

1 Start coding or [generate](#) with AI.

## MACHINE LEARNING E-COMMERCE ANALYTICS PROJECT

---

1 Start coding or [generate](#) with AI.

### \*MECE BREAKDOWN \*

1. Customer Data This category contains all user-related data that describes the attributes and behavior of customers.

Demographic Data: Age Gender Location (city, state, country) Income Level Behavioral Data: Session Duration Number of Visits Browsing History (viewed pages, search queries) Device Type (mobile, desktop, etc.)

2. Product Data All information related to the products sold on the platform. This is specific to the items available for purchase.

Product Attributes: Product ID Name Category Price Brand Stock Availability Ratings (average rating, number of ratings) Reviews (text of customer feedback) Product Metadata: Launch Date Product Description Manufacturer Details

3. Transaction Data Data that captures the actual transactions or purchases made by the customers.

Order Information: Order ID Transaction Date Product ID (linked to the product table) Quantity Ordered Order Value (total amount)

Discounts/Promotions Applied Payment Details: Payment Method (credit card, PayPal, etc.) Transaction Status (completed, pending, failed)

Shipping Address Billing Information Shipping Costs Taxes

4. Marketing and Engagement Data Data related to how customers interact with the website and marketing campaigns.

Campaign Information: Campaign ID Type of Campaign (email, social media, ads, etc.) Click-Through Rate (CTR) Conversion Rate Engagement Metrics (likes, shares, comments) User Engagement Metrics: Newsletter Subscriptions Coupon Usage Loyalty Program Participation Ad Impressions

5. Operational Data This category includes data related to the operations of the e-commerce platform, including logistics, inventory, and support.

Inventory Data: Product Stock Levels Restocking Date Warehouse Locations Shipping Information: Shipping Method Delivery Status (in-transit, delivered, returned) Shipping Provider Estimated Delivery Time Customer Support Data: Support Ticket ID Type of Issue (payment, product, delivery) Resolution Status Time to Resolution

6. Website Performance and Technical Data This category covers the technical data that tracks website performance, user interaction, and system health.

Traffic Data: Number of Visitors Page Views Bounce Rate Load Time Session Information: Session ID Session Duration Time Spent on Each Page Exit Page (last page visited before exiting)

1 Start coding or [generate](#) with AI.

### First all import the required Libraries

1 Start coding or [generate](#) with AI.

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
8 from sklearn.linear_model import LogisticRegression, LinearRegression
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.metrics import accuracy_score
11 from sklearn.cluster import KMeans
12
13
14 df=pd.read_csv("/content/kz.csv")
15 df.head()
16 df.info()
17 df.describe()
```

```
18 df.isnull().sum()
19 df.duplicated()

→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 2633521 entries, 0 to 2633520
Data columns (total 8 columns):
 #   Column      Dtype  
0   event_time   object  
1   order_id     int64  
2   product_id   int64  
3   category_id  float64 
4   category_code object  
5   brand        object  
6   price        float64 
7   user_id      float64 
dtypes: float64(3), int64(2), object(3)
memory usage: 160.7+ MB
```

	0
0	False
1	True
2	False
3	True
4	False
...	...
2633516	False
2633517	False
2633518	False
2633519	False
2633520	False

2633521 rows × 1 columns

1 Start coding or generate with AI.

1 Start coding or generate with AI.

#### Problem statements based on the MECE breakdown of the dataset

1. **Customer Segmentation:** How can we segment customers into distinct groups based on their demographics and purchasing behavior to improve targeted marketing?
- 2 . **Churn Prediction:** What are the key factors that lead to customer churn, and how can we predict which users are at risk of leaving the platform?
2. **Product Recommendations:** How can we build a recommendation engine that suggests relevant products to customers based on their browsing history and past purchases?
3. **Customer Lifetime Value (CLV) Prediction:** Can we accurately predict the long-term value of a customer based on their transaction history and engagement with the platform?
4. **Price Optimization:** How can product pricing be optimized to maximize revenue without negatively impacting sales volume?
5. **Sales Forecasting:** Can we develop a model to predict future sales trends based on historical sales data, seasonality, and marketing campaigns?
6. **Inventory Management:** How can we predict inventory needs to minimize stockouts and overstock situations, taking into account sales trends and product demand?
7. **Campaign Effectiveness:** Which marketing campaigns have the highest return on investment (ROI), and how can we improve engagement rates through more effective targeting?
8. **Fraud Detection:** How can we detect potentially fraudulent transactions by analyzing user behavior and payment data?

9. **Sentiment Analysis of Reviews:** How can we use natural language processing (NLP) techniques to classify the sentiment of product reviews and identify areas for product improvement? 11. **Conversion Rate Optimization:** What factors contribute most to conversion (product purchase) after customers interact with marketing campaigns or browse the platform?

10. \*\*Website Performance:\*\* How does website performance (e.g., page load time, bounce rate) impact customer satisfaction and sales, and how can we improve it?

11. **User Engagement Prediction:** Can we predict which users are most likely to engage with the platform regularly, and what factors drive this engagement?

12. **Delivery Time Optimization:** How can we predict delivery times more accurately based on shipping method, warehouse location, and customer destination to improve customer satisfaction?

13. **Customer Support Analysis:** What are the most common customer support issues, and how can we use this data to improve product offerings and the customer experience?

Generate
create a dataframe with 2 columns and 10 rows

1 Start coding or generate with AI.

Double-click (or enter) to edit

```

1 print(df.columns)
2
3 → Index(['event_time', 'order_id', 'product_id', 'category_id', 'category_code',
       'brand', 'price', 'user_id'],
       dtype='object')
4
5 df.head()
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

```

	event_time	order_id	product_id	category_id	category_code	brand	price	user_id
0	2020-04-24 11:50:39 UTC	2294359932054536986	1515966223509089906	2.268105e+18	electronics.tablet	samsung	162.01	1.515916e+18
1	2020-04-24 11:50:39 UTC	2294359932054536986	1515966223509089906	2.268105e+18	electronics.tablet	samsung	162.01	1.515916e+18
2	2020-04-24 14:37:43 UTC	2294444024058086220	2273948319057183658	2.268105e+18	electronics.audio.headphone	huawei	77.52	1.515916e+18

### 1 \*\*\*Customer Segmentation (Clustering)

