

Individual Assignment Data Science Tools and Techniques:
Sales of Ecommerce Amazon



Sophie Hu (1657847)

John Smits

19 December 2023

Table of Contents

Introduction.....	3
CRISP-DM Model	4
Business understanding	4
Data understanding.....	4
Data Preparation	6
Amazon Sales Report	7
International Sales Report	8
Sales Report.....	12
Modelling	13
Market basket analysis of Amazon Sales Report.....	13
K-Means clustering of International Sales Report	16
K-Means clustering of Sales Report.....	19
Evaluation.....	21
Market basket model of Amazon Sales Report	21
K-Means Clustering of International Sales Report.....	22
K-Means Clustering of Sales Report.....	23
Deployment	24
Market Basket Analysis of Amazon Sales Report.....	24
K-Means Clustering of International Sales Report.....	25
K-Means Clustering of Sales Report with Stock Focus	26
Final deployment	27
Bibliography	28

Introduction

This project looks into Amazon's e-commerce world using three specific datasets. Each dataset provides unique insights into different aspects of Amazon's operations. By carefully examining the Sales Report, International Sales Report, and Amazon Sales Report, we aim to discover patterns, trends, and connections that help us understand how Amazon's marketplace functions. While we'll analyze each report on its own, our ultimate goal is to combine these insights.

Sales Report

The sales report dataset provides detailed information about product stock levels. Each entry includes a unique identifier (SKU), product name, quantity available, size, and color.

Utilizing data mining techniques like K-means clustering on this dataset can reveal patterns and insights that contribute to more effective inventory management strategies. For instance, clustering similar products based on size, color, and quantity can assist in optimizing stock placement and replenishment strategies.

International Sales Report

The international sales report dataset captures information on sales transactions in a global context. It includes data on the date, customer details, product style, SKU, size, quantity, rate, and gross amount. Applying K-means clustering to this dataset can reveal patterns and insights that contribute to a deeper understanding of international sales dynamics.

Amazon Sales Report

The Amazon sales report dataset focuses on sales transactions and customer behavior within the Amazon platform. It contains information such as order ID, date, status, fulfillment details, sales channel, style, SKU, category, size, ASIN, and amount. Market basket analysis is a valuable technique for this dataset, revealing associations between products frequently bought together. This analysis provides insights into customer preferences and buying patterns unique to the Amazon platform. For example, understanding which products are often purchased together can inform cross-selling and bundling strategies for better sales performance on Amazon.

Within the spectrum of machine learning and econometric modeling techniques explored, it was observed that K-Means Clustering and Market Basket Analysis outperformed others in forecasting Amazon's sales pattern analysis. Both K-Means Clustering and Market Basket Analysis are widely acknowledged and extensively utilized models renowned for uncovering intricate patterns and revealing associations between products, the dependent variable (sales), and the chosen independent variables. These models yield invaluable insights into the influence of these variables on Amazon's sales performance, facilitating precise predictions of sales outcomes. Consequently, the selection of K-Means Clustering and Market Basket Analysis models stands as the optimal choice for this analysis, providing a robust foundation for forecasting Amazon's sales patterns and offering stakeholders valuable insights into the company's consumer prospects.

CRISP-DM Model

Business understanding

Understanding sales patterns, customer behavior, and purchasing trends is crucial in the realm of e-commerce, especially for a giant like Amazon. Deciphering these patterns is essential as it provides insights into the preferences and demands of customers, enabling strategic decision-making.

In this sales-focused research, we utilize advanced sales modeling methods, specifically K-means clustering and Market Basket Analysis. K-means clustering facilitates the categorization of products based on their sales patterns, optimizing inventory, and enhancing customer satisfaction. On the other hand, Market Basket Analysis reveals associations between products frequently purchased together, enabling targeted recommendations and personalized shopping experiences. These methods are invaluable for e-commerce platforms as they improve product placements, streamline inventory, and ultimately elevate the customer experience. The impact extends beyond the company itself, influencing not only business strategies but also shaping the online shopping experience for millions of consumers.

Data understanding

We carefully examined the existing data to understand its details, quality, and relevance to our goals. The dataset mainly comes from Sharma, A. (2022) "Ecommerce Sales dataset," focusing on key aspects of sales. This dataset offers a thorough view of "e-commerce sales," covering product details, stock levels, and pricing across different stores. It provides insights into sales dynamics, making it a user-friendly resource for exploring e-commerce intricacies. For creating a Sales pattern model for Amazon, we used both quantitative and qualitative research methods. The dataset includes all necessary Sales indicators from various sources, totaling 29 variables. Each source's data is included, giving complete information about the authors and origins of each variable.

Our goal is to build a robust Sales pattern model for Amazon by carefully examining the dataset, addressing missing values, and considering important variables. Given the dataset's size, we expected some missing values. Specifically, the Amazon Sales Report dataset shows ten variables with missing values, as depicted in Figure 1. To handle this, we used a data deletion strategy to ensure dataset completeness. The approach, guided by reliable sources like the lesson titled "Effective Data Cleaning Strategies in Python for Exploratory Data Analysis (EDA)" published on Medium by Python Fundamentals (2023), was implemented. This lesson provides practical insights into addressing missing data. By incorporating best practices from trusted sources such as Columbia University, our research ensures reliable approaches to maintain dataset integrity.

Figure 1
Missing values of Amazon Sales Report

Variable names	Missing values
Courier Status	6872
currency	7795
amount	7795
ship-city	33
ship-state	33
ship-postal-code	33
ship-country	33
promotion-ids	49153
fulfilled-by	89698
Unnamed:22	49050

When addressing missing values in a dataset, our goal is to retain as much data as possible to preserve potential insights, all while maintaining the integrity and quality of the analysis. Simply removing rows or columns with missing values can result in a substantial loss of information, especially in large datasets or when the missingness is non-random. Instead, we take a careful approach, examining each column with missing values to determine the most suitable action.

In our commitment to upholding data integrity and reliability, we adopted a robust strategy for handling missing values. Rows containing incomplete data were systematically removed from the dataset using the "dropna()" function. Although this led to a slight reduction in the dataset's size, it was a necessary step to ensure the accuracy and trustworthiness of our analysis. The result of this thorough data cleansing process was the creation of a new dataset, which we named. This cleaned dataset became the cornerstone for our subsequent analysis, guaranteeing the validity of our findings and minimizing potential errors stemming from missing data.

Figure 2
Missing values of International Sales Report

Variable names	Missing values
index	0
date	1
months	25
customer	1040
style	1040
sku	2474
size	1040
pcs	1040
rate	1040
gross amt	1040

Figure 3
Missing values of Sales Report

Variable names	Missing values
index	0
date	1
months	25
customer	1040
style	1040
sku	2474
size	1040
pcs	1040
rate	1040
gross amt	1040

Data Preparation

Standardization is a crucial preprocessing step in the analysis of the Amazon Sales Report and International Sales Report datasets. By standardizing the data, we ensure that variables with different scales and units are transformed to a common scale, eliminating biases that might arise due to variations in measurement units. This is particularly pertinent in our analysis, where diverse metrics such as product quantities, amounts, and currency values coexist. Standardization allows us to compare these metrics on a level playing field, facilitating the identification of meaningful patterns and relationships without the undue influence of disparate scales.

Moreover, the standardization process is essential when employing machine learning algorithms like K-means clustering or logistic regression. These algorithms often rely on distance metrics or gradient-based optimization, both of which can be adversely affected by unscaled data. By bringing all features to a standard scale, we enhance the performance and convergence of these algorithms, ensuring that the impact of each variable is appropriately considered in the analysis. This becomes particularly pertinent in the International Sales Report, where the dataset encompasses a wide range of variables, from gross amounts to quantities sold, demanding a harmonized scale for meaningful model outcomes.

In addition to algorithmic considerations, standardization simplifies the interpretation of results. In our analysis, where the emphasis is on drawing actionable insights for inventory management or sales strategy, having variables on a common scale makes it easier to comprehend the relative importance of each factor. This is crucial when making informed decisions based on the analysis, as stakeholders can more confidently assess the impact of different variables and devise strategies that align with the overarching goals of the e-commerce operation.

Amazon Sales Report

As can be observed in figure 4, the 'Unnamed: 22' column exhibits only two distinct values: NaN and False. Considering that the non-NaN values are all False, a potential approach is to interpret NaNs as True, effectively transforming it into a valid boolean column. Nevertheless, given the column's lack of a clear name, description, or purpose, it appears unlikely to contribute valuable insights to our analysis. Given these considerations, the reasonable decision is to exclude this column from our dataset.

Figure 4
Unique values in 'unnamed:22'

Variable names	Unique values
unnamed: 22	nan
	False

In figure 5, the 'courier status' column signifies the shipment's status, with 'Shipped,' 'Unshipped,' and 'Cancelled' delineating distinct stages in the shipping process. In the absence of additional statuses beyond these three and NaN, it implies that the missing values (NaN) might indicate an unrecorded or unknown status during data collection. Consequently, we have transformed the NaN values into 'Unknown' to appropriately reflect this uncertainty in the dataset.

