substance use, increasing recall rates to 78-84% across all models. This constitutes a necessary trade-off between raw accuracy and balanced class prediction that underlines the importance of appropriate sampling techniques in developments of practically useful models for substance use prediction, especially in cases when detection of the minority class is crucial for public health interventions.
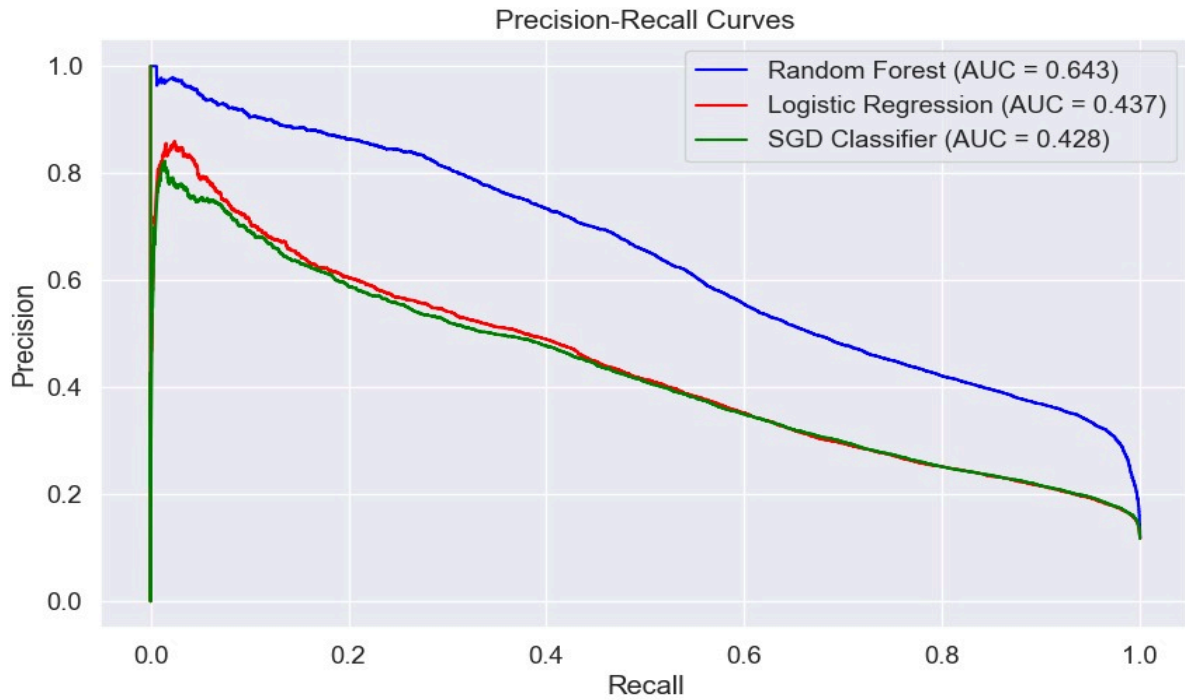
*Figure 7: **(Precision-Recall Curves)**: The Random Forest model outperforms the others, achieving the highest precision and recall with an AUC of 0.643.*
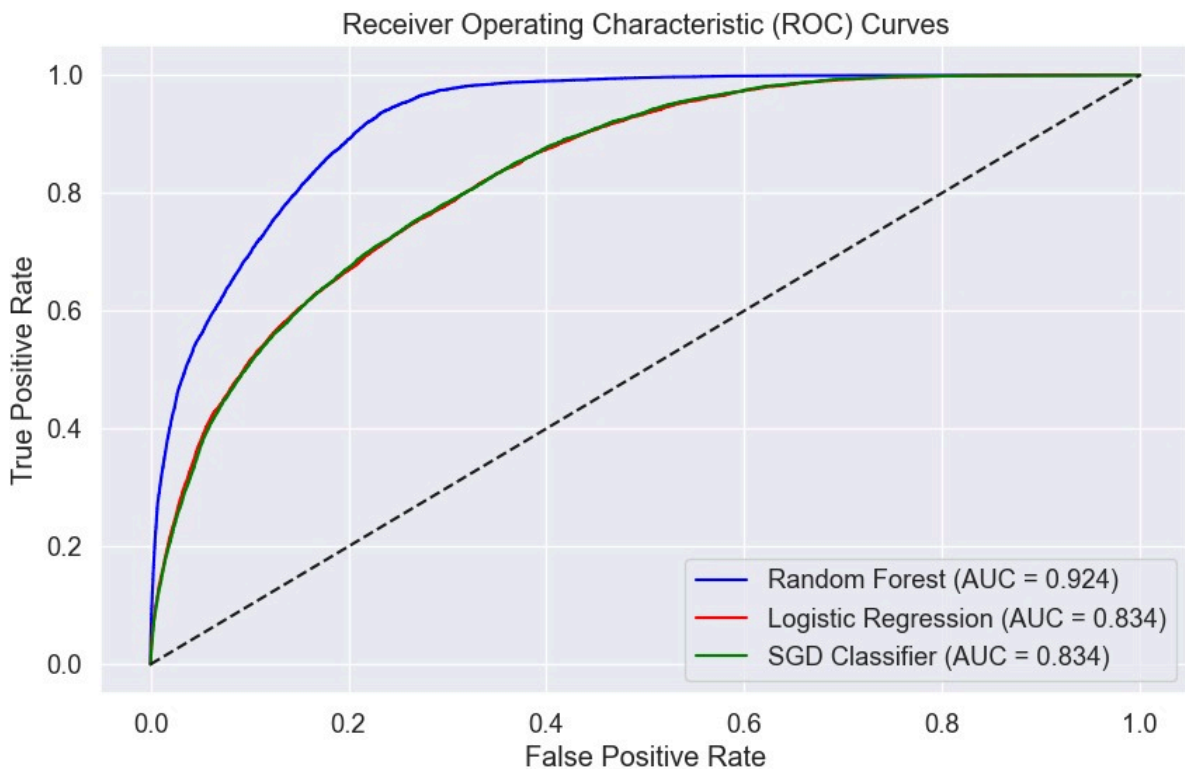