```

1 from sklearn.cluster import KMeans
2 import pandas as pd
3 import numpy as np # Import numpy for NaN handling
4
5 # Load dataset
6 df = pd.read_csv('/content/kz.csv')
7
8 # Select relevant features for customer segmentation
9 X = df[['user_id', 'order_id', 'price', 'category_id']]
10
11 # Handle NaN values: Replace with 0 or use other imputation methods
12 X = X.fillna(0) # Replace NaNs with 0.
13 # Consider other strategies like mean imputation if appropriate.
14
15 # Apply K-means clustering
16 kmeans = KMeans(n_clusters=5, random_state=42)
17 df['customer_segment'] = kmeans.fit_predict(X)
18
19 # View customer segments
20 df[['user_id', 'customer_segment']].head()

```

	user_id	customer_segment
0	1.515916e+18	2
1	1.515916e+18	2
2	1.515916e+18	2
3	1.515916e+18	2
4	1.515916e+18	2

## 2 Churn Prediction

```

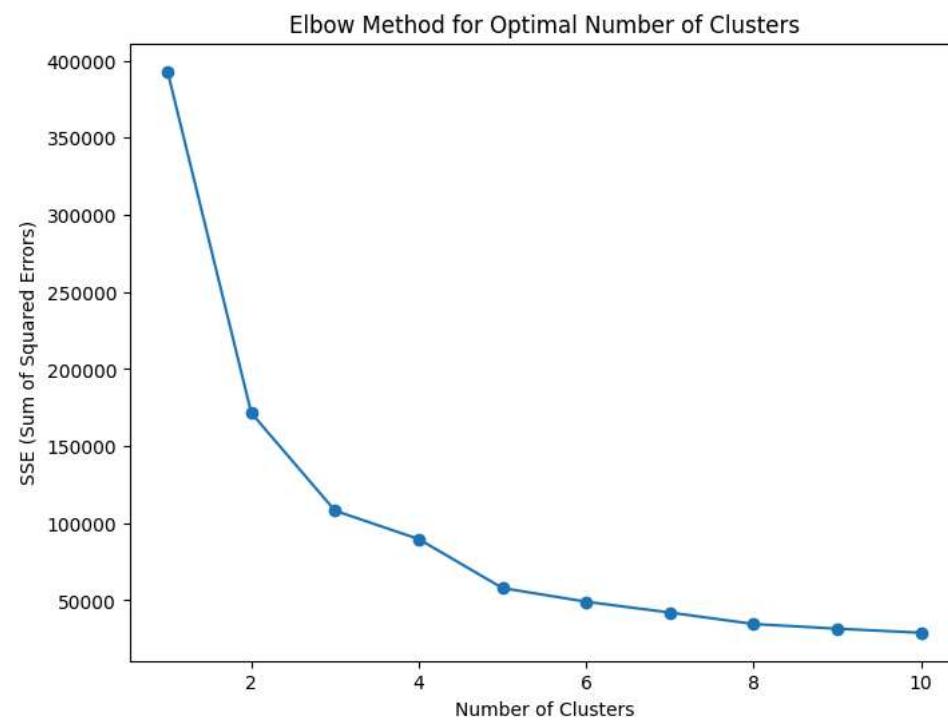
1 # Importing necessary libraries
2 import pandas as pd
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.cluster import KMeans
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 # Step 1: Load dataset
9 df = pd.read_csv('/content/kz.csv') # Replace with actual file path
10
11 # Step 2: Feature Engineering (Creating customer-level features)
12 # Group by 'user_id' to create customer behavior metrics
13 user_df = df.groupby('user_id').agg({
14     'order_id': 'nunique',          # Number of unique orders (order count)
15     'price': ['sum', 'mean'],       # Total spent and average price per order
16     'product_id': 'nunique'        # Number of unique products purchased
17 }).reset_index()
18
19 # Rename columns for easier access
20 user_df.columns = ['user_id', 'num_orders', 'total_spent', 'avg_spent', 'num_unique_products']
21
22 # Step 3: Data Preprocessing
23 # Standardize the data so that each feature has a similar scale
24 scaler = StandardScaler()
25 X = user_df[['num_orders', 'total_spent', 'avg_spent', 'num_unique_products']]
26 X_scaled = scaler.fit_transform(X)
27
28 # Step 4: Apply K-Means Clustering
29 # Using the Elbow method to determine the optimal number of clusters
30 sse = []
31 for k in range(1, 11):
32     kmeans = KMeans(n_clusters=k, random_state=42)
33     kmeans.fit(X_scaled)
34     sse.append(kmeans.inertia_)
35
36 # Plot the Elbow Curve
37 plt.figure(figsize=(8, 6))
38 plt.plot(range(1, 11), sse, marker='o')
39 plt.title('Elbow Method for Optimal Number of Clusters')
40 plt.xlabel('Number of Clusters')
41 plt.ylabel('SSE (Sum of Squared Errors)')
42 plt.show()
43
44 # Based on the elbow plot, we choose the optimal number of clusters (e.g., 4)
45 kmeans = KMeans(n_clusters=4, random_state=42)
46 user_df['cluster'] = kmeans.fit_predict(X_scaled)
47
48 # Step 5: Analyze the clusters
49 # Let's look at the characteristics of each cluster
50 cluster_analysis = user_df.groupby('cluster').agg({
51     'num_orders': ['mean', 'median'],
52     'total_spent': ['mean', 'median'],
53     'avg_spent': ['mean', 'median'],
54     'num_unique_products': ['mean', 'median'],
55     'user_id': 'count' # Number of customers in each cluster
56 }).reset_index()
57
58 print(cluster_analysis)
59
60 # Step 6: Visualize the Clusters
61 # Scatter plot of Total Spent vs Number of Orders, colored by cluster
62 plt.figure(figsize=(10, 6))

```

```

63 sns.scatterplot(data=user_df, x='num_orders', y='total_spent', hue='cluster', palette='Set1')
64 plt.title('Customer Segmentation based on Total Spend and Number of Orders')
65 plt.xlabel('Number of Orders')
66 plt.ylabel('Total Spent')
67 plt.show()
68