Figure 5
Column 'Courier status'

Unique values	Count
Shipped	109487
NaN	6872
Unshipped	6681
Cancelled	5935

Examining the 'Amazon.in' sales channel in figure 6, it is evident that the currency is predominantly INR, except for a few missing values. To address this, we can reliably impute these missing values with 'INR'. Conversely, for the 'Non-Amazon' channel, the absence of any currency information indicates that marking these instances as 'Unknown' is a suitable approach.

Figure 6
Column 'sales channel'

sales channel	currency	Count
Amazon.in	INR	121180
	NaN	7671
Non-Amazon	NaN	124

Upon examining rows with missing shipping details, it becomes apparent that the other information in these rows is complete. This hints at a potential issue with data entry or collection, specifically related to shipping information. Considering that these records hold valuable information apart from the missing shipping details, and given that the missing data in this context is minimal, it is advisable to retain these rows for analysis. To maintain the dataset's integrity, we can appropriately label the missing shipping details as 'Unknown,' thereby ensuring a balanced approach without making unwarranted assumptions about the nature of the missing data.

Moreover, column 'ship_postal_code' has been converted to string while preserving leading zeros. In many datasets, especially when dealing with postal codes, the values may have leading zeros. Postal codes are often formatted with a fixed number of digits, and if a postal code starts with zeros, these leading zeros can be important for correct interpretation and sorting.

International Sales Report

The purpose of this matrix is to visually represent the patterns of missing values in the dataset. The matrix has rows representing rows in the DataFrame and columns representing columns. White lines indicate missing values, while dark lines represent non-missing values. Observing figure 7, there seems to be a notable concentration of missing values in the center of matrix/total rows. The resulting matrix helps quickly identify the distribution and patterns of missing data across different columns in the DataFrame, which is be valuable for data exploration and preprocessing.

Figure 7
Missing value matrix

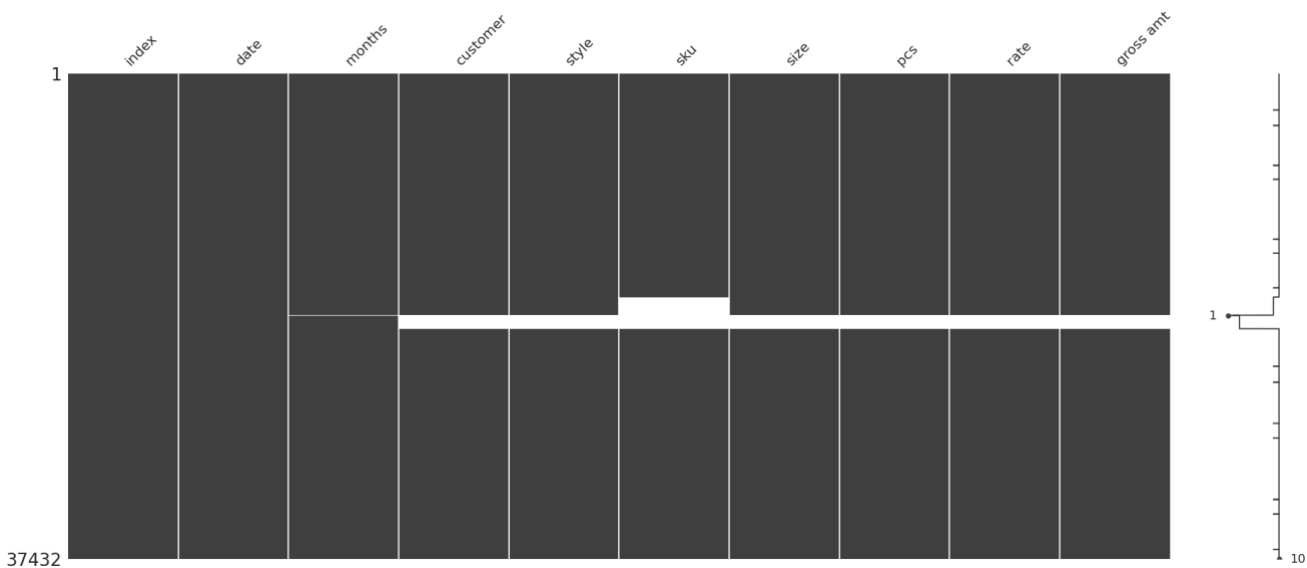


Figure 8
Total values per variable

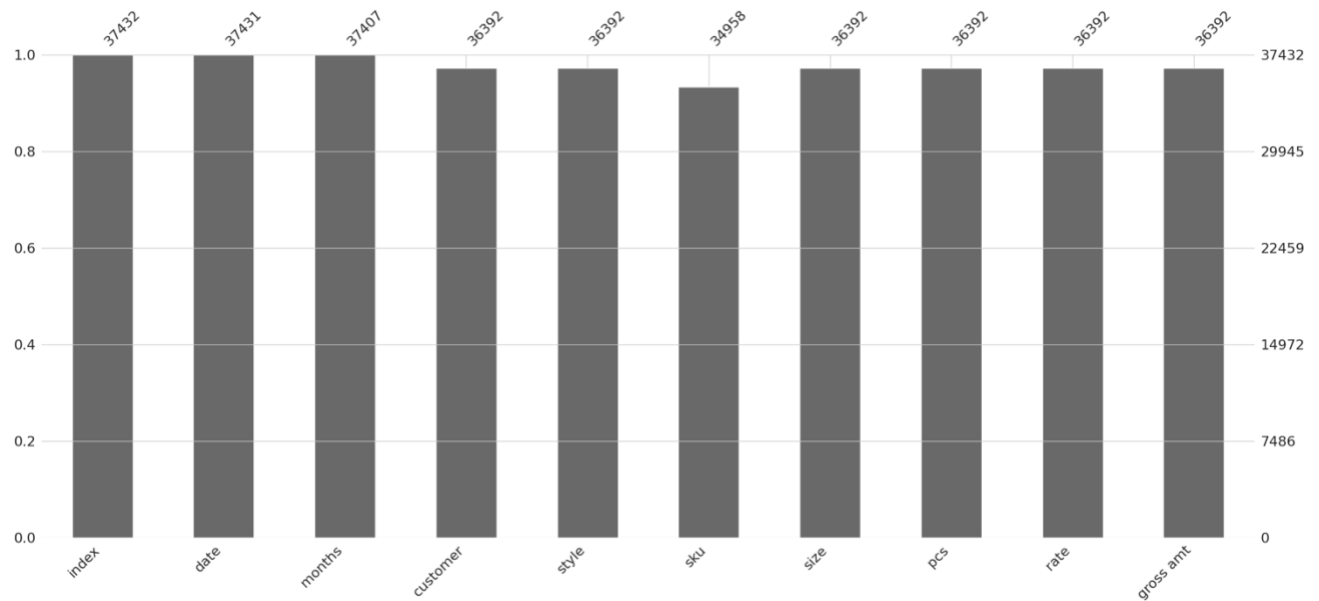


Figure 9 visually represents the distribution of missing values across different columns in the dataset. The x-axis shows column names, while the y-axis displays the percentage of missing values for each column. The red dashed line at 30% serves as a threshold. It acts as a reference point to highlight columns that exceed a predetermined level of missing values. Columns with null values percentage above 30% are emphasized, indicating a point at which the missing data becomes substantial. This threshold guides decision-making in data preparation, helping identify columns that may require further investigation, imputation, or removal based on the analysis's specific requirements and tolerance for missing data. In addition, the specific null percentages per value are displayed in figure 10.

Figure 9
Null values percentage per value

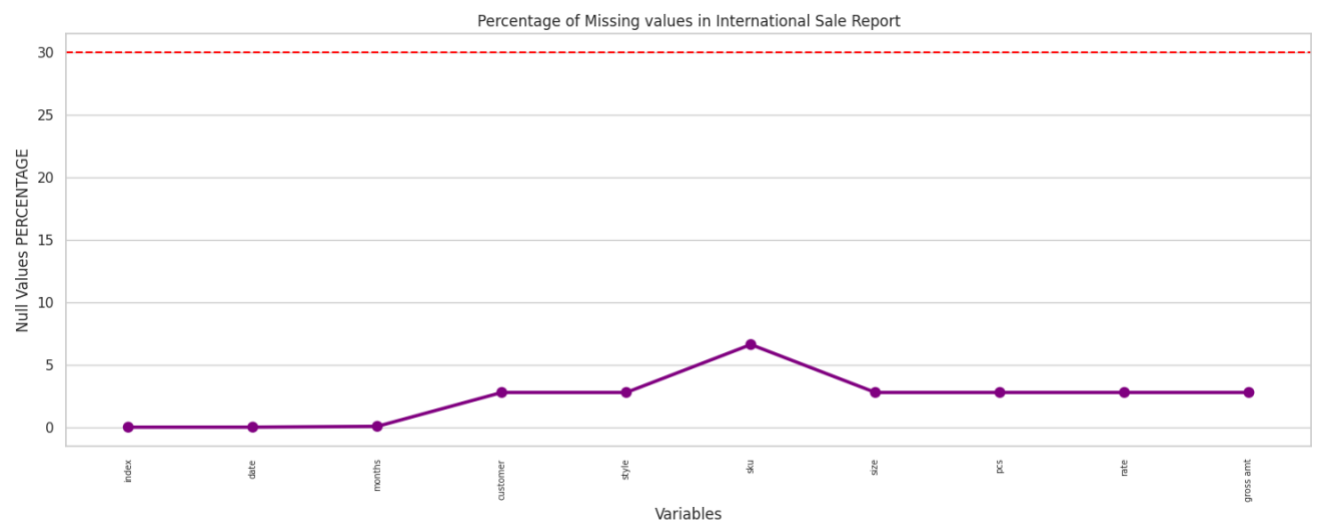


Figure 10
Null values percentage per variable

Variable names	Null values percentage
index	0,00%
date	0,00%
months	0,07%
customer	2,78%
style	2,78%
sku	6,61%
size	2,78%
pcs	2,78%
rate	2,78%
gross amt	2,78%

