PROGRAMMING FOR DATA ANALYTICS

PROJECT REPORT

TOPIC: CUSTOMER SEGMENTATION AND CLASSIFICATION USING MACHINE
LEARNING

CODE: 04-638A

Instructor: Dr. George Okeyo

Name: Ange Izabayo

Andrew ID: aizabayo

MS ECE 25'

Submission date on December 17, 2023

<h1 style="text-align:center">Technical Report</h1>

<h2>Abstract</h2>

For the financial banking analysis Credit card Customer segmentation play essential role tailoring Marketing strategies for diverse user groups. This project focuses on credit card customer segmentation, by addressing the need for personalizing financial services. The main purpose is to employ data-driven approach for effective customer segmentation, whereby studying on sample of 8950 credit card holders and 18 behaviour variables. The analysis involves handling missing values, outlier detection, and scaling data. Unsupervised learning, specifically KMeans clustering, is utilized to group customers. Supervised models, including Logistic Regression, Random Forest, and Support Vector Machine, are employed to determine the optimum accuracy within clusters. Hyperparameter tuning enhances model performance. The final model is deployed through a Flask web application for predicting customer segmentation, aiding in tailoring marketing strategies based on diverse user groups.

<h2>Background and Problem Description</h2>

This project focuses on credit card customer segmentation to enhance personalized financial services using a dataset from Kaggle. The primary challenge addressed is the ineffective categorization of customers, leading to low credit card usage rates and a prevalence of underutilized sleeping cards. Institutions often issue cards indiscriminately, neglecting consumer behaviour research and hindering the ability to provide tailored services, impacting customer loyalty. Customer segmentation is pivotal for targeted marketing. This project addresses the challenge of predicting customer segments for effective marketing strategies..[1], [2]

<h2 style="text-align:center">Approach</h2>

1. **Data preparation:** The initial step in data preparation involved reading the Kaggle dataset into a Pandas dataframe. To gain an insightful overview, the first and last 10 rows of the dataframe were displayed, offering a glimpse into the dataset's structure.
2. **Exploratory Data Analysis:** Perform the basic statistical measures and visualizations such as histograms and pairwise plots were used to understand feature distributions and relationships. A heatmap highlighted notable correlations, like 0.92 between purchases and one-off purchases, and identified independent features, e.g., -0.066 correlation between purchase transactions and cash advances.
3. **Pre-processing:** During data pre-processing I handled missing value by replacing with the mean, and then check the outlier in the dataset if they will not affect the result. And due to small dataset, I did not handle outlier. Then scale the data using standard scaler. Regarding the label encoder I did not use it because the categorical columns was dropped. Where the fitted scaler is saved as pickle.
4. **Unsupervised Model creation and Evaluation:** Utilized the Elbow Method and silhouette score in the KMeans unsupervised model to determine the optimal number of clusters for customer segmentation. The analysis identified three clusters, balancing within-cluster sum of squares and inter-cluster separation. The silhouette score of 0.251 affirmed the reasonable coherence and separation within these clusters. And save the label csv file.
5. **Supervised Model and Evaluation:** Implemented a supervised classification model employing Logistic Regression, Random Forest, and Support Vector Machine (SVM) algorithms on a labelled dataset. Feature scaling using the Standard Scaler preceded model evaluation, where cross-validation and accuracy served as the primary metric due to balanced classes. Learning curves were utilized to assess model fitting.

6.  **Feature Selection and Engineering**: Utilized SelectKBest and f_classif for feature selection, followed by PCA for dimensionality reduction. Integrated Logistic Regression, Random Forest, and Support Vector Classifier (SVC) with the selected features and PCA-transformed data. Cross-validation assessed model performance, and average accuracy scores were recorded. Trained models were saved as pickle files for future use.[3]

7.  **Hyperparameter tuning:** Comparison of model performance before and after feature selection highlighted Logistic Regression as the benchmark. Employed Grid Search for hyperparameter tuning, exploring parameters like regularization strength, penalty type, and solver type. Identified the best hyperparameter combination and saved the tuned Logistic Regression model. Evaluated the tuned model's performance and generated a visualization for comparative accuracy.[4]