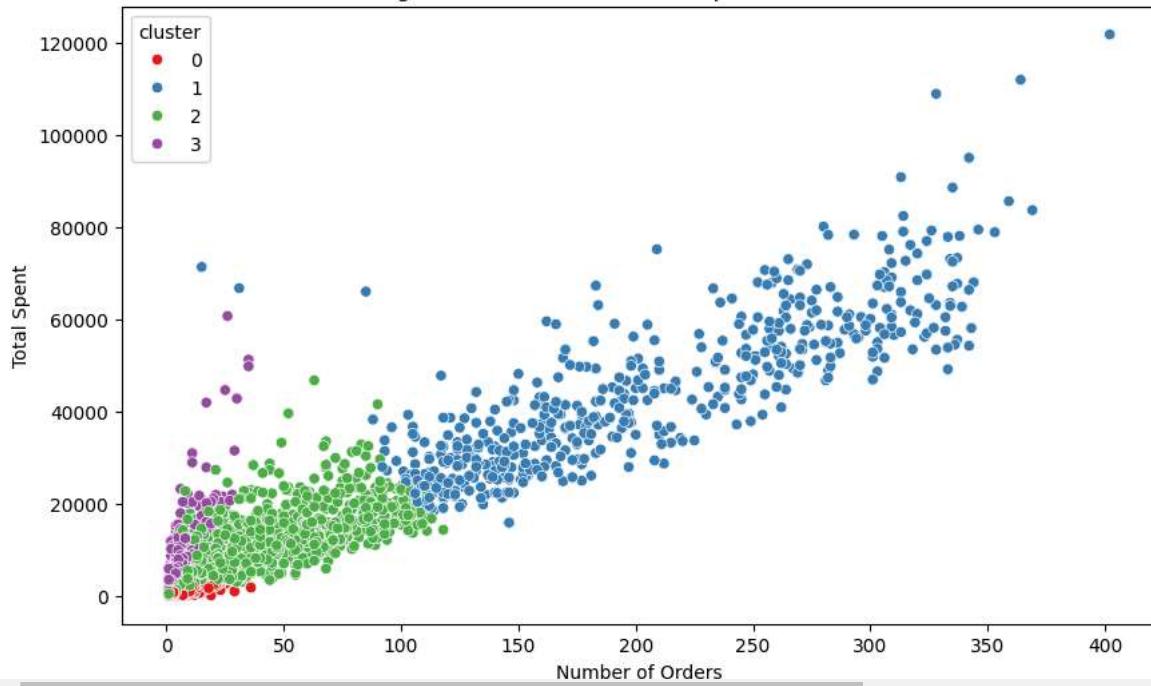
```



	cluster	num_orders		total_spent		avg_spent		\
		mean	median	mean	median	mean	median	
0	0	2.274618	2.0	355.802673	199.03	113.192697	105.413889	
1	1	204.930902	193.0	44284.516161	42047.52	142.248474	134.123143	
2	2	4.855551	2.0	1944.981297	1147.85	392.370223	370.350000	
3	3	2.088894	1.0	2148.241442	1388.84	996.301256	879.610000	

	num_unique_products		user_id	
	mean	median	count	
0	2.736793	2.0	61256	
1	279.180422	262.0	521	
2	5.380540	2.0	28273	
3	1.628714	1.0	8212	

Customer Segmentation based on Total Spend and Number of Orders



### 3 Product Recommendations

```

1 # Import necessary libraries
2 import pandas as pd
3 from sklearn.cluster import KMeans
4 from sklearn.preprocessing import StandardScaler
5
6 # Step 1: Load the Data
7 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
8
9 # Step 2: Preprocess the Data
10 # Group by user_id to get total spent and purchase frequency
11 customer_data = df.groupby('user_id').agg({
12     'price': 'sum', # Total spent
13     'order_id': 'count' # Purchase frequency
14 }).reset_index()
15
16 customer_data.columns = ['user_id', 'total_spent', 'purchase_frequency']
17
18 # Step 3: Standardize the Data
19 scaler = StandardScaler()
20 X = scaler.fit_transform(customer_data[['total_spent', 'purchase_frequency']])
21
22 # Step 4: Apply K-Means Clustering
23 kmeans = KMeans(n_clusters=3, random_state=42) # Choose number of clusters
24 customer_data['cluster'] = kmeans.fit_predict(X)
25
26 # Step 5: View Results
27 print(customer_data.head())
28
29 # Optional: Display cluster centroids
30 centroids = scaler.inverse_transform(kmeans.cluster_centers_)
31 print("Cluster centroids (total_spent, purchase_frequency):")
32 print(centroids)
33

```

	user_id	total_spent	purchase_frequency	cluster	
0	1.515916e+18	416.64		1	0
1	1.515916e+18	56.43		2	0
2	1.515916e+18	7530.34		14	0
3	1.515916e+18	5074.47		24	0
4	1.515916e+18	182.83		2	0

Cluster centroids (total\_spent, purchase\_frequency):

```

[[8.93674164e+02 3.66180716e+00]
 [5.63225826e+04 4.16614035e+02]
 [2.25255228e+04 1.38647975e+02]]
```