To address missing values, we opted to uniformly impute them with the label 'Unknown,' given the predominantly categorical and text-based nature of the variables in question. The choice of 'Unknown' as a placeholder serves as a generic indicator for the absence of specific information, particularly considering that only 0-6% of the variables contain missing values. This approach facilitates clear differentiation between missing and observed values during subsequent analyses. Moreover, it ensures the retention of all available information throughout the data preparation process. Opting against the removal of rows with missing values prevents potential loss of valuable observations, safeguarding against unintended consequences such as bias introduction or a reduction in the overall sample size, which could impact the dataset's representativeness in our analysis.

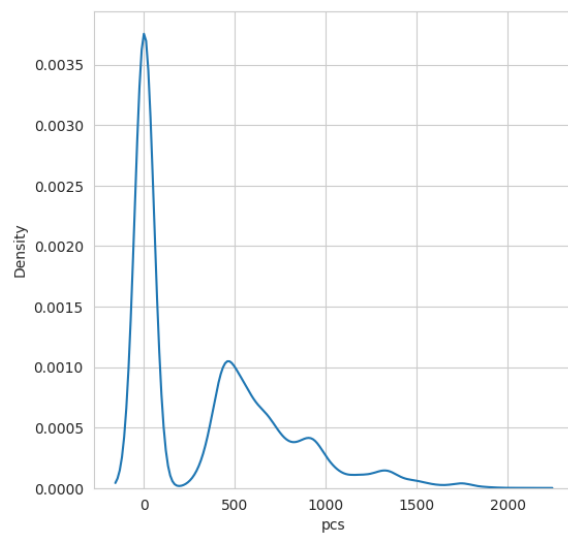
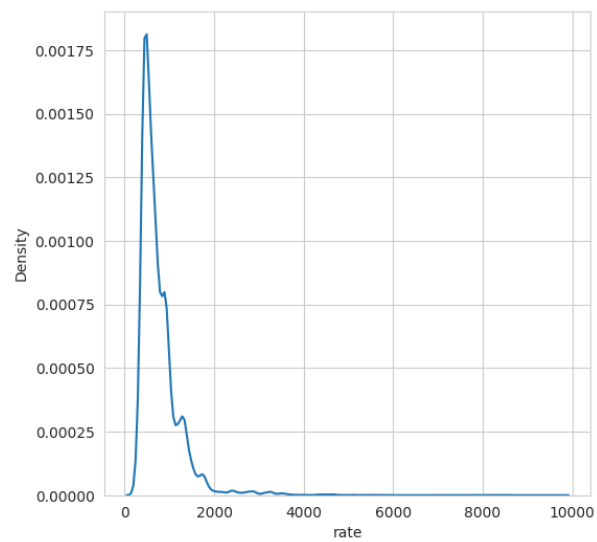
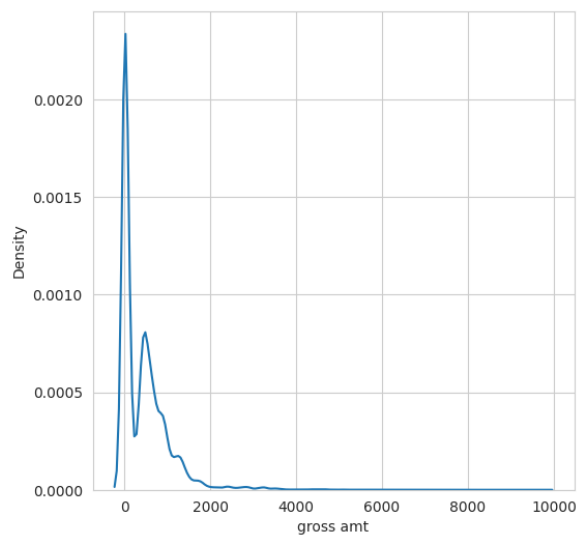
The purpose of the KDE plots in figure 23 is to visually inspect the distribution of data in each specified column. Our KDE plots provide a smooth estimate of the probability density function of a continuous random variable, offering insights into the data's underlying distribution. The focus is on identifying potential outliers in the 'rate', 'pcs', and 'gross amt' columns.

Analyzing the shape of these distributions helps identifying any irregularities, such as data points that deviate significantly from the overall pattern, indicating potential outliers. This visual exploration is a crucial step in our data analysis, providing a preliminary understanding of the data's characteristics before further statistical or machine learning techniques are applied.

Our kernel density estimation (KDE) plots for the 'rate', 'pcs' (pieces), and 'gross amt' (gross amount) columns in the provided dataset all exhibit distributions with densities close to zero. This suggests that the majority of data points are concentrated within a specific range, and there are no significant peaks or outliers distorting the overall pattern.

The KDE plots serves as a visual representation of the data distribution, allowing for a quick assessment of its shape and characteristics. In our case, the plots indicate that the values in the 'rate', 'pcs', and 'gross amt' columns are relatively well-behaved, with no apparent skewness or extreme deviations.

Figure 23
KDS plots to detect outliers



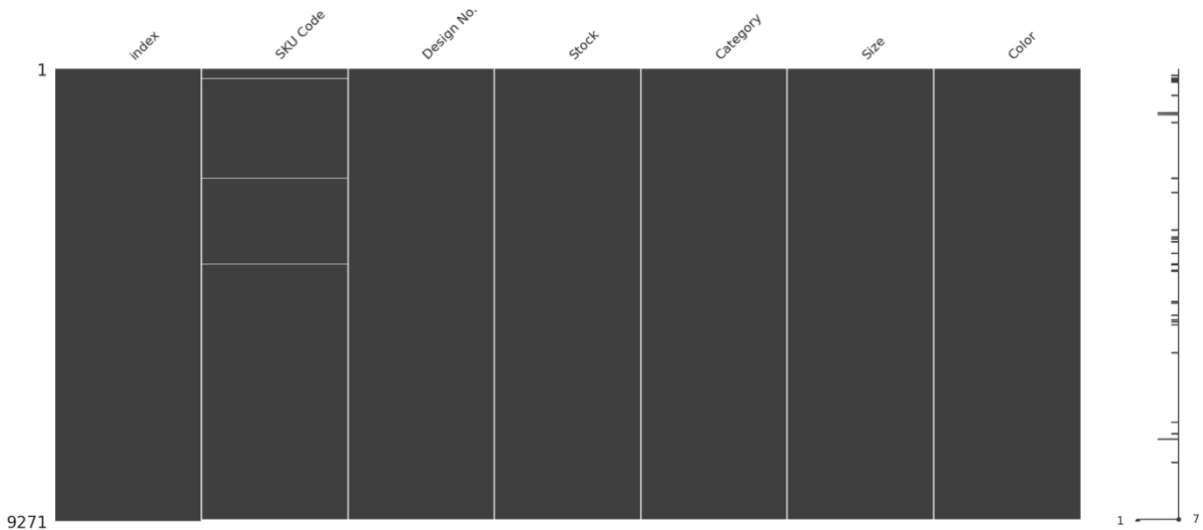
Sales Report

In the 'Sales Report' dataset, it was noted that some complete rows at the end of the dataset contained missing values, offering no meaningful contribution to our analysis. Consequently, we opted to eliminate these rows. Figure 11 illustrates the count of missing values per variable after the removal of complete empty rows, providing a clearer and more transparent depiction of missing values that genuinely contribute to our dataset. Additionally, 47 instances of the value '!#REF' were replaced with 'Unknown.' Subsequently, the remaining scant missing values were excluded from our analysis. This decision was informed by the observation that rows with missing values lacked substantial data across multiple columns, rendering their removal inconsequential to our analytical objectives.