8.  **Model Debugging:** Learning curves are employed for model debugging, providing insights into performance during training and validation. These curves facilitate the identification of overfitting or underfitting issues by visualizing trends in training and validation accuracy. This iterative debugging process ensures the delivery of accurate and reliable credit card customer segmentation results.[5]

9.  **Model deployment:** Flask, a Python web framework, is utilized to create an interactive interface. The trained model is seamlessly integrated into the Flask web application, enabling users to input features and receive real-time predictions. This approach prioritizes accessibility and usability for diverse stakeholders in the banking domain.
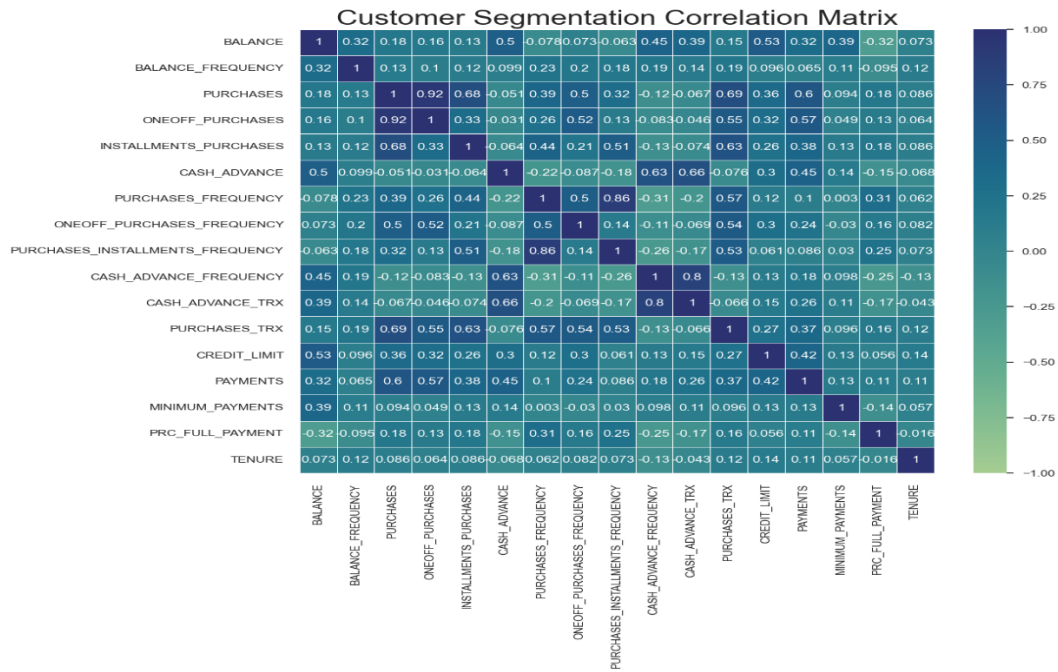
10. **Data Preparation**

    The exploration of the dataset revealed valuable information about its structure and content. By displaying the head and tail of the dataframe, I ensured a comprehensive understanding of the data's characteristics. Additionally, the examination of columns facilitated the identification ofkey features, setting the stage for further analysis. The process also involved checking for missing values, a crucial step in ensuring data integrity.