```

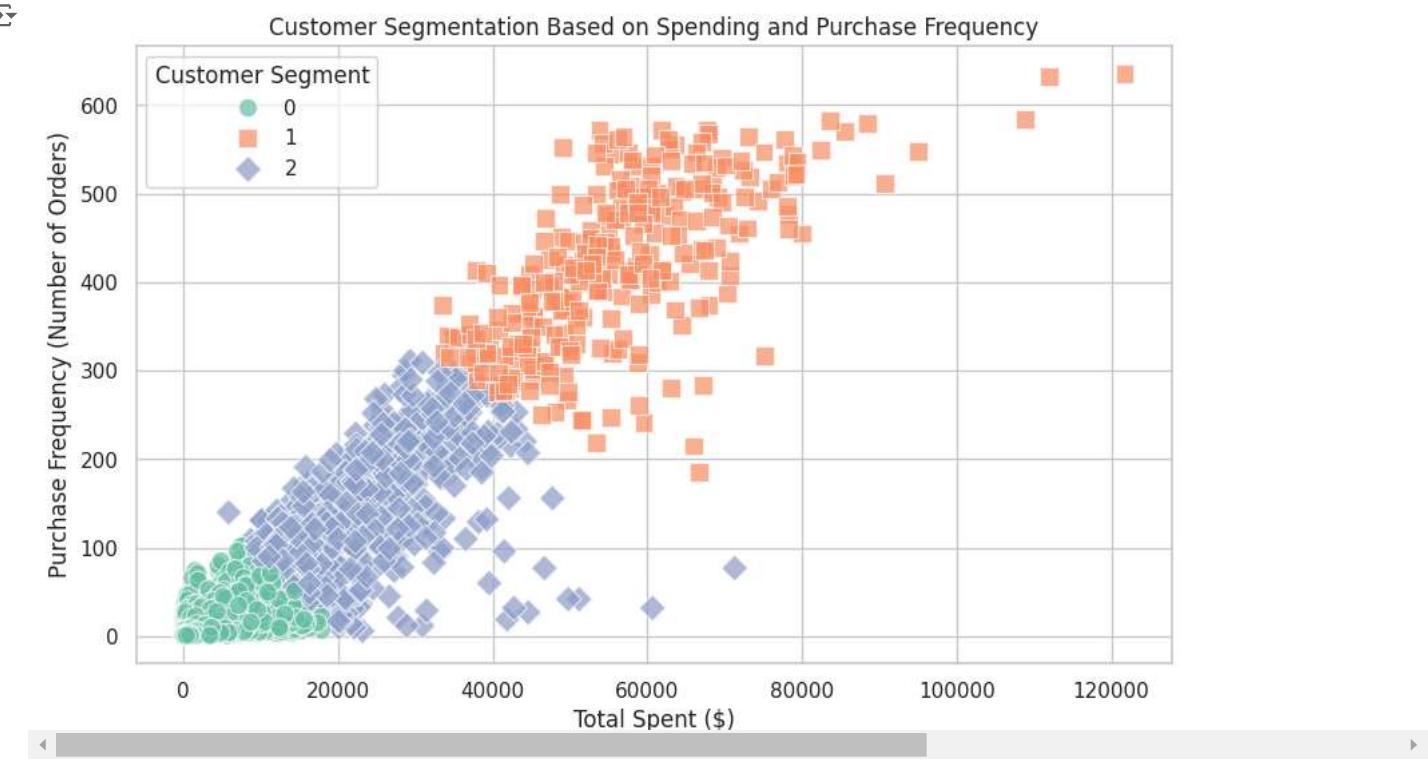
1
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5
6 sns.set(style="whitegrid")
7
8 # Create a scatter plot to visualize the clusters
9 plt.figure(figsize=(10, 6))
10 sns.scatterplot(
11     data=customer_data,
12     x='total_spent',
13     y='purchase_frequency',
14     hue='cluster', # Color by cluster
15     palette='Set2', # Set color palette
16     style='cluster', # Different marker styles for clusters
17     markers=["o", "s", "D"], # Specify different markers for each cluster
18     s=100, # Marker size
19     alpha=0.7 # Transparency
20 )
21
22 # 3: Customize the plot
23 plt.title('Customer Segmentation Based on Spending and Purchase Frequency')
24 plt.xlabel('Total Spent ($)')
25 plt.ylabel('Purchase Frequency (Number of Orders)')
26 plt.legend(title='Customer Segment')
27 plt.grid(True)
28

```

```

29 # Step 4: Show the plot
30 plt.show()
31

```



#### \*\* 4 Customer Lifetime Value (CLV) Prediction\*\*

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestRegressor
6 from sklearn.metrics import mean_absolute_error
7
8 # Step 1: Load the Data
9 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
10
11 # Step 2: Preprocess the Data
12 # Convert event_time to datetime
13 df['event_time'] = pd.to_datetime(df['event_time'])
14
15 # Calculate recency, frequency, and monetary (RFM) metrics for CLV prediction
16 current_date = df['event_time'].max()
17 rfm_df = df.groupby('user_id').agg({
18     'event_time': lambda x: (current_date - x.max()).days, # Recency
19     'order_id': 'count', # Frequency
20     'price': 'sum' # Monetary Value
21 }).reset_index()
22
23 # Rename columns for clarity
24 rfm_df.columns = ['user_id', 'recency', 'frequency', 'monetary']
25
26 # Step 3: Feature Engineering
27 # Create the target variable for CLV
28 rfm_df['CLV'] = rfm_df['monetary'] * rfm_df['frequency'] # Basic CLV calculation
29
30 # Step 4: Split the Data
31 X = rfm_df[['recency', 'frequency', 'monetary']] # Features
32 y = rfm_df['CLV'] # Target variable
33
34 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
35
36 # Step 5: Train a Model
37 model = RandomForestRegressor(n_estimators=100, random_state=42) # You can tune this model
38 model.fit(X_train, y_train)
39

```

```

40 # Step 6: Make Predictions
41 y_pred = model.predict(X_test)
42
43 # Step 7: Evaluate the Model
44 mae = mean_absolute_error(y_test, y_pred)
45 print(f'Mean Absolute Error: {mae:.2f}')
46
47 # Optional: Display the first few predictions
48 predictions_df = pd.DataFrame({'Actual CLV': y_test, 'Predicted CLV': y_pred})
49 print(predictions_df.head())
50

→ Mean Absolute Error: 2088.53
   Actual CLV  Predicted CLV
26774      274.11      274.1076
4148       194.25      194.2302
91537       138.63      138.6300
94695       1388.80     1388.8000
91645      5917.25     5913.5105

