

## **DATA SET 1 - CUSTOMER SEGMENTATION: CLUSTERING (KARNIKA KAPOOR)**

(The notebooks have been systematically critiqued, following the headings and subheadings used within them for easy cross-referencing. Cell numbers are also provided where deemed appropriate)

### **1. Project Objective:**

The primary objective of the project is to perform unsupervised clustering on customer records from a groceries firm's database for effective customer segmentation. The project is quite thorough in addressing the objective of customer segmentation. The step-by-step approach, starting from data cleaning and pre-processing to dimensionality reduction and clustering, is well-documented. The use of exploratory data analysis to understand the clusters and their characteristics is a good practice. Additionally, the visualizations provided give a clear insight into the distribution and patterns within the clusters. The conclusions drawn from profiling the clusters based on personal traits and spending habits add a valuable layer to understanding customer behavior.

Overall, the project appears to be successful in achieving its goal of customer segmentation and provides a foundation for informed marketing strategies.

### **2. Data Cleaning:**

#### **2.1 Handling missing values: (Code Line 5)**

Dropping missing values is deemed acceptable in this context, given its minimal impact on the dataset size. The reduction from 2240 to 2216 data points, further dwindling to 2212 post-outlier removal, remains within an acceptable range. **Marcelino<sup>1</sup>** underscores the significance of handling missing values diligently. The decision to **discard rows with missing income values** aligns with this principle, preventing the introduction of potentially biased or inaccurate information. This strategic approach contributes to a cleaner and more reliable dataset. An **alternative strategy for addressing missing values involves imputation based on the mean or median** of the column, avoiding row removal. Inspired by Pedro Marcelino's approach, efforts are invested in feature engineering to extract meaningful information.

#### **2.2 Feature Engineering: (Code Line 9)**

Pros of feature engineering emphasize the potential for richer insights and simplified representations. However, cons underscore the need for **cautious handling to mitigate risks associated with correlation, redundancy, and model sensitivity**. For instance, the introduction of "Customer\_For," representing the duration of customer enrolment, prompts an exploration of its correlation with features like "Recency" and "Dt\_Customer." Assessing whether prolonged customer tenure correlates with recent interactions or specific enrolment periods informs the contextual relevance of the new feature. Similarly, the creation of "Living\_With" and "Is\_Parent" warrants a thorough exploration of their correlation with demographic indicators like "Marital\_Status" and "Children." Ensuring consistency in living situations inferred from "Marital\_Status" and those deduced by "Living\_With" and evaluating the correlation between "Is\_Parent" and features like "Kidhome" and "Teenhome" validate their effectiveness in capturing parenthood status.

The critical analysis extends to investigating the impact of these engineered features on **predictive capabilities**. Understanding how "Customer\_For" influences response variables like "Response" provides insights into its predictive utility. Rigorous evaluation involves statistical measures, such as **correlation coefficients and significance tests**, to quantify relationships.

### **3. Data Pre-processing:**

The project aims to pre-process the data to enhance its suitability for analysis and modelling. It involves tasks such as encoding categorical variables, scaling numerical features, and splitting the dataset into training and testing sets. While the overall objective is addressed, there are areas that require attention.

#### **3.1 Identification of Categorical Variables: (Code Line 14)**

Guided by Pedro Marcelino's Kaggle exploration, crucial categorical variables ('Education' and 'Living\_With') are discerned through a thoughtful strategy. This lays the groundwork for subsequent transformations, but **misidentification carries the risk of skewed pre-processing**.

---

<sup>1</sup> Reference notebook given in coursework for guidance - "**Comprehensive data exploration with Python**" by **Pedro Marcelino**

### 3.2 Label Encoding: (Code Line 15)

Employing the LabelEncoder class in a loop, 'Education' and 'Living\_With' undergo label encoding to ensure numerical alignment for efficient machine learning. However, the implicit assumption of ordinal relationships in label encoding may not universally hold. **An alternative approach, one-hot encoding, preserves categorical nature without introducing a false order.**