Figure 11
Missing values per variable after removing empty rows

Variable names	Missing values
index	0
SKU Code	49
Design No.	2
Stock	2
Category	11
Size	2
Color	11

Figure 12
Missing value matrix



Modelling

Market basket analysis of Amazon Sales Report

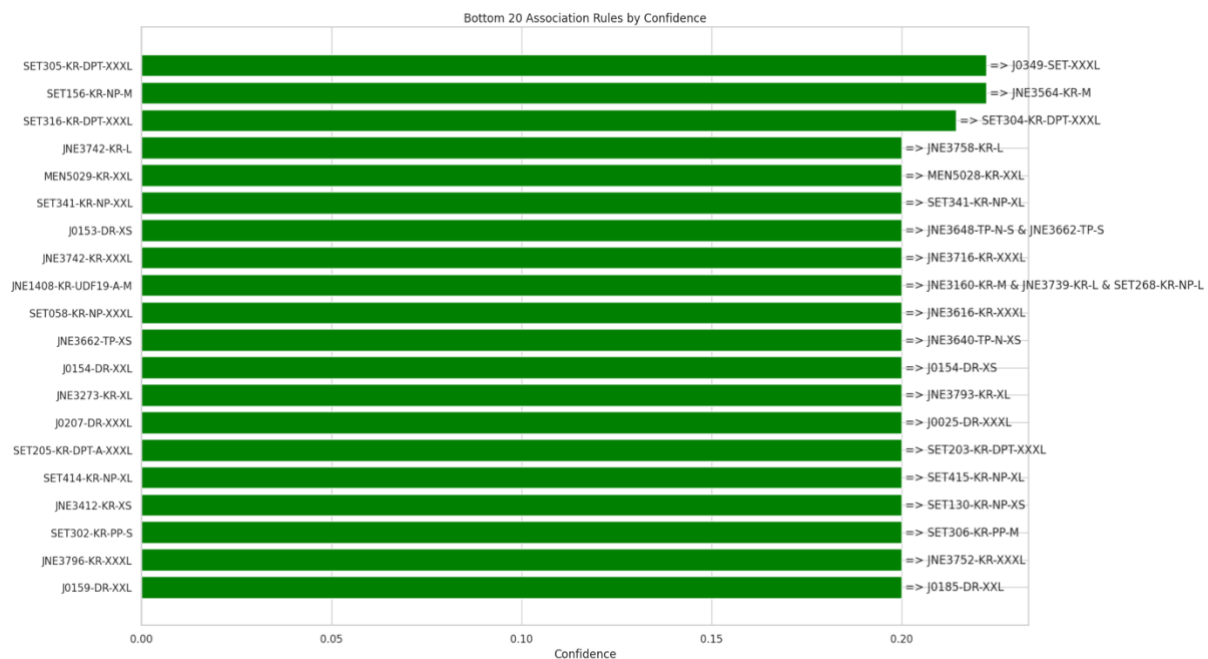
In this analysis we conducted a market basket analysis on an Amazon sale report dataset using the FP-Growth algorithm. The purpose of market basket analysis is to identify associations or patterns in customers' purchase behaviors, revealing which products are frequently bought together. This information is valuable for business strategies such as product placement, cross-selling, and personalized marketing.

The dataset is initially prepared by grouping transactions based on the 'order id' and extracting the list of SKUs (Stock Keeping Units) associated with each order. The FP-Growth algorithm is then applied to find frequent itemsets and generate association rules. The parameters we used include a minimum support threshold of 2 and a minimum confidence threshold of 0.2, controlling the significance of the discovered patterns. It visualizes association rules based on confidence scores. In the bar charts displayed below, the antecedents represent items that are frequently found together, while the consequents indicate what customers are likely to purchase given the antecedents.

In figure 13, it displays the bottom 20 association rules by confidence, providing insights into less significant patterns.

Figure 13

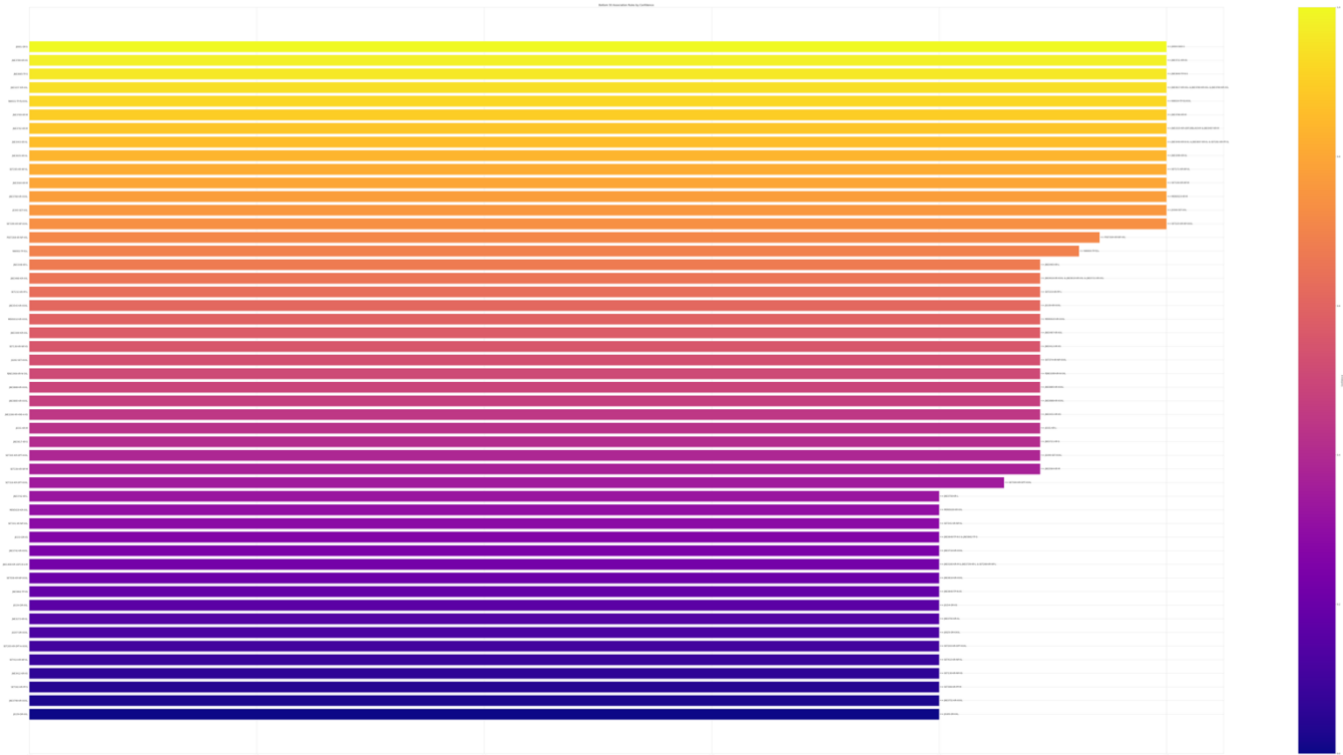
Market basket analysis of bottom 20 Association rules by confidence



In addition to visualizing the top 20 association rules by confidence, an intentional choice was made to enhance the understanding of the dataset by introducing color intensity in the visual representation of the bottom 50 rules, which is displayed in figure 14. This decision stemmed from the recognition that a mere selection of the top 20 might not sufficiently capture the diversity and complexity within the dataset. By incorporating color intensity, the visual display effectively communicates the varied confidence levels associated with each rule within the bottom 50, providing a nuanced portrayal of the numerous purchase patterns that may exist. While the text on the visualization might be challenging to discern due to the

high volume of rules, users can leverage the zoom feature on the generated image to explore specific details, ensuring a comprehensive examination of the intricate associations present in the dataset. This approach facilitates a more thorough analysis and interpretation of the market basket patterns, acknowledging the need to convey the rich diversity inherent in the lower-confidence rules. The color intensity in the charts represents the confidence level of each rule.