```

### Price Optimization:

```

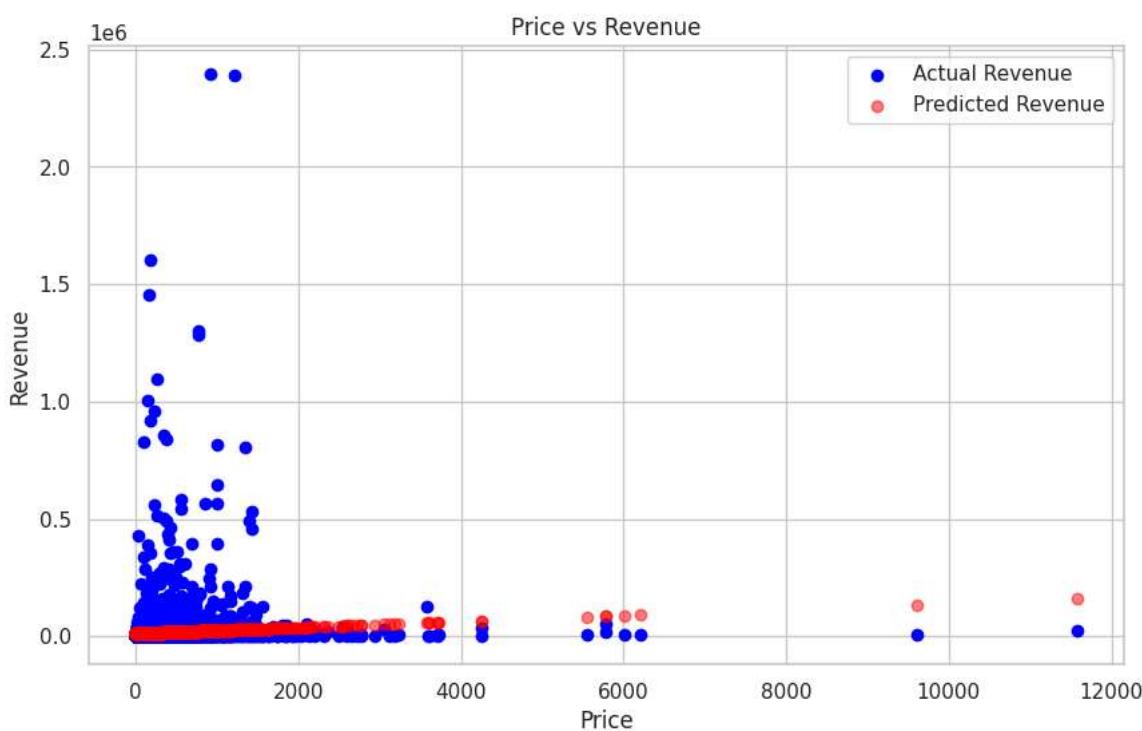
1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression
6 import matplotlib.pyplot as plt
7
8 # Step 1: Load the Data
9 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
10
11 # Step 2: Preprocess the Data
12 # Convert event_time to datetime
13 df['event_time'] = pd.to_datetime(df['event_time'])
14
15 # Check if the dataset has required columns
16 required_columns = ['event_time', 'order_id', 'product_id', 'category_id', 'category_code', 'brand', 'price', 'user_id']
17 for col in required_columns:
18     if col not in df.columns:
19         raise ValueError(f'Missing required column: {col}')
20
21 # Calculate total revenue and total sales volume by product_id and price
22 revenue_df = df.groupby(['product_id', 'price']).agg({
23     'order_id': 'count' # Sales volume (number of orders)
24 }).reset_index()
25
26 # Calculate revenue for each price point
27 revenue_df['revenue'] = revenue_df['order_id'] * revenue_df['price']
28
29 # Step 3: Train a Linear Regression Model
30 X = revenue_df[['price']] # Features
31 y = revenue_df['revenue'] # Target variable
32
33 # Split the data
34 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
35
36 # Train the model
37 model = LinearRegression()
38 model.fit(X_train, y_train)
39
40 # Step 4: Make Predictions
41 predicted_revenue = model.predict(X_test)
42
43 # Step 5: Visualize Results
44 plt.figure(figsize=(10, 6))
45 plt.scatter(X_test, y_test, color='blue', label='Actual Revenue')
46 plt.scatter(X_test, predicted_revenue, color='red', label='Predicted Revenue', alpha=0.5)
47 plt.xlabel('Price')
48 plt.ylabel('Revenue')
49 plt.title('Price vs Revenue')
50 plt.legend()
51 plt.grid(True)
52 plt.show()
53
54 # Step 6: Determine Optimal Price

```

```

55 # Generate a price range for optimization
56 price_range = np.linspace(revenue_df['price'].min(), revenue_df['price'].max(), 100).reshape(-1, 1)
57 optimal_revenue = model.predict(price_range)
58
59 # Find optimal price point
60 optimal_price_index = np.argmax(optimal_revenue)
61 optimal_price = price_range[optimal_price_index][0]
62 max_revenue = optimal_revenue[optimal_price_index]
63
64 print(f'Optimal Price: ${optimal_price:.2f} with a Maximum Revenue of ${max_revenue:.2f}')
65

```



Optimal Price: \$50925.90 with a Maximum Revenue of \$666310.43

/usr/local/lib/python3.10/dist-packages/sklearn/base.py:493: UserWarning: X does not have valid feature names, but LinearRegression was warnings.warn(

Double-click (or enter) to edit

1 Start coding or generate with AI.

### Inventory Management:

```

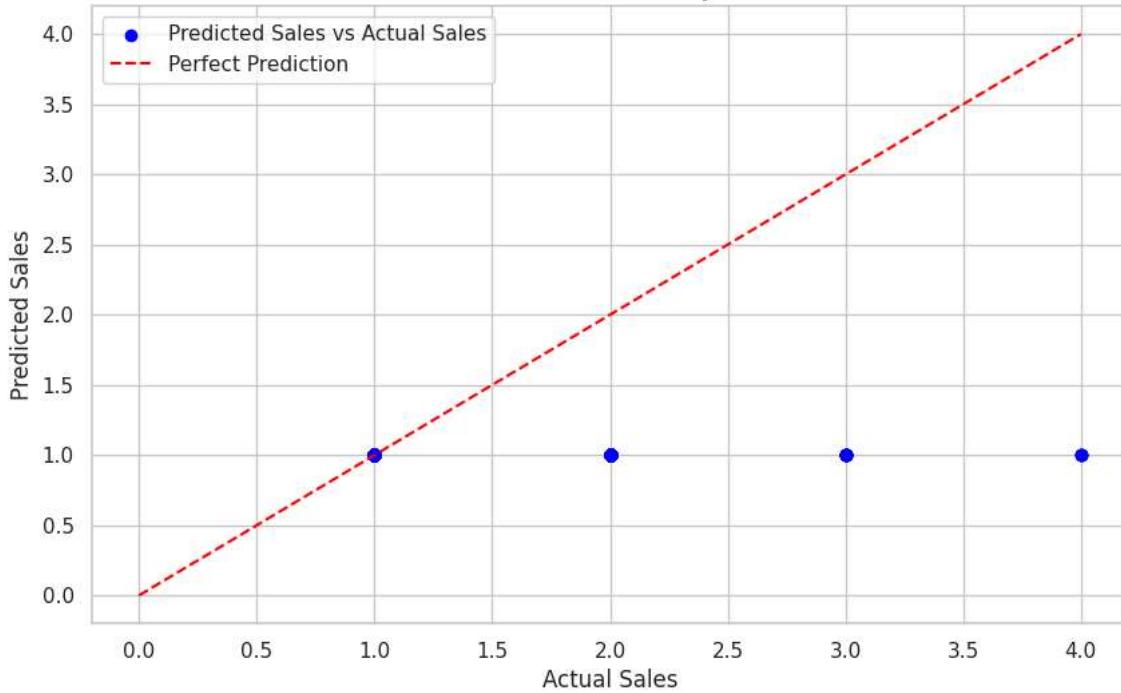
1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression
6 import matplotlib.pyplot as plt
7
8 # Step 1: Load the Data
9 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
10
11 # Step 2: Preprocess the Data
12 # Convert event_time to datetime
13 df['event_time'] = pd.to_datetime(df['event_time'])
14
15 # Check if the dataset has required columns
16 required_columns = ['event_time', 'order_id', 'product_id', 'category_id', 'category_code', 'brand', 'price', 'user_id']
17 for col in required_columns:
18     if col not in df.columns:
19         raise ValueError(f'Missing required column: {col}')
20
21 # Step 3: Create Sales Data
22 # Count orders per product per day