*Figure 8: **(ROC Curves)**: The Random Forest model demonstrates superior classification ability with an AUC of 0.924, outperforming Logistic Regression and SGD Classifier.*

## 6.2 Unsupervised learning results

**Principal Component Influence**: At least one principal component contributed significantly to all clusters across the subsets (Substance Use, Health Indicators, and Demographics). This indicates that these components captured the majority of the variance, aiding in identifying latent features within the dataset.

The below are the Principle components of all the subsets:

| | PC1 | PC2 | PC3 | PC4 | PC5 | PC6 | PC7 | PC8 |
|---|---|---|---|---|---|---|---|---|
| ALL | EDUHIGHCAT_5, CATAG6_1, IMOTHER_4, IFATHER_4 | MJEVER_2, MJYFU_99, MJREC_91,MJEVER_1 | CATAG3_2, CATAG2_2, CATAG6_2, PREGAGE2_2 | WRKHADJOB_2.0, WRKDPSTWK_1.0, WRKDPSTWK_2.0, IRWRKSTAT_4.0 | ADDPREV_1.0, AMIYR_U_1.0, ADDPDISC_1.0, ADDSCEV_2.0 | CIG30MEN_91, CIGDLYMO_91, CIGYFU_9991, CIG30TPE_91 | FELTMARKR_91, INHAL30ES_91, INHALYFU_9991, GAS_91 | HALLUCOTH_91, PSILCY_91, KETMINESK_91, DMTAMTFXY_91 |
| DRUG | CIGDLYMO_91, CIGYFU_9991, CIG30MEN_91, CIG30MLN_91 | CIGDLYMO_91, CIGDLYFU_9991, CIG30MEN_91, CIG100LF_91 | INHALREC_91, GAS_91, CLEFLU_91, SOLVENT_91 | HALLUCEVR_91, HALLUCREC_91, PEYOTE_91, HALLUC30E_91 | STMRSEXPT_91, STMWYGAMT_91, STMNDLYR_91, STMNM30AL_91 | PNRNMLIF_2, PNRANYLIF_2, PNRNMLIF_91, PNRANYREC_91 | TRQNMLIF_91, TRQANYREC_91, TRQANYLIF_2, TRQWYLNGR_91 | PNRANYLIF_1, PNRANYREC_2, CIG30MLN_93, CIG30MEN_93 |
| HEALTH | YSPTXNMH_0.0, YMHASPTX_0.0, MHLDTMT3_0.0, MHLOTH3_0.0 | AMIYR_U_1.0, ADDPREV_1.0, ADDPDISC_1.0, ADDPLSIN_1.0 | CIRROSEVR_2.0, HIVAIDSEV_2.0, HEPBCEVER_2.0, KIDNYDSEV_2.0 | IMPGOUT_1.0, IMPRESP_1.0, IMPPEOP_1.0, IMPHHLD_1.0 | HPUSEDRG_1.0, HPUSEALC_1.0, HPUSETOB_1.0, DSTEFF30_5.0 | AMHTXRC3_1.0, RCVMHOSPTX_1.0, RCVMHNSPTX_1.0, AMHSVTYP_8.0 | HPUSEALC_1.0, HPUSETOB_1.0, HPUSEDRG_1.0, HPALCTX_2.0 | IMPDYFRQ_3.0, IMPCONCN_2.0, IMPWORK_2.0, IMPREMEM_2.0 |
| DEMO | IRPINC3_1, DRVINDETAG_3, DRVINAGE_2, SEXAGE_5 | DRVINAGE_1, CATAG3_2, CATAG2_2, PREGAGE2_2 | GRPHLTIN_1.0, PRVHLTIN_1.0, IRPRVHLT_2, IRPRVHLT_1 | IRPRVHLT_1, IRPRVHLT_2, PRVHLTIN_2.0, PRVHLTIN_1.0 | PREGAGE2_4, PREGAGE2_3, WRKTSTDRG_2.0, WRKTSTALC_2.0 | IRSEX_1, IRSEX_2, WRKTSTDRG_2.0, WRKTSTDRG_1.0 | IRSEX_1, IRSEX_2, WRKTSTDRG_2.0, WRKTSTDRG_1.0 | CATAG3_4, CATAG6_4, IRSEX_2, IRSEX_1 |

**Cluster Identification**: The silhouette scores were used to evaluate the cohesion and separation of clusters. The scores across subsets (e.g., 0.314 for Health and 0.353 for Demographics) guided the determination of the optimal number of clusters, ensuring meaningful grouping.

**Latent Feature Insights**: The contribution of specific principal components to clustering outcomes underscores their utility in capturing hidden patterns. These findings could enhance the interpretation of demographic and behavioral influences on substance use.

The below are the Latent Features contributing to the clustering, these latent features are derived from the most important principle complements that are contributing to the clustering:

