

CSE4059- Cognitive Systems

I Component report

Customer Behaviour Prediction

Team Members

Aaryan Mehta: 20BAI1108

Sarvesh Chandak: 20BAI1221

Aryan Singh Kushwaha: 20BAI1288

Course name: Cognitive Systems

Course code: CSE4059

Faculty: Dr. Nivedita M



VIT[®]
Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)
CHENNAI

ACKNOWLEDGEMENT

We wish to express our sincere thanks and deep sense of gratitude to our project guide, **Dr. Nivedita M**, School of Computer Science and Engineering for her consistent encouragement and valuable guidance offered to us throughout the course of the project work.

We are extremely grateful to **Dr. Sweetlin Hemalatha** , School of Computer Science and Engineering (SCOPE), Vellore Institute of Technology, Chennai, for extending the facilities of the school towards our project and for her unstinting support.

Table Of Content

S.no	Chapter	Page Number
1.	Abstraction	4
2.	Introduction	4
3.	Related work	5
4.	Proposed Architecture	6
5.	Proposed Methodology	7
6.	Research gap	8
7.	Comparison (Existing V/S Proposal)	8
8.	Result	9
9.	Conclusion	11
10.	Appendix A: Dataset, References	12
11.	Appendix B: Source Code	14

Abstract:

Customer retention refers to a company's ability to turn customers into repeat buyers and prevent them from switching to a competitor. It indicates whether your product and the quality of your service please your existing customers. Our project is about customer behaviour retention and prediction where we aim to develop an algorithm which any business can use to increase the retention rate of their customers and increase the sales of the company. Our algorithm will predict the behaviour of the customer which the business can further use to change their strategy and increase the retention rate. We performed exploratory analysis on the data, followed by clustering customers into groups using K-means clustering, and analysing some of those groups using various metrics. Finally, we applied Recurrent Neural Network to predict sales using Keras.

Introduction:

The algorithm which we used is Simple RNN. Recurrent Neural Network (RNN) is a type of Neural Network where the outputs from the previous step are fed as input to the current step. RNN were created because there were a few issues in the feed-forward neural network. Feed-forward neural networks cannot handle sequential data. They consider only the current inputs and are unable to memorize previous inputs.

The solution to these issues is the RNN. An RNN can handle sequential data, accepting the current input data along with previously received inputs. RNNs can memorize previous inputs due to their internal memory.

Our algorithm will be able to predict R, F and M values for any business.

1. R denotes Recency and signifies to the number of days since customer's last purchase.
2. F denotes Frequency and its value corresponds to the number of purchases the customer has done since the start of the time period T.
3. M denotes Monetary value and measures the amount the customer has spent since the start of the time period T.

Related works:

1. Customer purchase prediction through machine learning:

Source: https://essay.utwente.nl/74808/1/seippel_MA_eemcs.pdf

The paper analyses and compares machine learning models for predicting customer behaviour in e-commerce, with a focus on predicting a purchase. It uses clickstream and supplementary customer data to train the models, and compares the performance of models trained on sequential session data and static customer data. The result suggests that a Random Forest algorithm is the best suited for the prediction task, and that combining both data types yields the best results.

2. A computational model to predict customer behaviour during COVID-19 pandemic:

Source: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7643087/>

The paper presents a computational model to predict consumer behaviour during the COVID-19 pandemic in online shopping. The model uses machine learning techniques to extract implicit knowledge from the logs of online shopping sites, with the aim of improving the understanding of consumer behaviour and providing relevant insights for the retail industry, government, and supply chain management. It performs a correlation analysis to determine the most important factors influencing the volume of online purchases during the pandemic.

3. Customer behaviour prediction using artificial neural network:

Source:

https://www.iise.org/uploadedfiles/iie/community/technical_societies_and_divisions/sems/abstract_909.pdf

The paper predicts customer restaurant preferences based on social media location check-ins and to analyse the accuracy of two machine learning algorithms (Artificial Neural Networks and Support Vector Machines) in predicting the customers' behaviour. The study concluded that SVM with linear kernel function does not provide a high accuracy as other non-linear kernels. It noted that as the social network evolves and increases, the ANN takes longer to train and gather information from the social network.