```

```
23 sales_data = df.groupby(['event_time', 'product_id']).agg({
24     'order_id': 'count' # Count of orders per day per product
25 }).reset_index()
26
27 # Rename columns
28 sales_data.columns = ['event_time', 'product_id', 'daily_sales']
29
30 # Step 4: Feature Engineering
31 # Create time-based features
32 sales_data['day_of_week'] = sales_data['event_time'].dt.dayofweek
33 sales_data['month'] = sales_data['event_time'].dt.month
34 sales_data['year'] = sales_data['event_time'].dt.year
35
36 # Create lag features for sales
37 sales_data['lag_sales'] = sales_data.groupby('product_id')['daily_sales'].shift(1)
38
39 # Drop rows with NaN values created by lag
40 sales_data = sales_data.dropna()
41
42 # Step 5: Prepare Data for Modeling
43 X = sales_data[['lag_sales', 'day_of_week', 'month', 'year']] # Features
44 y = sales_data['daily_sales'] # Target variable
45
46 # Split the data into training and testing sets
47 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
48
49 # Step 6: Train a Linear Regression Model
50 model = LinearRegression()
51 model.fit(X_train, y_train)
52
53 # Step 7: Make Predictions
54 predicted_sales = model.predict(X_test)
55
56 # Step 8: Evaluate the Model
57 plt.figure(figsize=(10, 6))
58 plt.scatter(y_test, predicted_sales, color='blue', label='Predicted Sales vs Actual Sales')
59 plt.plot([0, max(y_test)], [0, max(y_test)], color='red', linestyle='--', label='Perfect Prediction')
60 plt.xlabel('Actual Sales')
61 plt.ylabel('Predicted Sales')
62 plt.title('Actual vs Predicted Daily Sales')
63 plt.legend()
64 plt.grid(True)
65 plt.show()
66
67 # Calculate and print performance metrics
68 mae = np.mean(np.abs(y_test - predicted_sales)) # Mean Absolute Error
69 print(f'Mean Absolute Error: {mae:.2f}')
70
71 # Optional: Display predicted vs actual sales for the first few entries
72 predictions_df = pd.DataFrame({'Actual Sales': y_test, 'Predicted Sales': predicted_sales})
73 print(predictions_df.head())
74
```



Actual vs Predicted Daily Sales



```
Mean Absolute Error: 0.00
      Actual Sales   Predicted Sales
1713235           1       1.000890
342283            1       1.000946
811021            1       1.000917
1479238           1       1.000979
...
```

### Fraud Detection:

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import classification_report, confusion_matrix
7 import matplotlib.pyplot as plt
8
9 # Step 1: Load the Data
10 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
11
12 # Step 2: Preprocess the Data
13 # Convert event_time to datetime
14 df['event_time'] = pd.to_datetime(df['event_time'])
15
16 # Check if the dataset has required columns
17 required_columns = ['event_time', 'order_id', 'product_id', 'category_id', 'category_code', 'brand', 'price', 'user_id']
18 for col in required_columns:
19     if col not in df.columns:
20         raise ValueError(f'Missing required column: {col}')
21
22 # Step 3: Create Features for Fraud Detection
23 # Create a feature: Total purchases per user
24 user_purchase_count = df.groupby('user_id')['order_id'].count().reset_index()
25 user_purchase_count.columns = ['user_id', 'total_purchases']
26
27 # Create a feature: Average price per user
28 user_avg_price = df.groupby('user_id')['price'].mean().reset_index()
29 user_avg_price.columns = ['user_id', 'avg_price']
30
31 # Merge features back to the main dataframe
32 df = df.merge(user_purchase_count, on='user_id', how='left')
33 df = df.merge(user_avg_price, on='user_id', how='left')
34
35 # Create a synthetic label for fraud detection (this is just for demonstration)
36 # In practice, you would have a labeled dataset indicating fraud

```