## Result

### 1.  Exploratory Data Analysis

The statistical overview, histograms, pairwise plots, and heatmap collectively contributed to a nuanced exploration of the dataset. Key observations included notable correlations between specific features and instances of independence, enhancing the understanding of feature dynamics. Example purchase and on/off purchases with coefficient of 0.92. For negative correlation Purchase _TXR and Cash Advanced TR with -0.066 which shows that there are not correlated each is independent.
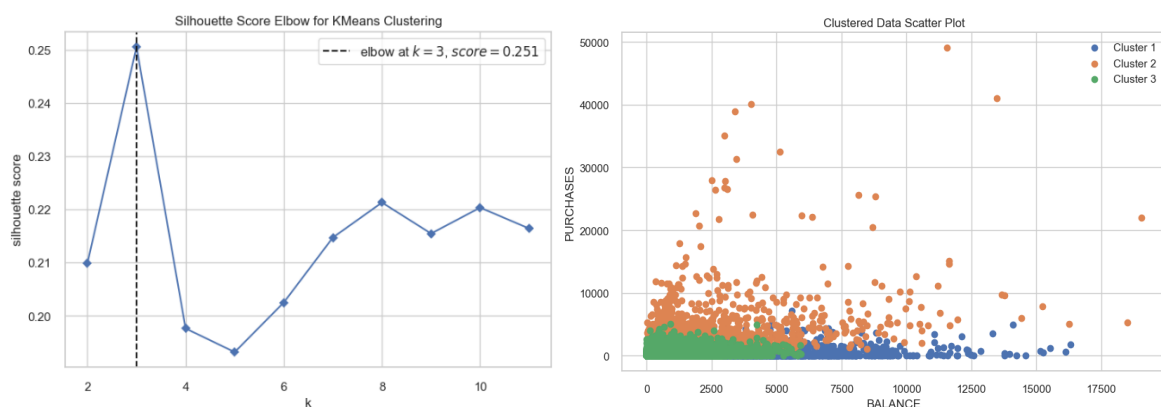
Customer Segmentation Correlation Matrix

## 2. Pre-processing

During data pre-processing I handled missing value by replacing with the mean, and then check the outlier in the dataset if they will not affect the result. And due to small dataset I did not handle outlier. Then scale the data using standard scaler. Regarding the label encoder I did not use it because the categorical columns was dropped. Where the fitted scaler is saved as pickle.

## 3. Unsupervised Model creation and Evaluation

The selection of three clusters as shown in figure below is pivotal for effectively grouping customers based on underlying patterns in the data, providing a foundation for targeted marketing strategies or personalized services tailored to the distinct needs of each cluster. The selection of three clusters forms a crucial basis for effective customer grouping, facilitating targeted marketing strategies and personalized services aligned with the unique characteristics of each cluster.



Example plotting to visually represent the clustering results in a two-dimensional space defined by the features BALANCE and PURCHASES. The code iterates through each cluster, extracts the data points belonging to that cluster, and plots them with distinct colours for each cluster.

## 4. Supervised model and Evaluation

Supervised classification model was built and evaluated using Logistic Regression, Random Forest, and Support Vector Machine (SVM) algorithms on a labelled dataset. The features were first scaled using the Standard Scaler. For model evaluation, cross-validation was employed, and accuracy was chosen as the primary metric for its suitability in balanced classes. Learning curves were generated to assess model fitting, indicating that all three models: Logistic Regression, Random Forest, and SVM

```
Logistic Regression Accuracy: 0.9956424581005587
Random Forest Accuracy: 0.9681564245810055
Support Vector Classifier - Cross-Validation Accuracy: 0.8927374301675977
```

The output shows that the logistic regression has high accuracy compared to other models.
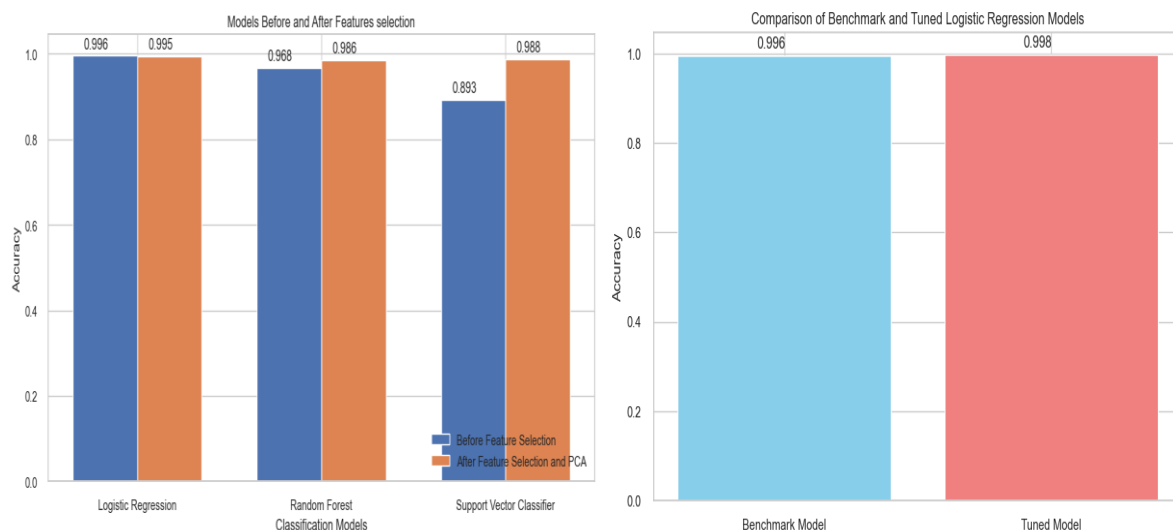
## 5. Feature Selection and Engineering