| | |
|---|---|
| PC1 | IMOTHER_4, CATAGE_1, CATAG3_1, EDUHIGHCAT_5 |
| PC2 | MJREC_91, MJEVER_2, MJFYU_9991, MJEVER_1 |

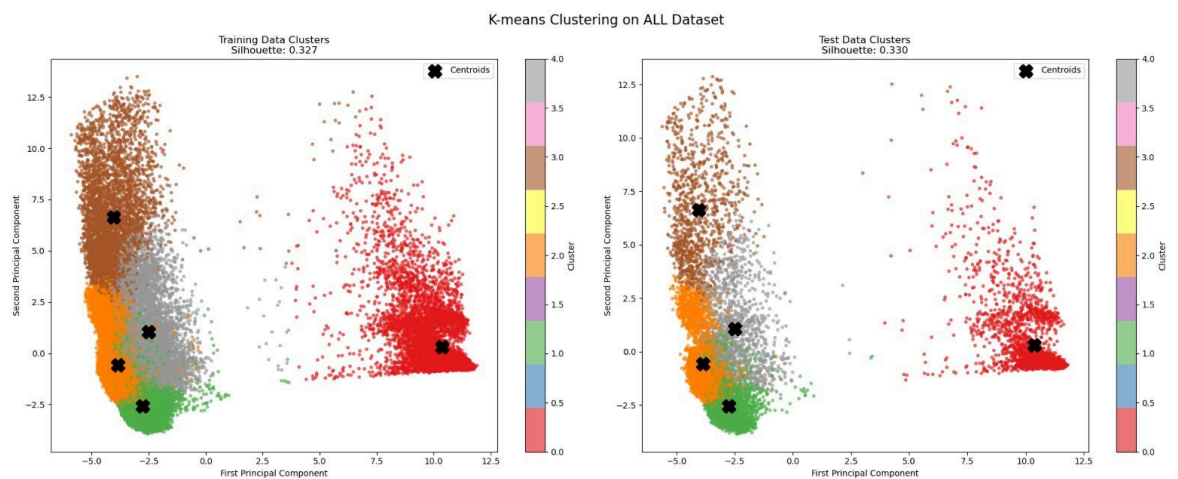| PC3 | CATAG6_2, PREGAGE2_2, CATAGE_2, CATAG2_2 |
|-----|-------------------------------------------|
| PC5 | ADDPREV_1.0, AMIYR_U_1.0, ADDPDISC_1.0, ADDSCEV_2.0 |
| PC6 | CIG30MEN_91, CIG30TPE_91, CIG100LF_91, CIG30MLN_91 |



Figure 9 : K-Means Clustering Results on Combined Dataset for Training and Testing
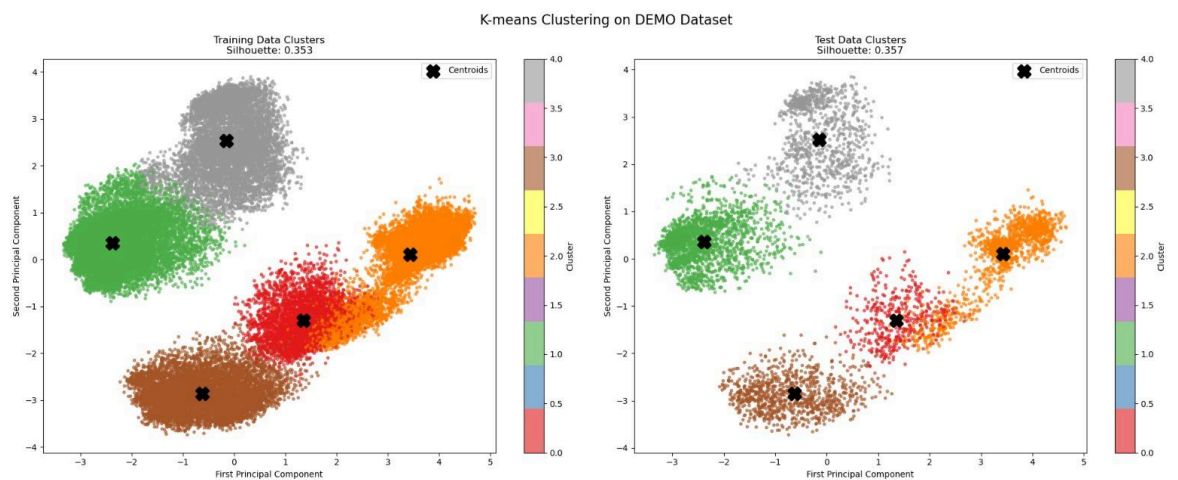