```
37 np.random.seed(42) # For reproducibility
38 df['is_fraud'] = np.random.choice([0, 1], size=len(df), p=[0.95, 0.05]) # 5% fraud for simulation
39
40 # Step 4: Prepare Data for Modeling
41 X = df[['total_purchases', 'avg_price']] # Features
42 y = df['is_fraud'] # Target variable
43
44 # Split the data into training and testing sets
45 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
46
47 # Step 5: Train a Random Forest Classifier
48 model = RandomForestClassifier(random_state=42)
49 model.fit(X_train, y_train)
50
51 # Step 6: Make Predictions
52 y_pred = model.predict(X_test)
53
54 # Step 7: Evaluate the Model
55 print("Confusion Matrix:")
56 print(confusion_matrix(y_test, y_pred))
57 print("\nClassification Report:")
58 print(classification_report(y_test, y_pred))
59
60 # Optional: Feature Importance
61 importances = model.feature_importances_
62 feature_names = X.columns
63
64 plt.figure(figsize=(10, 6))
65 plt.barh(feature_names, importances)
66 plt.xlabel('Feature Importance')
67 plt.title('Feature Importance for Fraud Detection')
68 plt.show()
69
70
```

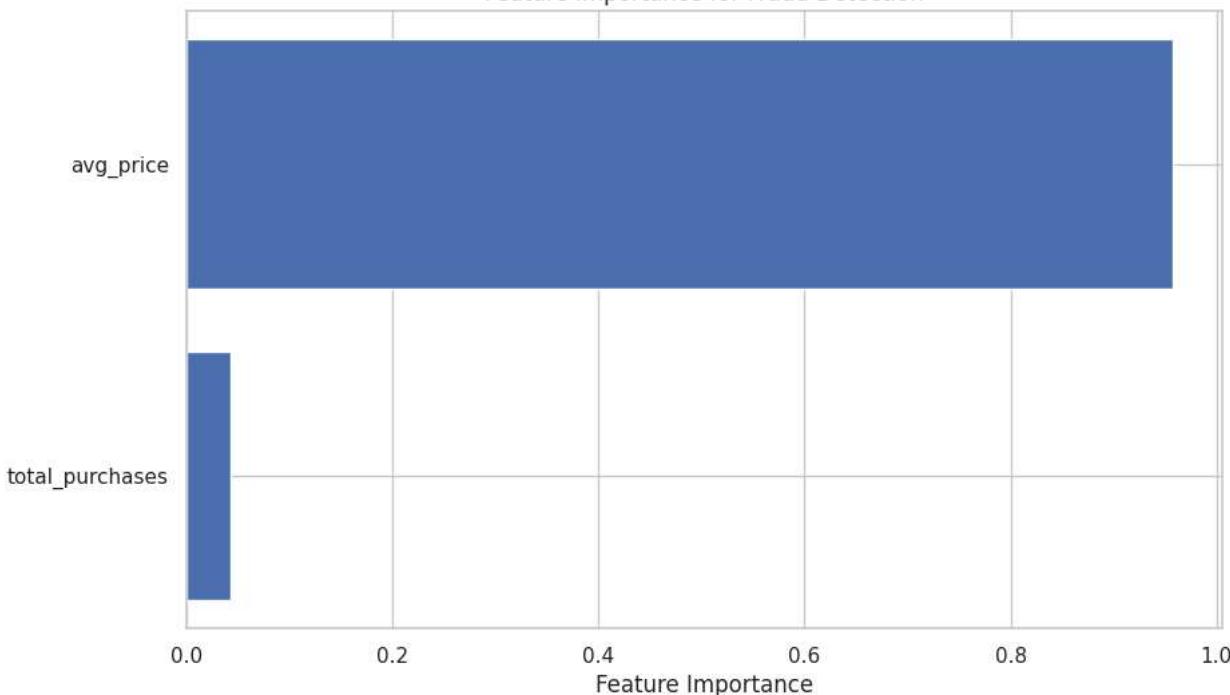
Confusion Matrix:

```
[[500049  199]
 [ 26447   10]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.95	1.00	0.97	500248
1	0.05	0.00	0.00	26457
accuracy			0.95	526705
macro avg	0.50	0.50	0.49	526705
weighted avg	0.90	0.95	0.93	526705

Feature Importance for Fraud Detection



Sentiment Analysis of Reviews:

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.feature_extraction.text import TfidfVectorizer
6 from sklearn.naive_bayes import MultinomialNB
7 from sklearn.metrics import classification_report, confusion_matrix
8
9 # Step 1: Load the Data
10 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
11
12 # Step 2: Create Mock Reviews
13 # For demonstration, let's create a mock reviews DataFrame based on product IDs
14 # In practice, you would replace this with actual product reviews
15 np.random.seed(42)
16 product_ids = df['product_id'].unique()
17 reviews = [f"This is a great product {pid}" if np.random.rand() > 0.5 else f"This product {pid} is bad" for pid in product_ids]
18 sentiments = [1 if "great" in review else 0 for review in reviews] # 1 for positive, 0 for negative
19
20 # Create a DataFrame with mock reviews and sentiments
21 reviews_df = pd.DataFrame({'product_id': product_ids, 'review': reviews, 'sentiment': sentiments})
22
23 # Step 3: Text Preprocessing and Feature Extraction
24 # Use TF-IDF to convert text data into numerical format
25 tfidf = TfidfVectorizer(stop_words='english')
26 X = tfidf.fit_transform(reviews_df['review'])
27 y = reviews_df['sentiment']
28
29 # Step 4: Split the Data into Training and Testing Sets
30 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
31
```

```
31
32 # Step 5: Train a Naive Bayes Classifier
33 model = MultinomialNB()
34 model.fit(X_train, y_train)
35
36 # Step 6: Make Predictions
37 y_pred = model.predict(X_test)
38
39 # Step 7: Evaluate the Model
40 print("Confusion Matrix:")
41 print(confusion_matrix(y_test, y_pred))
42 print("\nClassification Report:")
43 print(classification_report(y_test, y_pred))
44
```

→ Confusion Matrix:

```
[[2496  0]
 [ 0 2527]]
```

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	2496
1	1.00	1.00	1.00	2527
accuracy			1.00	5023
macro avg	1.00	1.00	1.00	5023
weighted avg	1.00	1.00	1.00	5023

\*\* Conversion Rate Optimization:\*\*

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import classification_report, confusion_matrix
7
8 # Step 1: Load the Data
9 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
10
11 # Step 2: Data Preprocessing
12 # Convert event_time to datetime
13 df['event_time'] = pd.to_datetime(df['event_time'])
14
15 # Check if the dataset has required columns
16 required_columns = ['event_time', 'order_id', 'product_id', 'category_id', 'category_code', 'brand', 'price', 'user_id']
17 for col in required_columns:
18     if col not in df.columns:
19         raise ValueError(f'Missing required column: {col}')
20
21 # Step 3: Create a Conversion Label
22 # Create a label where 1 indicates a conversion (purchase made), and 0 indicates no conversion
23 df['conversion'] = df['order_id'].notnull().astype(int) # 1 if order_id exists, else 0
24
25 # Feature Engineering: Create features for the model
26 # For simplicity, let's use only product price and the brand as features
27 X = df[['price']] # Features
28 y = df['conversion'] # Target variable
29
30 # Step 4: Split the Data into Training and Testing Sets
31 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
32
33 # Step 5: Train a Random Forest Classifier
34 model = RandomForestClassifier(random_state=42)
35 model.fit(X_train, y_train)
36
37 # Step 6: Make Predictions
38 y_pred = model.predict(X_test)
39
40 # Step 7: Evaluate the Model
41 print("Confusion Matrix:")
42 print(confusion_matrix(y_test, y_pred))
43 print("\nClassification Report:")
44 print(classification_report(y_test, y_pred))
45

→ Confusion Matrix:
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:409: UserWarning: A single label was found in 'y_true' and 'y
    warnings.warn(
[[526705]]

Classification Report:
precision    recall    f1-score   support
      1       1.00      1.00      1.00      526705
accuracy                           1.00      526705
  macro avg       1.00      1.00      1.00      526705
weighted avg       1.00      1.00      1.00      526705