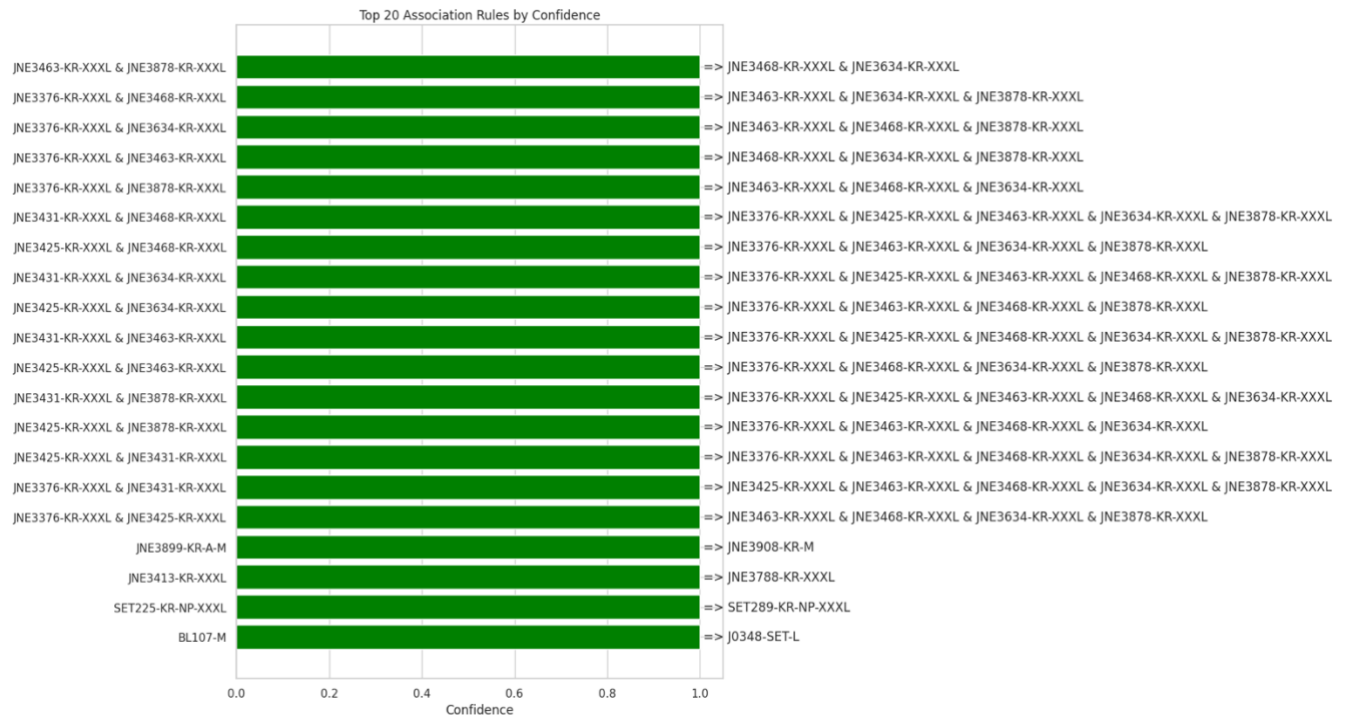
Figure 14
Market basket analysis of bottom 50 Association rules by confidence



In figure 15, it puts focus on the top 20 rules by confidence, highlighting the most significant purchase patterns.

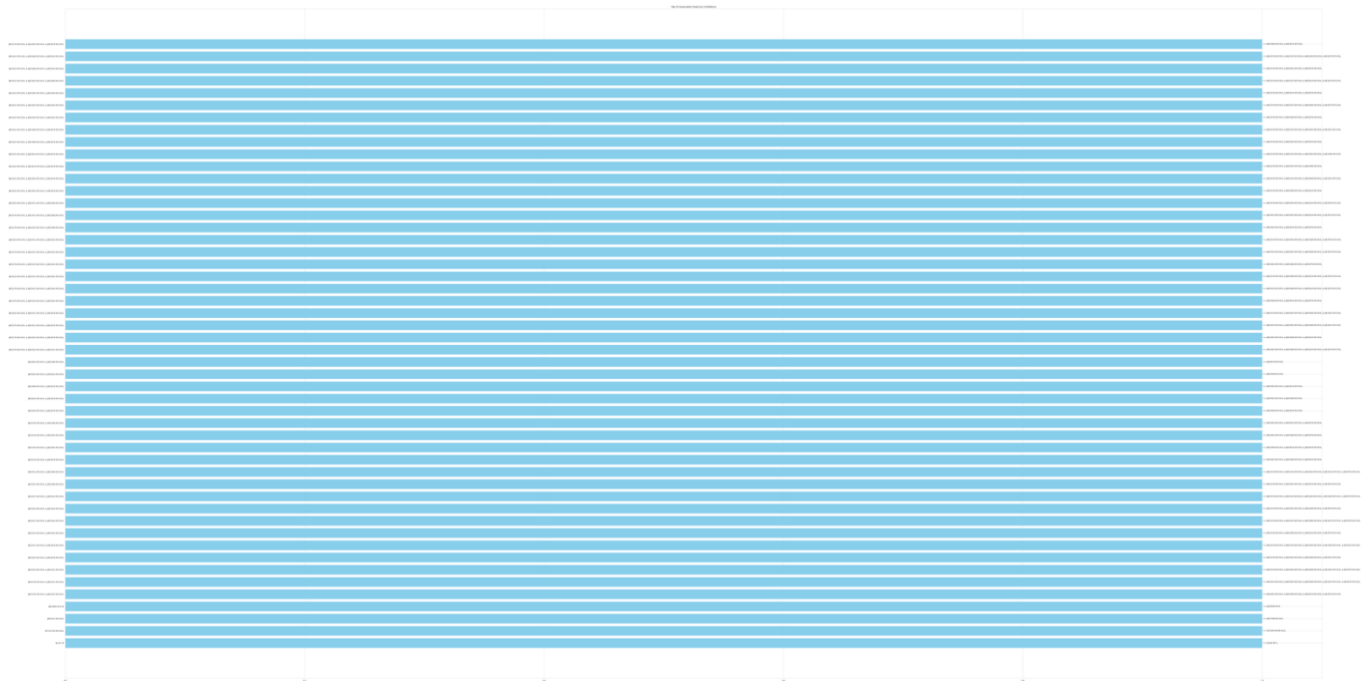
Interpreting the results, businesses can identify product combinations that are frequently bought together, enabling them to optimize product placements or design targeted promotions. For example, if customers often purchase a specific set of products simultaneously, the company might bundle these items for a discounted price to encourage further sales. Overall, our market basket analysis provides actionable insights into our customer behavior, aiding businesses in making informed decisions to enhance customer satisfaction and maximize revenue through strategic product offerings and marketing strategies.

Figure 15
Top 20 Association rules by confidence



In contrast to the nuanced color intensity applied to the bottom 50 rules, a deliberate decision was made to maintain a uniform color scheme for the top 50 association rules by confidence in figure 20. Our choice is grounded in the understanding that the top rules, particularly at least the top 200, share a confidence interval of 1, indicating a highly reliable association. The absence of color gradation in the top 50 visualization streamlines the representation, emphasizing the consistency and robustness of these prominent rules. While the clarity of individual text entries on the visualization may be compromised due to the sheer volume of rules, users are encouraged to leverage the zoom functionality in the generated image. This intentional approach aims to accentuate the significant associations in the top rules and allows for a closer examination of specific details. Overall, this visual strategy serves to underscore the smaller disparity in confidence levels across different segments of association rules, offering a comprehensive view of the diverse patterns inherent in our market basket analysis.

Figure 16
Top 50 association rules by confidence



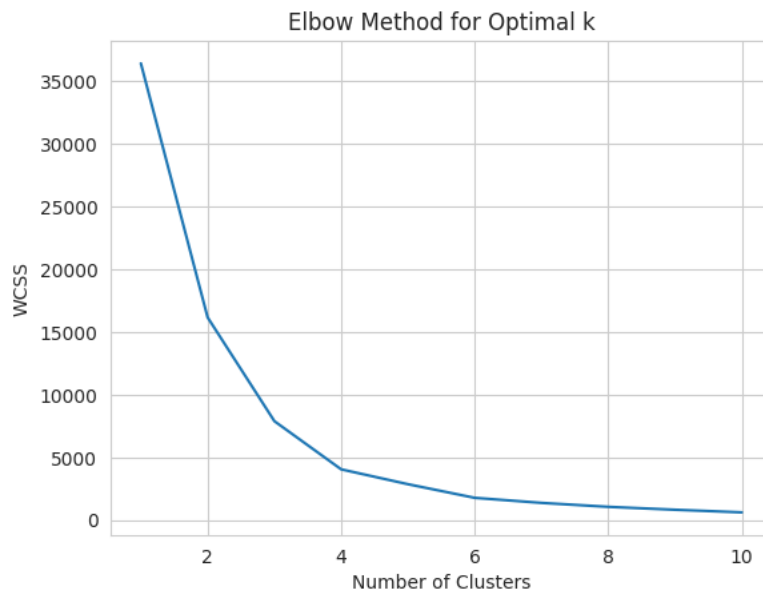
K-Means clustering of International Sales Report

In our K-Means Clustering analysis, the primary objective was to uncover inherent patterns and groupings within the dataset related to our sales (gross amt). K-Means clustering is a machine learning technique used for unsupervised learning, which means the algorithm identifies patterns in the data without the need for labeled outcomes. Our model included both categorical and numerical columns, such as months, customer details, style, SKU, size, category, rate, pieces sold (pcs), and gross amount.

To prepare our data for clustering, categorical columns were converted to categorical data types, and non-numeric values in numerical columns were handled appropriately. Negative values in numerical columns were transformed to their absolute values to ensure consistency. Additionally, string values were capitalized to maintain uniformity in the dataset. The 'gross amount' column, a key numerical feature for clustering, was standardized using the StandardScaler to bring all features to the same scale, preventing any particular feature from dominating the clustering process due to differences in magnitude.

Our following crucial step was determining the optimal number of clusters (k) for the K-Means algorithm. The Elbow Method in figure 17 was employed, plotting the within-cluster sum of squares (WCSS) against the number of clusters. The "elbow" point on the plot represents an optimal k value where adding more clusters does not significantly reduce the WCSS. For this analysis, we initially chose k value of 3. However, after generating our Elbow Method an optimal k value of 10 was selected.