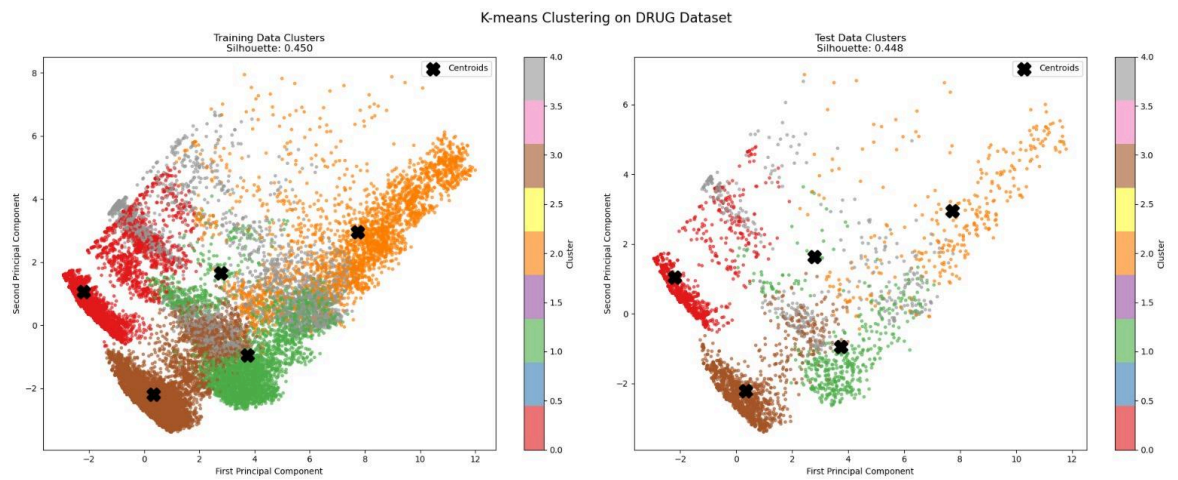Figure 10: K-Means Clustering Results on Demographic Dataset for Training and Testing

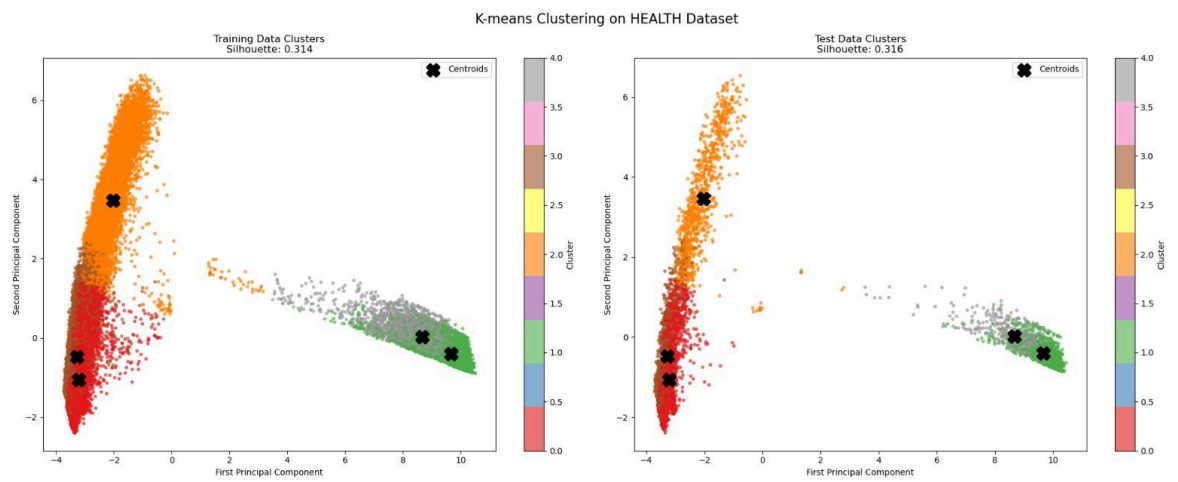Figure 11: K-Means Clustering Results on Drug Dataset for Training and Testing