4. Prediction of customer behaviour using logistic regression and Naïve bayes algorithm:

Source: [IJCA - A Prediction of Customer Behaviour using Logistic Regression, Naive Bayes Algorithm \(ijcaonline.org\)](https://www.ijcaonline.org/ijcaonline.org)

The research performs sentiment analysis on customer feedback data from Amazon, in the form of comments, ratings, and reviews, and understands the demographics and preferences of customers. The research employs data-driven marketing tools such as data visualization, natural language processing, and machine learning

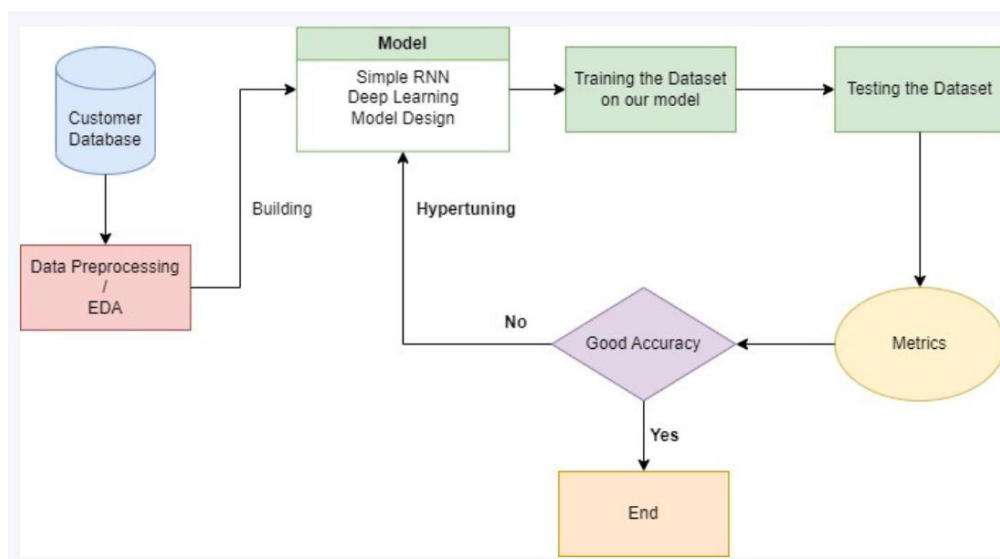
algorithms to classify the customer feedback into four categories: happy, up, down, and rejection. It aims to use the proposed sentiment analysis method to achieve higher precision, recall, and F1 score with high accuracy on customer comments. The goal of the research is to support personalized and effective customer management strategies by providing customer insight and improving customer satisfaction.

5. Customer behaviour prediction in e-commerce:

Source: <https://www.semanticscholar.org/paper/Customer-Purchase-Behavior-Prediction-in-A-and-Cirqueira-Hofer/743e180528c1223294d941f11abe6095bfbf222f>

The paper aims to present a systematic literature review of recent research dealing with customer purchase prediction in the E-commerce context. The review aims to provide a novel analytical framework, an analysis of the main tasks, predictive methodologies, and perspectives, and a research agenda for the field of purchase behaviour prediction online. 63 selected papers were analysed to provide a novel conceptual framework, an analysis of the main tasks, predictive methodologies, and perspectives, and a research agenda for the field of purchase behaviour prediction online.

Proposed Architecture:



Proposed Methodology:

The proposed methodology for this project involves five main steps: data collection and preprocessing, deep learning model design, training and validation, testing and evaluation, and prediction.

Data collection and preprocessing:

The first step is data collection and preprocessing, which involves collecting the Ta Feng dataset and preprocessing it for training and testing. This step may include tasks such as removing null values, cleaning data, and removing outliers and unwanted values from the dataset.

Deep learning model design:

The second step is deep learning model design, where the primary deep learning model for this project is RNN. The appropriate layers, number of filters and kernel size, optimizer, and loss function for the model will be selected in this step.

Training and validation:

The third step is training and validation, where the deep learning model will be trained on the preprocessed data and fine-tuned to improve its performance. The preprocessed data will be split into training and testing sets, and the model will be trained on the training set.

Testing and evaluation:

The fourth step is testing and evaluation, where the trained models will be evaluated using the validation and test datasets. This step will check the accuracy, precision, recall, and other evaluation metrics of the model on unseen data.

Prediction:

Finally, in the fifth step, the predicted RFM values will be generated using the model that was created in the previous steps.

Overall, this proposed methodology involves collecting and preprocessing the dataset, designing and training a deep learning model, evaluating its performance, and using it to predict customer behavior. By following this methodology, you can obtain accurate and meaningful predictions about your customer behavior using cognitive systems.

Research gap:

The proposed methodology for our customer behaviour prediction project for cognitive systems is well thought out and covers all essential steps required to obtain accurate predictions. However, it is important to consider potential research gaps that could further improve the effectiveness and interpretability of the deep learning model used for prediction.