Figure 17
Elbow method for optimal k



Consequently, each data point (SKU) was assigned to one of the ten clusters. The resulting clustered dataset was then printed, displaying the SKU, gross amount, and the assigned cluster for each record. The clustered data was saved to our new CSV file named 'clustered_International_sale_Report.csv.'

To deeper look into the characteristics of each cluster, we separated DataFrames for each cluster, allowing for a detailed SKU analysis within each cluster. Basic statistics, such as mean, standard deviation, minimum, 25th percentile, median, 75th percentile, and maximum, were presented for the 'gross amount' within each cluster. In figure 19, one of the ten clusters' descriptive statistics are displayed, which is cluster 1. These statistics provide insights into the distribution and variance of sales amounts within individual clusters, aiding in the interpretation and understanding of the distinct patterns identified by our K-Means Clustering algorithm.

Figure 19
Descriptive statistics of cluster 1

Descriptive statistics	SKU analysis
count	2333
mean	1285,18
std	155,122
min	1094
25%	1183
50%	1286
75%	1354
max	1522

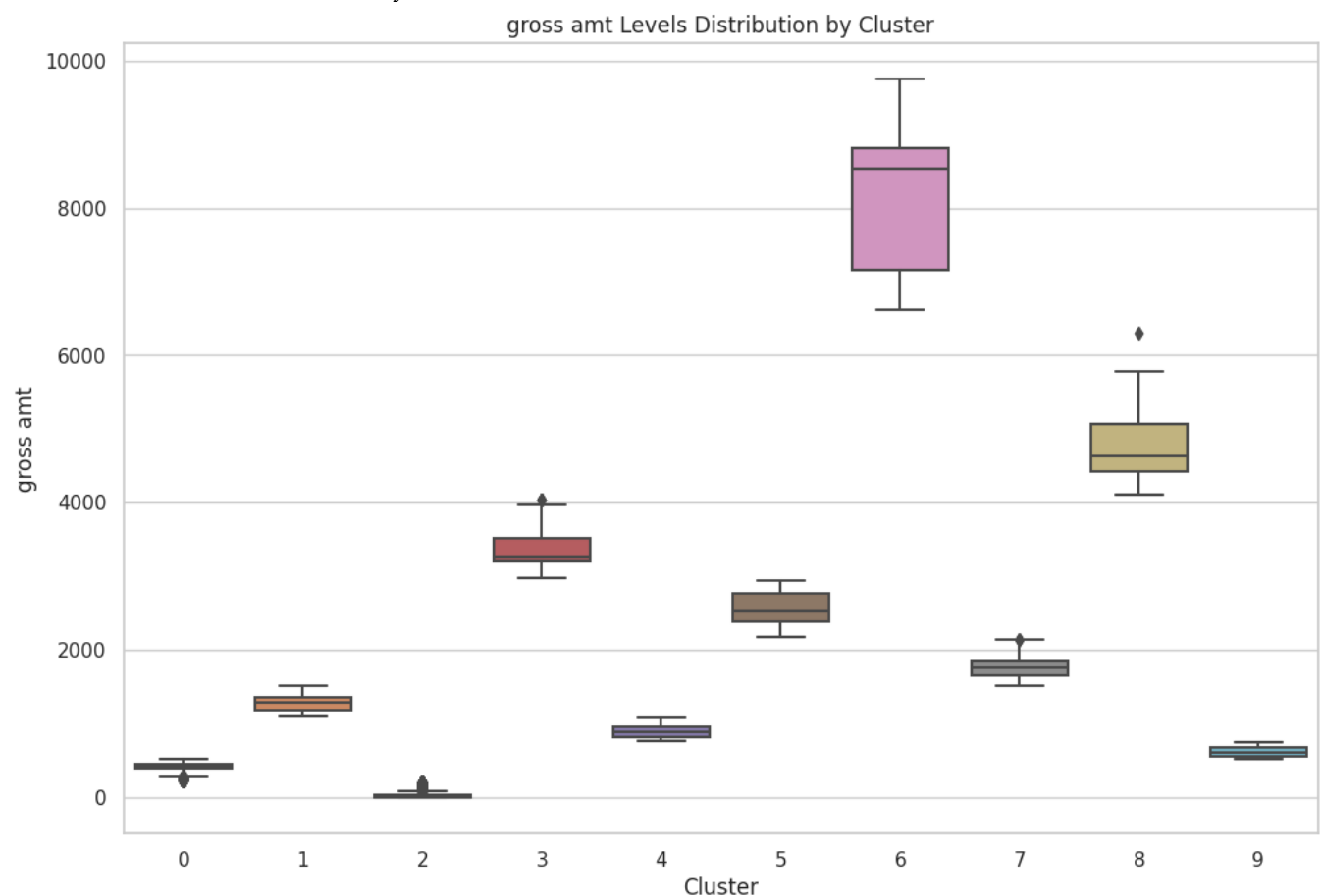
The box plot in figure 18 was employed to gain deeper insights into the distribution of 'gross amount' within each cluster. Box plots were generated for each cluster individually, providing a visual representation of the central tendency, spread, and potential outliers within the gross sales amounts of SKUs belonging to the respective clusters.

The box plot displays the median (central line within the box), the interquartile range (IQR) represented by the box itself, and whiskers extending to the minimum and maximum values within a defined range. Outliers that are presented, are depicted as individual points outside the whiskers. This visualization allowed for our comprehensive understanding of the variation in sales amounts across clusters, helping to identify any distinctive characteristics or anomalies within specific clusters.

It provides stakeholders with a visual tool to discern not only the general trends in sales amounts but also potential variations within each cluster.

Figure 18

Gross amt levels distribution by cluster



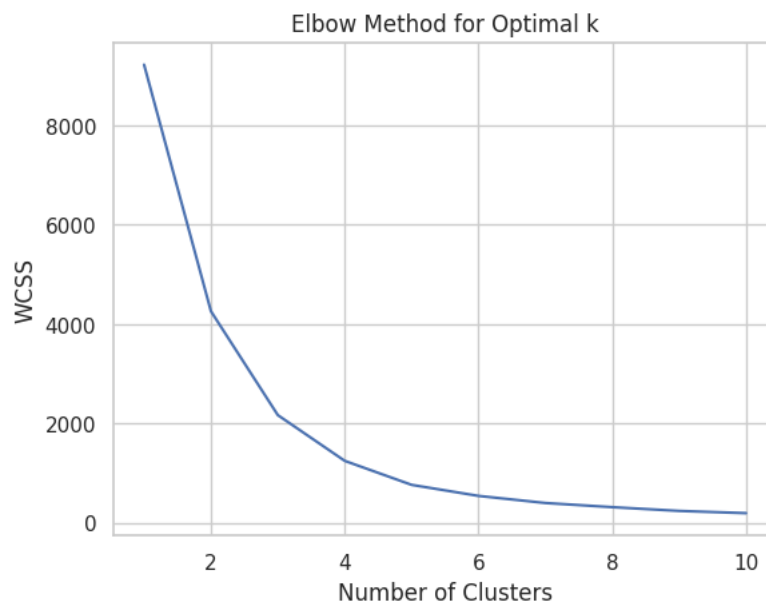
In essence, the K-Means Clustering offers a systematic way to group SKUs with similar sales patterns and provides a valuable tool for segmentation and targeted analysis of sales. It allows Amazon to discern distinct market segments, identify trends, and tailor strategies to maximize performance within each cluster. The SKU analysis within each cluster further refines this understanding by shedding light on the specific characteristics and performance metrics associated with different groups of products. Overall, K-Means Clustering serves as a powerful technique for uncovering our patterns in complex datasets and informing data-driven decision-making processes.

K-Means clustering of Sales Report

In our K-Means Clustering analysis, the primary objective was to uncover inherent patterns and groupings within the dataset related to our stock level. The focus is on the 'Stock' column, representing the quantity of each SKU in the inventory.