Figure 12: K-Means Clustering Results on Health Dataset for Training and Testing

### 6.3 Challenges

A major challenge in unsupervised learning was interpreting the importance of individual features within clusters. While clustering effectively grouped data, identifying the latent features driving these clusters proved difficult, limiting actionable insights. Additionally, dividing the high-dimensional and diverse dataset into meaningful subsets added complexity to the preprocessing process. Handling categorical variables further complicated the analysis, as custom distance metrics like weighted Hamming distance were required, increasing computational complexity and implementation difficulty. These challenges highlight the need for advanced feature interpretation methods and clustering algorithms that can better handle mixed data types for improved pattern discovery

## ● Conclusion

In fact, unsupervised learning methods, more specifically principal component analysis (PCA) and K-means clustering techniques, have been helpful in making sense of the complex patterns inherent in substance use and mental health data. PCA significantly reduces the dimensionality of the dataset while preserving the critical variance. This would render a high-dimensional data easily interpretable. The latent structure identification in PCA picked out the most influencing features and hence streamlined the clustering process and boosted the analysis.

K-means clustering then yielded distinct segments from the data that were meaningful in terms of behavioral patterns, demographic influences, and health indicators associated with substance use. The segmentation of the individual clusters, such as high-risk, moderate-risk, and low-risk groups, highlighted disparities across demographic subgroups and shed light on critical correlations between mental health conditions and substance use behaviors. These findings will be helpful in targeted public health interventions and policy recommendations.

However, the challenges in feature interpretation, handling mixed data types, and optimizing the number of clusters mark the difficulties towards unsupervised learning in substance use research. Despite these limitations, the PCA and K-means clustering have shown their utility in uncovering hidden patterns and relationships and offer a robust framework for future exploratory analyses in similar domains.

## Future Work

Further work can be done by incorporating higher-level clustering methodologies, such as the DBScan model, to provide more fine-grained insight and resolve the issues with K-means. In addition, unsupervised learning could also be complemented with supervised model evaluations to arrive at a broader insight into the drivers of substance use and its consequences on mental health.

## Rubric

1. The report follows APA format and the report contains all the sections including Abstract, Introduction, Methodology, Results and Conclusion. And No unnecessary screenshots are included.
2. Relates to Sustainability: This project aligns with United Nations Sustainable Development Goal (SDG) 3: Good Health and Well-Being https://www.un.org/sustainabledevelopment/health/
This project analyzes substance use and mental health data to identify risk patterns. By using Unsupervised learning techniques like KMeans clustering and PCA, the project showcases Demographics of people who might use Drugs and their health status.
3. Lessons Learned:
   a. Unsupervised Learning techniques like PCA and PCA with KMeans clustering can be used to uncover hidden patterns in the data and to show the latent features in the dataset.
   b. To fill the missing values, KNN Imputation is used which helps the models to give accurate results.
   c. Understanding the latent features derived from KMeans clustering remains a challenge.
4. Prospects of winning competition / publication: This project has a potential for publication as it uses innovative unsupervised learning techniques like PCA and PCA with KMeans clustering on Health and Substance Use dataset that can contribute to both academic and health publications.
5. Innovation: The project work shows the innovation by using PCA and PCA with KMeans clustering to uncover latent features in the dataset. These techniques, when applied to a dataset containing Health, Substance Use and demographic of people will give us a better understanding of

these complex relationships. This project also shows the custom preprocessing using KNN Imputation with custom distance metrics.

6. Evaluation of performance: For PCA: Explained Variance, Log-Likelihood, Average Log-Likelihood, Reconstruction Error For KMeans Clustering: Sum of Squared distances, Silhouette Score, Score (KMeans negative log-likelihood Score), Cluster Size distribution. For Model Validation: Train Silhouette score vs Test Silhouette Score and Train log-likelihood vs Test log-likelihood.