```

\*\* Website Performance:\*\*

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestRegressor
6 from sklearn.metrics import mean_squared_error, r2_score
7
8 # Step 1: Load the Data
9 df = pd.read_csv('/content/kz.csv') # Replace with the actual file path
10
11 # Step 2: Create Mock Website Performance Metrics
12 # For demonstration, we create random page load times and bounce rates
13 np.random.seed(42)

```

```
14 df['page_load_time'] = np.random.uniform(1, 10, size=len(df)) # Page load time in seconds
15 df['bounce_rate'] = np.random.uniform(0, 1, size=len(df)) # Bounce rate (0 to 1)
16
17 # Step 3: Create a Conversion Label
18 # Create a label where 1 indicates a purchase made (conversion), and 0 indicates no conversion
19 df['conversion'] = df['order_id'].notnull().astype(int) # 1 if order_id exists, else 0
20
21 # For the sake of this analysis, let's assume that sales can be influenced by conversions and performance metrics
22 # Calculate sales as a function of conversion and price
23 df['sales'] = df['conversion'] * df['price']
24
25 # Step 4: Feature Engineering
26 # Define features and target variable
27 X = df[['page_load_time', 'bounce_rate', 'price']] # Features
28 y = df['sales'] # Target variable
29
30 # *** Handle NaN values in the target variable 'y' ***
31 # Remove rows with NaN values in 'sales' column
32 df = df.dropna(subset=['sales'])
33
34 # Update X and y after removing NaN values
35 X = df[['page_load_time', 'bounce_rate', 'price']]
36 y = df['sales']
37
38
39 # Step 5: Split the Data into Training and Testing Sets
40 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
41
42 # Step 6: Train a Random Forest Regressor
43 model = RandomForestRegressor(random_state=42)
44 model.fit(X_train, y_train)
45
46 # Step 7: Make Predictions
47 y_pred = model.predict(X_test)
48
49 # Step 8: Evaluate the Model
50 mse = mean_squared_error(y_test, y_pred)
51 r2 = r2_score(y_test, y_pred)
52
53 print(f"Mean Squared Error: {mse:.2f}")
54 print(f"R^2 Score: {r2:.2f}")
55
56 # Optional: Feature Importance
57 importances = model.feature_importances_
58 feature_names = X.columns
59
60 # Visualization of feature importance
61 import matplotlib.pyplot as plt
62
63 plt.figure(figsize=(10, 6))
64 plt.barh(feature_names, importances)
65 plt.xlabel('Feature Importance')
66 plt.title('Feature Importance for Website Performance Impact on Sales')
67 plt.show()
```

Mean Squared Error: 0.18  
R^2 Score: 1.00

### Feature Importance for Website Performance Impact on Sales

price

#### User Engagement Prediction:

```

1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import classification_report
7
8 # Step 1: Load the Data
9 # Load your dataset; adjust the file path accordingly
10 # For example, using a local file: df = pd.read_csv('path/to/ecommerce_data.csv')
11 # Simulating a DataFrame for demonstration purposes:
12 np.random.seed(42)
13 data_size = 1000 # Number of entries in the simulated dataset
14 df = pd.DataFrame({
15     'event_time': pd.date_range(start='2023-01-01', periods=data_size, freq='H'),
16     'order_id': np.random.choice([None, 1, 2, 3], size=data_size, p=[0.7, 0.1, 0.1, 0.1]),
17     'product_id': np.random.randint(1, 50, size=data_size),
18     'category_id': np.random.randint(1, 10, size=data_size),
19     'category_code': np.random.choice(['A', 'B', 'C', 'D'], size=data_size),
20     'brand': np.random.choice(['BrandX', 'BrandY'], size=data_size),
21     'price': np.random.uniform(10, 100, size=data_size),
22     'user_id': np.random.randint(1, 100, size=data_size)
23 })
24
25 # Step 2: Feature Engineering
26 # Create a user engagement metric
27 user_engagement = df.groupby('user_id').agg({
28     'order_id': 'count', # Count of orders per user
29     'price': 'sum', # Total spent per user
30 }).reset_index()
31
32 # Rename columns for clarity
33 user_engagement.columns = ['user_id', 'order_count', 'total_spent']
34
35 # Define engagement label: engaged if they made more than 5 orders
36 user_engagement['engaged'] = (user_engagement['order_count'] > 5).astype(int)
37
38 # Step 3: Define Features and Target Variable
39 X = user_engagement[['order_count', 'total_spent']] # Features
40 y = user_engagement['engaged'] # Target variable
41
42 # Step 4: Split the Data into Training and Testing Sets
43 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
44
45 # Step 5: Train a Random Forest Classifier
46 model = RandomForestClassifier(random_state=42)
47 model.fit(X_train, y_train)
48
49 # Step 6: Make Predictions
50 y_pred = model.predict(X_test)
51
52 # Step 7: Evaluate the Model
53 print(classification_report(y_test, y_pred))
54
55 # Optional: Display feature importance
56 importances = model.feature_importances_
57 feature_names = X.columns
58 print("Feature importances:", dict(zip(feature_names, importances)))
59

```

→ <iopvthon-input-14-63defed0790e>:15: FutureWarning: 'H' is deprecated and will be removed in a future version. please use 'h' instead.