### 3.3 Subset DataFrame Creation and Scaling: (Code Line 16)

A targeted subset DataFrame ('ds') is meticulously crafted, excluding irrelevant features related to deals and promotions. Concurrently, the StandardScaler normalizes this subset. While enhancing model focus, it's crucial to consider potential information loss. Exploring **scaling alternatives like MinMax or robust scaling,** based on data characteristics, could be beneficial. Additionally, careful thought is needed when reducing the dataset, as it may lead to information loss. **An alternative approach could involve clustering on the entire dataset with subsequent dimensionality reduction.**

## 4. Dimensionality Reduction and PCA: (Code Line 18/19)

Acknowledging the complexity introduced by numerous correlated features, a proactive choice is made to tackle this challenge through dimensionality reduction. Inspired by Pedro Marcelino's methodology, the project opts for **Principal Component Analysis (PCA), a potent technique that effectively minimizes redundancy while retaining crucial information.** Despite its power, reducing dimensions to three poses a risk of information loss, potentially impacting the dataset's comprehensiveness. The decision to limit dimensions should be carefully weighed against the desired balance between interpretability and variance preservation. Additionally, while visualizing the dataset in three dimensions offers valuable insights, exploring **alternative techniques like t-SNE could provide valuable points of comparison.** While PCA proves robust, a comprehensive understanding of its assumptions and limitations is imperative for making informed decisions in the modelling process.

## 5. Clustering: (Code Line 20-22)

The application of the **hierarchical clustering method involves iteratively merging** examples until the desired number of clusters is achieved. An **alternative approach like DBSCAN** could be more robust for the customer segmentation dataset. The dataset's characteristics, encompassing customer information with columns like 'Income,' 'Recency,' and 'MntWines,' may exhibit non-linear patterns in customer behavior. These non-linearities may arise due to complex interactions between different features, and **DBSCAN's adaptability to varying cluster densities and dynamic determinations makes it promising for capturing intricate patterns in the data.** (Alternative code for clustering given at the end of notebook)

### 5.1 Elbow Method for Determining Clusters:

The current method, using the **Elbow Method with KMeans clustering, is effective but lacks clear annotations** in its visual representation. Considering dataset characteristics and potential variations in customer behavior patterns, an **alternative approach using the silhouette score could better capture diverse cluster shapes and patterns.**

### 5.2 Clustering via Agglomerative Clustering:

Agglomerative Clustering aligns with the hierarchical approach but lacks detailed rationale behind selecting four clusters, hindering transparency. An alternative like **Gaussian Mixture Models (GMM) provides a probabilistic clustering approach** and flexibility in capturing complex data distributions.

### 5.3 Examining Clusters via 3-D Scatter Plot:

The **3-D scatter plot** aids in interpreting cluster distribution but lacks axis labels and a legend, diminishing clarity. Exploring alternatives like **t-Distributed Stochastic Neighbor Embedding (t-SNE) for a nuanced view of cluster relationships could enhance the visualization.**

While the hierarchical clustering method demonstrates sound principles, improvements in visualization annotations and rationale documentation are needed. Considering alternatives like DBSCAN, Silhouette analysis, and GMM for specific steps might provide a more tailored and insightful approach for the customer segmentation project.

## 6. Evaluating Models: (Code Line 23)

The current approach in evaluating models for customer segmentation clustering demonstrates a comprehensive exploration of cluster patterns and customer behavior. The countplot effectively visualizes the distribution of clusters, providing a quick overview. The scatter plot of 'Income' vs. 'Spending' skillfully identifies distinct spending patterns

within each cluster, aiding in interpretation. The detailed product-wise spending distributions and analysis of campaign responses contribute to a holistic understanding of customer preferences.

There are areas where the current approach could be enhanced. **The lack of annotations and clarity in the countplot and scatter plot diminishes interpretability, affecting the communicative value of the visualizations. Additionally, the absence of axis labels and legends in the scatter plot reduces its overall effectiveness.**

#### 6.1 Quantitative Analysis: (Code Line 24-28)

Quantitatively, the joint plots offer valuable insights into the relationship between various purchasing styles and spending. Still, the analysis could benefit from statistical measures to quantify the observed patterns and relationships, providing a more robust foundation for decision-making. As for alternatives, incorporating statistical tests, such as correlation coefficients, significance tests, or regression analyses, could strengthen the quantitative analysis. Utilizing advanced visualization techniques, like **heatmaps or dendrograms, might offer a clearer representation of relationships within the clusters.** Furthermore, considering **machine learning models for pattern recognition, such as decision trees or random forests, could provide additional depth to the analysis.**

#### 7. Profiling: (Code Line 29)

The existing profiling approach delves into customer spending habits based on personal traits within formed clusters, providing valuable insights into diverse behaviours. However, to refine this analysis, incorporating **advanced clustering algorithms or ensemble methods tailored to the retail domain** could offer a more nuanced understanding of customer segments.