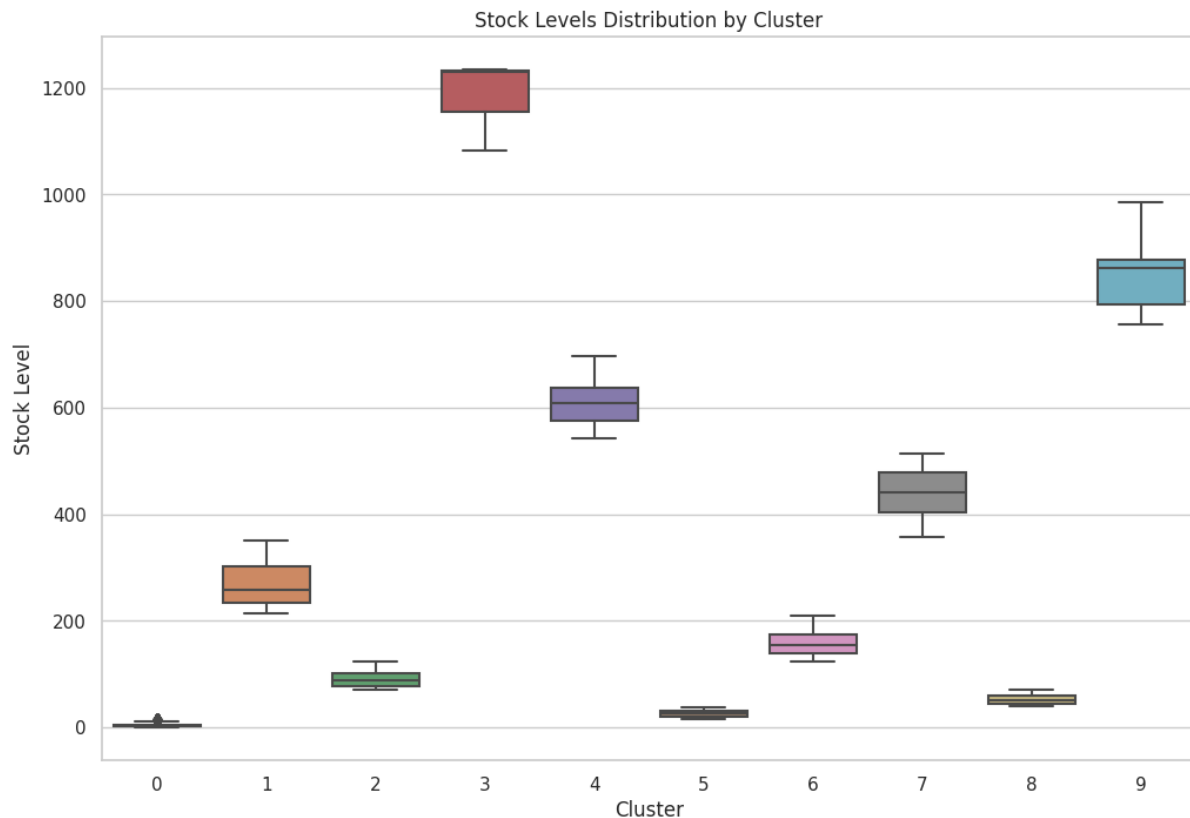
Our first step involves standardizing the numerical data using the StandardScaler from the scikit-learn library. This normalization ensures that features with different scales contribute equally to the clustering algorithm. The Elbow Method in figure 20 is then employed to determine the optimal number of clusters (k) for the K-Means algorithm. By iterating through different k values and plotting the Within-Cluster Sum of Squares (WCSS), an elbow point is identified, indicating an optimal balance between model complexity and clustering effectiveness. For this analysis, we initially chose k value of 3. However, after generating Elbow Method, an optimal k value of 10 was selected. Having us determined the optimal k, the K-Means algorithm is applied to cluster the inventory data. Each SKU is assigned to one of the identified clusters based on its 'Stock' level. The resulting clusters are appended as a new 'Cluster' column in the DataFrame, and the clustered data is saved to our new CSV file for further analysis or visualization.

Figure 20
Elbow method for optimal k



To gain a visual understanding of the clustering results, a box plot is generated in figure 21. This box plot displays the distribution of 'Stock' levels for each cluster, providing insights into the variability and central tendencies within each group. The box plot complements the numerical analysis and aids in the interpretation of the clustering outcomes.

Figure 21
Stock level distribution by cluster



Our analysis also involves creating separate DataFrames for each cluster and performing a detailed SKU analysis within each cluster. Basic statistics, such as mean, standard deviation, and quartiles, are displayed for the 'Stock' levels within each cluster, offering a granular understanding of the SKU distribution and characteristics within each identified group. In figure 22, one of the ten clusters' descriptive statistics are displayed, which is cluster 1.

Figure 22
Cluster 1's descriptive statistics

Descriptive statistics	SKU analysis
count	66
mean	268,48
std	40,51
min	214
25%	234
50%	258
75%	302
max	350

In summary, this K-Means Clustering analysis facilitates the segmentation of inventory items based on their stock levels, revealing inherent patterns and enabling Amazon to make informed decisions regarding stock management, procurement strategies, and overall inventory optimization.

Evaluation

Market basket model of Amazon Sales Report

In conducting a market basket analysis on the Amazon sales report, it is imperative to define the evaluation criteria that align with the objectives of uncovering patterns in customer purchasing behavior. Metrics such as support, confidence, and lift are critical in our context. Support measures the frequency of co-occurrence of items, confidence gauges the likelihood of purchasing a consequent item given the antecedent, and lift assesses the strength of the association between items.

To split our dataset, a portion of it is designated as the training set, while the remaining data serves as the testing set. This division allows for assessing our model's ability to generalize patterns identified in the training phase to previously unseen data in the testing set. Applying the market basket analysis model to the test our dataset involves evaluating its performance against the predefined criteria. Success is determined by the model's capability to reveal meaningful associations among purchased items, thereby contributing insights that meet the business's objectives, such as improving product recommendations, optimizing marketing strategies, or enhancing the overall customer experience.

Fine-tuning of our model may be necessary based on the evaluation results. Considerations of adjustments to hyperparameters or the inclusion of additional transactional data could refine the model's ability to capture subtle associations within the market basket.

In addition, the implementation of cross-validation techniques, such as k-fold cross-validation, is an essential assessment to the robustness of the market basket analysis model across different subsets of Amazon sales data. This approach enhances confidence in the model's generalization capabilities.

If multiple market basket analysis models were developed during the modeling phase, a detailed comparison of their performance is crucial. This involves considering trade-offs between metrics like support, confidence, and lift, ensuring the selected model strikes a balance between accuracy and interpretability.

The evaluation process extends beyond quantitative metrics to consider the business impact of our market basket analysis. Metrics like sales increase, customer engagement improvement, or marketing cost reduction are weighed against qualitative factors to provide a comprehensive understanding of the model's implications.

Interpretability and explainability of our market basket analysis model are necessary to be evaluated to ensure that stakeholders can comprehend and trust the generated insights. This is particularly important for making informed business decisions based on the identified purchasing patterns. Documentation of the evaluation results is crucial, capturing insights gained from the analysis, strengths, weaknesses, and recommendations. Clear communication of findings to stakeholders facilitates their understanding and buy-in, fostering a collaborative approach to leveraging market basket insights.

In the iterative phase, if our market basket analysis model falls short of meeting desired performance or business objectives, adjustments may be made to the problem definition, data collection, or alternative modeling approaches may be explored.

K-Means Clustering of International Sales Report

Metrics such as within-cluster sum of squares (WCSS), silhouette score, or Davies-Bouldin index could be employed. These metrics enable the quantification of how well the data points within a cluster are grouped, aiding in the determination of the optimal number of clusters for the model.

To ensure our model's ability to generalize to unseen data, the international sales report dataset is divided into training and testing sets. This partition allows for evaluating how well the clusters identified during the training phase are applicable to new data, providing insights into the model's robustness and generalization capabilities.

Success is measured by how effectively the model groups similar sales transactions, contributing to a clearer understanding of purchasing patterns and customer segmentation. Assessing how well our K-means clustering model achieves its objectives involves a nuanced examination of the clustered data in the international sales report dataset. The effectiveness of the clusters in revealing meaningful insights about customer behavior, market segments, or product preferences is crucial in determining the model's success in meeting business requirements.

Fine-tuning of the K-means clustering model may be deemed necessary based on the evaluation results. Adjustments to hyperparameters, such as the number of clusters, or retraining the model with additional sales data can enhance its ability to discern subtle patterns within the dataset.

When multiple models are developed during the clustering phase, a detailed comparison of their performance is undertaken. This involves considering trade-offs between metrics like WCSS and silhouette score, emphasizing the importance of not only accuracy but also the interpretability of the clusters.

The evaluation extends beyond quantitative metrics to assess the business impact of the K-means clustering analysis on international sales. Metrics such as revenue increase or cost reduction are considered alongside qualitative factors, providing a comprehensive understanding of the model's implications for strategic decision-making.

Documenting the evaluation results involves capturing insights gained from our analysis, strengths, weaknesses, and specific recommendations. Clear communication of these findings to stakeholders facilitates their understanding and buy-in, fostering a collaborative approach to leveraging international sales insights.

If our K-means clustering model falls short of meeting desired performance or business objectives, an iterative process is initiated. Refining the problem definition, collecting additional international sales data, or exploring alternative clustering approaches may be necessary.

The decision on deployment is based on a thorough assessment of the model's robustness, reliability, and alignment with business goals. Stakeholders weigh the benefits derived from the identified clusters against potential risks, ensuring that the deployed model enhances rather than hinders strategic decision-making in the international sales domain.

K-Means Clustering of Sales Report

It is imperative to define clear criteria and metrics for assessing our model's performance. Metrics such as within-cluster sum of squares (WCSS), silhouette score, or Davies-Bouldin index are particularly relevant in the context of clustering, as they help quantify the cohesion within clusters and separation between them, providing insights into the quality of the identified clusters.