With the selected features are 12 in 17 features where below is the cross-validation score. After feature selection, Logistic Regression showed a slight decrement in performance, while other models exhibited performance improvement. The bar graph below illustrates the comparative accuracy of each model.

```
Logistic Regression Accuracy with PCA: 0.9950837988826816
Random Forest Accuracy with PCA: 0.986145251396648
SVC Accuracy with PCA: 0.9884916201117319
```

Comparing the performance after feature selection logistic regression decremented a bit but other model performance increases and the bar graph below or comparing accuracy of each.
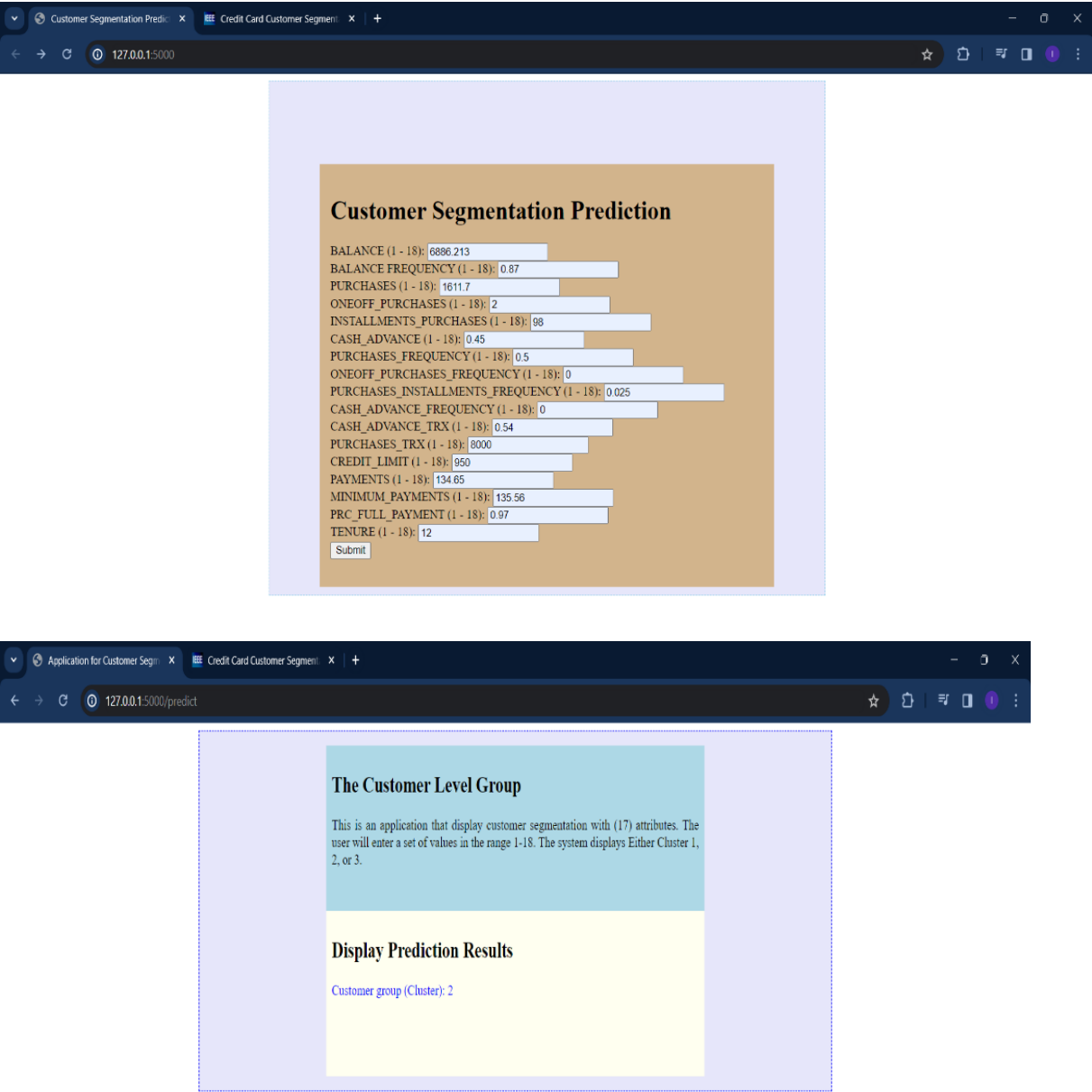
## 6. Hyperparameter tuning