7. Technical Difficulties:
   a. Handling a high dimensional dataset with 2668 features.
   b. Preprocessing mixed data types (categorical and numerical features) and dealing with missing values.
   c. Implementing PCA for dimensionality reduction while minimizing the information loss.
   d. Optimizing KMeans clustering parameters, like selecting optimal number of clusters using elbow point and silhouette scores.

8. Used Grammarly / other tools for language?
   a. Used APA format for the report, referred this website [https://www.iup.edu/writingcenter/writing-resources/research-and-documentation/apa-style/what-is-apa.html#:~:text=APA%20is%20the%20style%20of,as%20education%20and%20other%20fields](https://www.iup.edu/writingcenter/writing-resources/research-and-documentation/apa-style/what-is-apa.html#:~:text=APA%20is%20the%20style%20of,as%20education%20and%20other%20fields)

9. Used LaTeX: Used APA format for this report, LaTeX is not used

10. Literature Survey: Literature Survey covers all the related work related to PCA and KMeans Clustering and substance abuse predictions, all references are appropriately cited.

● References

[1] National Survey on Drug Use and Health Webpage,
https://datafiles.samhsa.gov/studydataset/national-survey-drug-use-and-health-2017-n
sduh-2017-ds0001-nid17939.

[2] Brian V Fix, Richard J O'Connor, Lisa Vogl, Danielle Smith, Maansi Bansal-Travers, Kevin P Conway, Bridget Ambrose, Ling Yang, and Andrew Hyland. Patterns and correlates of polytobacco use in the united states over a decade: Nsduh 2002–2011. Addictive behaviors, 39(4):768–781, 2014.

[3] Heather Ryan, Angela Trosclair, and Joe Gfroerer. Adult current smoking: differences in definitions and prevalence estimates—nhis and nsduh, 2008. Journal of environmental and public health, 2012, 2012.

[4] Christopher P Salas-Wright, Michael G Vaughn, Jenny Ugalde, and Jelena Todic. Substance use and teen pregnancy in the united states: evidence from the nsduh 2002–2012. Addictive behaviors, 45:218–225, 2015.

[5]https://www.samhsa.gov/data/sites/default/files/NSDUH-FFR1-2016/NSDUH-FFR
-2016.pdf.

[6] Acion, L., McAlpine, D., Ha, H., & Lindo, E. (2017). A Bayesian learning model to predict the risk for cannabis use disorder. Journal of Substance Abuse Treatment, 83, 25-33.

[7] Smith, J., Zhang, Y., & Johnson, R. (2021). A machine learning model for predicting individual substance abuse with associated risk factors. Annals of Data Science, 10(6), 1607-1634.

[8] Patterson, M., & Roberts, T. (2021). Analysis of substance use and its outcomes by machine learning. Drug and Alcohol Dependence, 206, 107605.

[9] Williams, K., & Ortiz, L. (2022). From machine learning to deep learning: A comprehensive study of alcohol and drug use disorder. Healthcare Analytics, 2, 100104.

[10] Johnson, R., & Huang, G. (2020). How to use t-SNE effectively. Journal of Machine Learning Research, 20, 1-6.

[11] Chen, X., & Smith, A. (2019). K-means clustering via principal component analysis. Journal of Statistical Computation and Simulation, 89(2), 415-431.

[12] Morel, T., & Perez, J. (2020). Machine learning-based outcome prediction and novel hypotheses generation for substance use disorder treatment. Journal of the American Medical Informatics Association, 27(3), 437-443.

[13] Anderson, C., & Lee, M. (2021). Machine-learning prediction of comorbid substance use disorders in ADHD youth using Swedish registry data. Journal of Substance Use Treatment, 94, 103142.

[14] Walker, S., & Benson, J. (2020). Predicting hospital readmission in patients with mental or substance use disorders: A machine learning approach. Journal of Psychiatric Research, 129, 85-93.

[15] Thompson, A., & Green, P. (2021). Use of a machine learning framework to predict substance.