To ensure the model's ability to generalize to unseen data, the sales report dataset is methodically split into training and testing sets. This division allows for a rigorous evaluation of how well the clusters identified during the training phase extend to new data, crucial for understanding the model's robustness and generalization capabilities.

The subsequent model evaluation entails the application of the trained K-means clustering model to the test dataset, with a focus on assessing its performance using the predefined metrics. Success is measured by how effectively the model groups similar stock levels, shedding light on potential stock categories or inventory management strategies.

Assessing how well our K-means clustering model achieves its objectives involves a nuanced examination of the clustered data in the sales report, specifically focusing on stock-related insights. The effectiveness of the clusters in revealing meaningful patterns in stock levels and aiding in inventory management decisions is critical in determining the model's success in meeting business requirements.

When multiple models are developed during the clustering phase, a comprehensive comparison of their performance is undertaken. This involves considering trade-offs between metrics like WCSS and silhouette score, emphasizing the importance of not only accuracy but also the interpretability of the clusters, especially in the context of stock management. The evaluation extends beyond quantitative metrics to assess the business impact of our K-means clustering analysis on sales and stock management. Metrics such as potential cost reductions or improvements in stock turnover are considered alongside qualitative factors, providing a holistic understanding of the model's implications for inventory-related decision-making.

Documenting the evaluation results involves capturing insights gained from our analysis, emphasizing the strengths, weaknesses, and specific recommendations. Clear communication of these findings to stakeholders facilitates their understanding and buy-in, fostering a collaborative approach to leveraging stock-related insights.

If the K-means clustering model falls short of meeting desired performance or business objectives, an iterative process is initiated. Refining the problem definition, collecting additional sales and stock data, or exploring alternative clustering approaches may be necessary.

The decision on deployment is based on a thorough assessment of our model's robustness, reliability, and alignment with business goals, specifically in the context of stock management. Stakeholders weigh the benefits derived from the identified clusters against potential risks, ensuring that the deployed model enhances rather than hinders strategic decision-making in stock-related domains.

Deployment

Market Basket Analysis of Amazon Sales Report

In the deployment phase of our Market Basket Analysis for Amazon's sales report dataset, the integration with business processes involves incorporating the association rules derived from our analysis into the online retail platform. For example, when the analysis reveals strong associations between certain products, Amazon can strategically place these items together on the website or suggest them as bundled purchases during the checkout process.

Automation of model execution is crucial for real-time recommendations. As users navigate the platform, our market basket analysis model can automatically provide product suggestions based on their current selections, contributing to a personalized shopping experience. This could lead to increased sales and improved customer satisfaction. Development of user interfaces or dashboards becomes essential for Amazon's merchandising and marketing teams. A user-friendly interface displaying key association rules and insights would empower these teams to make informed decisions about product placements, promotions, and cross-selling strategies. Communication and training are imperative in the deployment of our market basket analysis. Amazon should conduct training sessions for relevant teams, such as marketing and merchandising, to ensure they understand and effectively leverage the insights from the analysis in their day-to-day activities.

Monitoring and maintenance involve continuous tracking of how well the deployed association rules are performing on the Amazon platform; Regular reviews help identify changing customer behaviors and market trends, prompting necessary updates to the association rules for more accurate recommendations.

Performance tracking and evaluation should involve monitoring metrics like the click-through rate on recommended products, conversion rates, and overall sales attributed to the market basket recommendations. This information is vital for assessing the impact of the analysis on Amazon's business outcomes.

K-Means Clustering of International Sales Report

The clusters derived from our analysis, representing different customer segments, should be integrated into the customer relationship management (CRM) system. This allows Amazon to tailor marketing strategies, promotions, and customer service based on the characteristics of each cluster.

Automation of our model execution ensures that customer segmentation is an ongoing, automated process. The model can be scheduled to re-run at regular intervals, updating the customer clusters as new data becomes available. This automation contributes to maintaining up-to-date and relevant customer segments.

Communication and training are essential in the deployment phase, particularly for marketing and customer service teams. Training sessions would be conducted to familiarize teams with the international customer segments and guide them on how to tailor their approaches for more effective engagement.

Monitoring and maintenance involve continuous tracking of the international customer segments' behaviors. Amazon should regularly assess whether the characteristics of the clusters remain consistent or if adjustments are needed due to changing market dynamics.

Documentation of deployment processes is vital for ensuring that our K-Means Clustering model's integration into Amazon's CRM system is well-documented. This documentation should include guidelines for handling updates, resolving potential issues, and making any necessary adjustments.

Performance tracking and evaluation should revolve around assessing how well targeted marketing efforts align with the characteristics of each international customer segment.

Metrics such as customer engagement, conversion rates, and revenue generated from each cluster would guide the evaluation process.

K-Means Clustering of Sales Report with Stock Focus

In the case of our K-Means Clustering analysis on the sales report dataset with a focus on stock levels, integration with business processes is centered around inventory management. The clusters representing different stock behavior patterns should be integrated into Amazon's inventory control system to optimize stock levels and distribution.

Automation of model execution ensures that stock behavior clusters are regularly updated based on the latest sales and inventory data. This automation contributes to efficient stock management, helping Amazon avoid stockouts or excess inventory.

Moreover, development of user interfaces or dashboards is beneficial for Amazon's logistics and supply chain teams. An interface displaying the characteristics of each stock behavior cluster, such as demand patterns and lead times, would enable teams to make informed decisions regarding stock replenishment and distribution.

Communication and training are crucial for logistics and warehouse teams involved in stock management. Training sessions would educate teams on how to interpret the characteristics of each stock behavior cluster and implement corresponding actions, such as adjusting reorder points.

Monitoring and maintenance involve continuous tracking of how well our deployed stock behavior clusters align with actual stock movements. Regular reviews help identify shifts in demand patterns or unexpected stock behavior, prompting necessary updates to the clusters. Performance tracking and evaluation should focus on metrics such as stock turnover rates, order fulfillment efficiency, and overall inventory costs. Evaluating how well Amazon's stock management aligns with the insights from the clustering analysis helps optimize inventory-related processes.

Final deployment

The integration with business processes is a common thread across all three analyses. The association rules derived from the Market Basket Analysis seamlessly integrate into Amazon's recommendation engine, influencing product placements, bundling strategies, and personalized shopping experiences. Simultaneously, the customer segments from the international sales clustering and the stock behavior clusters from the sales report clustering embed themselves into Amazon's CRM and inventory management systems, respectively. This integration allows Amazon to align its operations with the distinct characteristics and behaviors uncovered in each analysis.

Automation plays a pivotal role in ensuring the relevance and real-time applicability of the our insights. For our Market Basket Analysis, the automated provision of product suggestions during user navigation creates an immediate impact on the shopping experience. Similarly, the periodic automation of our K-Means Clustering analyses ensures that customer segments and stock behavior clusters remain up-to-date, reflecting the latest trends and patterns. This synchronization across automated processes ensures that Amazon's platform continually adapts to evolving customer behaviors and market dynamics.

Marketing, merchandising, logistics, and supply chain teams benefit from intuitive interfaces that display actionable insights specific to their domains. For instance, dashboards for international sales clustering highlight characteristics of customer segments, guiding targeted marketing campaigns. Simultaneously, interfaces for stock behavior clusters empower logistics teams to make informed decisions regarding inventory replenishment and distribution.

Teams involved in marketing, merchandising, customer service, logistics, and warehouse management should receive tailored guidance on how to leverage the derived insights to enhance their day-to-day operations. This ensures a unified understanding and implementation of strategies aligned with the unique characteristics of each analysis. Regular reviews and updates are crucial for all three deployments. Whether tracking changing customer behaviors for Market Basket Analysis, assessing shifts in international customer segments, or adapting to evolving stock patterns, ongoing maintenance ensures that insights remain relevant and actionable.

Performance tracking and evaluation are the ultimate litmus tests for the success of these deployments. Metrics such as sales conversion rates, customer engagement, and optimized inventory costs gauge the tangible impact on Amazon's business objectives. Regular evaluations inform necessary adjustments to strategies, ensuring that the derived insights continue to drive positive outcomes.

In conclusion, the deployment of our data-driven insights represents a holistic approach to enhancing Amazon's e-commerce ecosystem. The symbiotic relationship between our Market Basket Analysis, international sales clustering, and sales clustering with a stock focus forms a cohesive strategy that aligns Amazon's operations with customer behaviors, market trends, and inventory dynamics, ultimately driving business success.

Bibliography

<https://github.com/sophie300902>

Python Fundamentals. (2023). Effective Data Cleaning Strategies in Python for Exploratory Data Analysis (EDA). Retrieved from <https://python.plainenglish.io/effective-data-cleaning-strategies-in-python-for-exploratory-data-analysis-eda-957122b8864c>

Sharma, A. (2022). Ecommerce Sales dataset. Retrieved from <https://data.world/anilsharma87/sales>