Hyperparameter tuning significantly improved model performance, enhancing accuracy from 0.996 in the benchmark Logistic Regression to 0.988 after tuning, as illustrated in the figure above. This iterative debugging process ensures the delivery of accurate and reliable credit card customer segmentation results.



Model debugging using learning Curve: To ensure the robustness and effectiveness of the models, learning curves are utilized for debugging. Learning curves provide insights into the model's performance during training and validation phases, aiding in identifying issues such as overfitting or

underfitting. By visualizing the trends in training and validation accuracy. This iterative process of debugging ensures the delivery of accurate and reliable credit card customer segmentation results.

**Model deployment**: Flask, a web framework for Python, is employed for creating an interactive and user-friendly interface. The trained model is integrated into the Flask web application, allowing users to input relevant features and receive real-time predictions. This deployment strategy ensures accessibility and usability for diverse stakeholders in the banking domain.





The Figure above illustrate a user-friendly Flask web application designed for credit card customer segmentation prediction. Users input values for 17 financial attributes, ranging from balance and purchase frequency to credit limit and tenure. The application promptly displays the predicted customer group or cluster, exemplified by Cluster 2 in the provided instance. This interactive platform, with its intuitive design and real-time predictions, offers stakeholders valuable insights into customer behaviour, fostering informed decision-making for tailored marketing strategies and improved customer engagement.

## Conclusion

In Conclusion, the project utilized a data-driven approach for credit card customer segmentation, employing unsupervised KMeans and supervised models and different means for accurate predictions. The deployment of the tuned model using Flask ensures a user-friendly interface for real-time customer group predictions. The results offer financial institutions valuable insights for tailored marketing strategies and customer-level decision-making. For more clarification markdown in notebook details the work.

## Reference

[1] W. Li, X. Wu, Y. Sun, and Q. Zhang, "Credit Card Customer Segmentation and Target Marketing Based on Data Mining," in *2010 International Conference on Computational Intelligence and Security*, Dec. 2010, pp. 73–76. doi: 10.1109/CIS.2010.23.

[2] "Segmenting credit card customers for insight and profit." Accessed: Dec. 17, 2023. [Online]. Available: https://www.electronicpaymentsinternational.com/features/segmenting-credit-card-customers-for-insight-and-profit/

[3] "sklearn.feature_selection.SelectKBest," scikit-learn. Accessed: Dec. 17, 2023. [Online]. Available: https://scikit-learn/stable/modules/generated/sklearn.feature_selection.SelectKBest.html

[4] M. Gusarov, "Do I need to tune logistic regression hyperparameters?," CodeX. Accessed: Dec. 17, 2023. [Online]. Available: https://medium.com/codex/do-i-need-to-tune-logistic-regression-hyperparameters-1cb2b81fca69

[5] "Underfitting vs. Overfitting," scikit-learn. Accessed: Dec. 17, 2023. [Online]. Available: https://scikit-learn/stable/auto_examples/model_selection/plot_underfitting_overfitting.html