The dataset features diverse customer attributes, including spending patterns, demographics, and responses to marketing campaigns. With objectives centred around understanding and profiling customer segments for targeted marketing, the dataset's retail-centric nature emphasizes customer behaviour and preferences.

To evaluate alternatives, it's crucial to consider the necessity for accurate segmentation that reflects the nuances of retail dynamics. The project's goals revolve around effective customer profiling and informed marketing strategies within the retail domain.

##### 7.1 Exploring Alternatives:

Experimenting with **hierarchical clustering or Gaussian mixture models can enhance cluster granularity, capturing subtler patterns in customer behaviour.** Integrating demographic segmentation aligns seamlessly with the retail context, providing a more precise understanding of customer profiles.

The exploration of alternative methods should be tailored to enhance the project's objectives, emphasizing effective customer profiling and the formulation of informed marketing strategies within the specific context of the retail domain.

#### 8. Conclusion:

In conclusion, while the methods employed in the project exhibit solid foundations, there is **room for enhancement and strategic exploration of alternatives.** The data cleaning and feature engineering approaches showcase diligence, but considering alternatives such as **imputation strategies** could further refine the dataset. The preprocessing steps and dimensionality reduction through PCA are commendable, yet a more thorough exploration of scaling alternatives and dimensionality reduction techniques like **t-SNE** could offer additional insights. Clustering methods, specifically hierarchical clustering, form a strong basis, but alternative algorithms like **DBSCAN and GMM** should be carefully considered for their potential to capture intricate non-linear patterns in customer behavior. The model evaluation process provides valuable visualizations, but incorporating quantitative measures and exploring advanced visualization techniques could strengthen the analysis. Profiling customer segments demonstrates insightful patterns, yet refining with advanced clustering methods tailored to the retail domain could offer a more nuanced understanding. Therefore, a critical and strategic consideration of alternative methods stands poised to elevate the project's overall efficacy and insights.

## **DATA SET 2 : DISASTER TWEETS – WORD2VEC AND TF-IDF (MAXIM KNYAZEV)**

### **1. Objective of the project:**

The primary objective of this project is to explore and implement quantitative methods and machine learning models for text classification, specifically focusing on disaster and non-disaster tweet categorization. Overall, the project successfully addresses its objective by systematically applying various techniques, such as text pre-processing, TF-IDF vectorization, Word Embedding with Word2Vec, and model building with evaluation. Including a section to **evaluate the model's predictions through tools such as confusion matrices or ROC-AUC** could have enhanced the understanding of the model's accuracy on the test data.

### **2. EDA - Text Pre-processing: (Code Line 12)**

In the text pre-processing section, the code adopts a basic strategy by using **Gensim's predefined STOPWORDS** to remove common stopwords, aiming to reduce dimensionality in natural language processing tasks. While this method is straightforward, it has a **potential limitation in neglecting domain-specific terms**. However, the project acknowledges this and attempts to address it by creating a **custom list**. This adaptation seeks to ensure that domain-specific terms are not overlooked during the pre-processing stage, adding a layer of customization to the process. The effectiveness of this solution is highlighted as it allows the project to tailor the text pre-processing step to the dataset's specific nuances. The custom list is presented as a pragmatic solution, enhancing the overall robustness of the pre-processing phase. This approach aligns with the project's objective of achieving reasonably accurate and contextually relevant text representations for subsequent analysis. **An additional consideration for robust text pre-processing involves handling other forms of noise in the data, such as special characters or URLs. The code, as provided, doesn't explicitly address these elements.**

### **3. Term Frequency-Inverse Document Frequency - TF-IDF: (Code Line 35)**

TF-IDF is a numerical statistic reflecting a word's importance to a document in a collection or corpus, calculated based on its frequency in a document compared to the entire corpus. In the TF-IDF Vectorization section, the code utilizes the TfidfVectorizer from scikit-learn to convert text data into numerical features. This method captures the significance of words in individual documents relative to their occurrence in the entire corpus.

However, it's crucial to note a **limitation associated with the TF-IDF approach—the assumption that each term is independent of others. TF-IDF inherently treats terms as isolated units, neglecting semantic relationships and word order**. This limitation might impact the model's understanding of context, as it doesn't consider the nuanced interplay between words in a sequence.

While TF-IDF is effective for certain tasks, particularly when the importance of individual terms is paramount, its use might lead to suboptimal results in tasks where contextual information and semantic relationships play a significant role. **This critique doesn't negate the utility of TF-IDF but encourages consideration of alternative approaches like Word Embeddings or transformer-based models for tasks where capturing context is crucial.**

### **4. Model Building and Evaluation: (Code Line 39)**

The section presents various machine learning models like **Support Vector Machines (SVM), Logistic Regression, K-Nearest Neighbors (KNN), Multinomial Naive Bayes, Decision Tree, and Random Forest** for text classification. A more comprehensive critique could consider **potential class imbalances in the dataset, impacting performance metrics**. It's essential to discuss strategies, if any, to mitigate these imbalances, ensuring a fair evaluation across classes.

While default settings offer a quick baseline, a deeper critique should explore how models perform under different hyperparameter configurations. Though **exhaustive tuning may be computationally intensive**, experimenting with

key configurations for each model can provide valuable insights into their sensitivity and potential areas for improvement.

The project wisely opts for GridSearchCV for Logistic Regression but should extend this approach, even in a limited fashion, to other models. This pragmatic strategy can uncover performance gains without an exhaustive search. In summary, addressing class imbalances and exploring varied hyperparameter configurations would enhance the understanding of model strengths and limitations in text classification.

## 5. Word Embedding with Word2Vec section: (Code Line 63)

In the Word Embedding with Word2Vec section, the textual data undergoes a transformation into numerical vectors using the Word2Vec technique. This involves leveraging the `en_core_web_lg` model from spaCy for vectorization, followed by normalization to ensure consistent scaling. The normalized Word2Vec vectors are then employed to train and evaluate a set of machine learning models, including Support Vector Machines, Logistic Regression, K-Nearest Neighbors, Decision Tree, and Random Forest.

This section explores the semantic relationships embedded in Word2Vec representations, potentially offering more nuanced and context-aware features compared to traditional TF-IDF vectorization. However, the effectiveness of Word2Vec relies on the quality and relevance of pre-trained word vectors, and normalization mitigates biases toward specific dimensions during model training.

## 6. Hyperparameter Tuning: (Code Line 82)

In the realm of Hyperparameter Tuning, the use of GridSearchCV for Logistic Regression reflects an approach to enhancing model performance. The method systematically explores a hyperparameter grid, providing an evaluation through cross-validation. While this exhaustive search can yield optimal configurations, it comes at the cost of increased computational intensity, particularly with larger parameter grids. The code snippet uses scikit-learn's GridSearchCV, streamlining the tuning process.

Alternative methods, such as RandomizedSearchCV, offer efficiency gains by sampling a subset of the hyperparameter space, potentially achieving similar performance with less computational burden. Bayesian optimization, a more sophisticated alternative, adapts the search space based on previous evaluations, showing promise in finding optimal configurations with fewer iterations.

Critically, the effectiveness of hyperparameter tuning hinges on the judicious selection of the hyperparameter grid. A poorly chosen grid might lead to suboptimal results, highlighting the importance of domain knowledge and experimentation. Moreover, the computational resources available play a role in determining the feasibility of exhaustive searches.

This section emphasizes the broader significance of quantitative methods in refining models for improved predictive accuracy. It underscores a trade-off between computational intensity and optimization benefits, encouraging practitioners to consider their specific constraints and goals.

## 7. Test Set Prediction phase: (Code Line 41/44/47/50/53/56)

In the Test Set Prediction phase, the application of the trained Logistic Regression model with optimized parameters on the test set is a step in evaluating the model's predictive performance. This process relies on statistical inference, specifically assessing the model's ability to generalize from the training data to new, unseen instances. The utilization of the `predict` method involves translating the learned patterns into predictions based on the test set features.

Quantitatively, this phase emphasizes metrics such as accuracy, precision, recall, and F1 score, essential statistical measures to gauge the model's classification performance. The saved results in the CSV file serve as a tangible output for conducting further statistical analyses and visualizations, allowing for a deeper understanding of the model's behavior on the test set.

Alternative approaches to this step may involve ensemble methods, boosting techniques, or other classification algorithms, each with its own set of quantitative metrics for evaluation. Weighing these alternatives requires a nuanced consideration of the dataset characteristics, computational resources, and the desired trade-offs between interpretability and predictive accuracy. Exploring diverse algorithms and tuning hyperparameters provides a quantitative basis for selecting the most suitable model for the specific application context.

## 8. Check Predictions:

The "Check Predictions" section should be incorporated at the end to gauge model's performance by loading and examining saved predictions from the test set by comparing predicted target values to the actual distribution.

Confusion matrices or ROC-AUC may be used to provide detailed insights into true/false positives/negatives and visualize trade-offs, which is beneficial in scenarios with significant consequences for misclassifications.

## 9. Conclusion:

In conclusion, this project addresses text classification using basic quantitative methods and machine learning models. The strengths lie in the straightforward text preprocessing and the use of TF-IDF vectorization, providing a foundational approach. The array of machine learning models in the Model Building and Evaluation section offers a simple baseline, while the Word Embedding with Word2Vec section introduces semantic richness, underscored by careful normalization. Hyperparameter Tuning showcases an approach, acknowledging the trade-off between computational intensity and optimization benefits. The Test Set Prediction phase emphasizes statistical metrics for evaluation, yet the simplicity of the Check Predictions section prompts consideration of alternatives. While the project provides a basic understanding of text classification, it encourages exploration of additional metrics and visualization methods, acknowledging the intricacies of the field and the continual pursuit of refinement in both methodology and results.



## DATA SET 3 : CRITIQUE OF THE IEEE FRAUD DETECTION COMPETITION KERNEL

### 1. Project Objectives:

Explore and analyze the IEEE Fraud Detection dataset to improve fraudulent transaction alerts, helping businesses reduce fraud loss and increase revenue. Successful exploration of features, visualization of key insights (e.g., time series split, distribution of target, transaction amount analysis, categorical features), and identification of potential patterns (e.g., higher transaction amounts for fraudulent charges). Ongoing updates and clear, reusable code for broader dataset analysis. Overall, the project is on track, providing valuable insights for fraud detection. However, there is room for improvement, particularly in areas such as in-depth feature engineering discussions, statistical testing for time series analysis, and detailed interpretation of certain features.

### 2. Introduction:

The provided notebook aims to conduct exploratory data analysis (EDA) on the IEEE Fraud Detection dataset, focusing on the prediction of online transaction fraud. The critique will be organized into three main sections: Time Series Analysis, Categorical Features, and Numeric Features.

### 3. Time Series Analysis:

Pros:

- **Clear Temporal Split:** The analysis correctly identifies and visualizes the temporal split between the training and test datasets using the "TransactionDT" feature. *(Code Line 6)*
- **Visualization of Transaction Amounts:** The scatter plots effectively illustrate transaction amounts over time, highlighting potential patterns related to fraud. *(Code Line 7)*

Cons:

- **Limited Feature Engineering Discussion:** The notebook lacks an in-depth discussion on feature engineering techniques related to time series analysis, which could enhance model performance.
- **Missing Statistical Testing:** The absence of statistical testing, such as hypothesis tests or ANOVA, limits the robustness of the time series analysis.

Recommendations:

- **Feature Engineering Exploration:** Discuss potential feature engineering techniques specific to time series data, such as lag features, rolling statistics, or time-based aggregations.
- **Statistical Testing:** Integrate statistical tests to validate observed patterns in transaction amounts over time.

### 4. Categorical Features:

Pros:

- **Categorical Feature Analysis:** The notebook effectively explores categorical features like "ProductCD" and "card" columns, providing insights into their distribution and fraud percentages. *(Code Line 13)*
- **Visualizations for DeviceType and DeviceInfo:** The visualizations for "DeviceType" and "DeviceInfo" offer a clear understanding of the distribution and potential relationships with fraud. *(Code Line 33)*

Cons:

- **Limited Interpretation of Card Columns:** While histograms for "card" columns are presented, there is a lack of detailed interpretation, making it challenging to extract meaningful insights.
- **Address Features Ambiguity:** The ambiguity in interpreting "addr1" and "addr2" as categorical features is acknowledged but not further explored or clarified.

Recommendations:

- **Detailed Interpretation for Card Columns:** Provide additional interpretation and analysis for the "card" columns, addressing challenges in feature engineering.
- **Clarification on Address Features:** Further explore and discuss the potential role of "addr1" and "addr2" in fraud detection, clarifying their significance.

### 5. Numeric Features:

Pros:

- **Comprehensive Pair Plots:** The pair plots for numerical features like "C1-C14," "D1-D9," and "V1-V339" offer a visual exploration of feature interactions and distributions. *(Code Line 22/24/31)*
- **Log Transformation for TransactionAmt:** The log transformation of "TransactionAmt" effectively addresses skewness and aids in better visualizing its distribution. *(Code Line 11)*

#### Cons:

- **Limited Discussion on Feature Interaction:** While pair plots are presented, there is a limited discussion on the observed interactions and their potential impact on fraud prediction.
- **Sparse Exploration of Numeric Features:** The analysis briefly mentions numeric features like "C1-C14," "D1-D9," "M1-M9," and "V1-V339" having a sparse matrix structure with many empty values. It lacks depth in exploring the significance of these sparse values for fraud detection.

#### Recommendations:

- **In-Depth Analysis of Feature Interactions:** Provide a more detailed analysis of observed feature interactions, discussing their implications for fraud detection.
- **Extended Exploration of Numeric Features:** Elaborate on the exploration of specific numeric features, offering insights into their potential significance - considering strategies like feature engineering or imputation to extract meaningful patterns from sparse values.

#### 6. General Critiques:

##### a. Model Building and Evaluation:

The notebook lacks details on model building and evaluation, which are crucial components of the fraud detection task. Including information on the chosen model, training process, and evaluation metrics would enhance the completeness of the analysis.

##### b. Outlier Detection and Data Imbalance:

There is no explicit discussion or handling of outliers, which can significantly impact model performance. Additionally, addressing the issue of data imbalance, especially in the context of fraud detection, is essential.

##### c. Advanced Visualizations and Tools:

While the notebook uses standard visualizations, incorporating advanced visualizations like pair plots, heatmaps, and interactive plots could provide a deeper understanding of complex relationships. Moreover, leveraging automated EDA tools like pandas-profiling or Sweetviz could streamline the analysis.

##### d. Cross-Validation and Model Complexity:

The recommendation to implement cross-validation techniques for a more reliable estimate of model performance is valuable. Moreover, exploring more sophisticated models beyond simple linear regression, such as ensemble methods or gradient boosting, could improve predictive performance.

##### e. Addressing Limitations:

The notebook mentions limitations but does not provide concrete steps to address them. Specifically, incorporating hypothesis testing, advanced statistical techniques, and addressing the imbalance issue would strengthen the analysis.

#### 7. Conclusion:

In conclusion, the notebook lays a sturdy foundation for exploratory data analysis (EDA) on the IEEE Fraud Detection dataset. While it successfully navigates through the temporal intricacies, categorical nuances, and numeric complexities, there exists untapped potential for additional depth and clarity. The critique identifies commendable aspects, such as the adept visualization of temporal splits and transaction amounts, along with comprehensive pair plots for numeric features. However, the analysis falls short in certain areas, notably in-depth discussions on feature engineering, statistical testing for time series patterns, and detailed interpretations of categorical and numeric components.

The general critiques shed light on pivotal aspects like model building, outlier detection, data imbalance, advanced visualizations, cross-validation techniques, and the proactive addressing of limitations. In essence, while the notebook serves as a commendable starting point, embracing the outlined recommendations would undoubtedly elevate its standing, fostering a more comprehensive and insightful analysis of the IEEE Fraud Detection dataset.