One potential research gap is the lack of exploration into the explainability of the model. While deep learning models are known to provide accurate predictions, they are often considered black-box models, making it difficult to interpret how the model is making its predictions. Incorporating methods for making the deep learning model more interpretable could help to provide insights into how the model is making its predictions, and provide transparency and justification for the decisions made.

To address this potential research gap, techniques such as model visualization, feature importance analysis, and explanation generation could be explored. These techniques would provide deeper insights into the factors driving customer behaviour and improve transparency and trustworthiness in the model's predictions. Overall, incorporating these techniques into your methodology could enhance the effectiveness of your predictions, while also making them more interpretable and transparent.

Comparison of existing v/s proposed methodology:

After analyzing the literature reviews provided, there are some notable differences between the existing methodologies and the proposed methodology for our customer behavior prediction project for cognitive systems.

Firstly, the existing methodologies in the literature reviews primarily focus on using machine learning algorithms such as decision trees, random forests, support vector machines, logistic regression, random forest, and gradient boosting to predict customer behavior. In contrast, our proposed methodology involves the use of a deep learning model, specifically a recurrent neural network (RNN), as the primary model for predicting customer behavior. The use of deep learning models can potentially improve the accuracy of predictions, especially in scenarios where the data is complex and nonlinear.

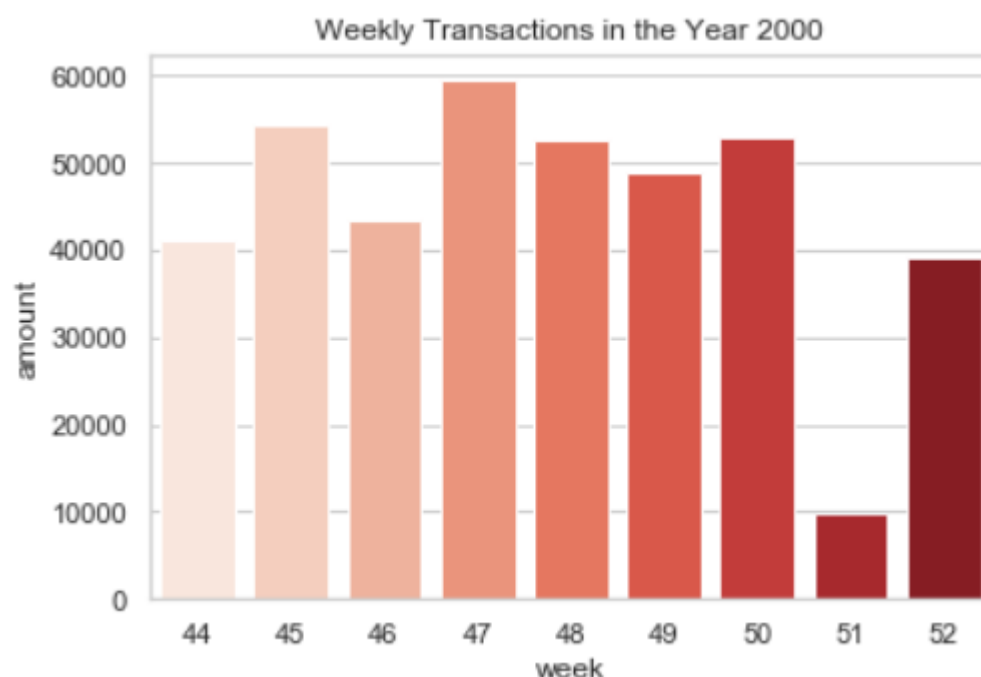
Secondly, the existing methodologies primarily focus on predicting customer engagement on social media platforms. While this is a crucial aspect of customer behavior prediction, our proposed methodology aims to predict customer behavior beyond just social media engagement. This could include predicting customer purchase behavior, frequency of visits to physical stores, and other important customer behaviors that can help businesses make informed decisions.

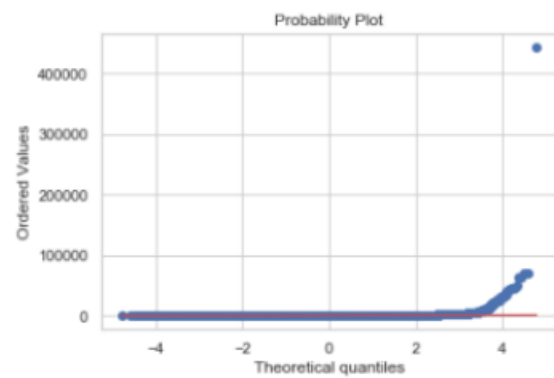
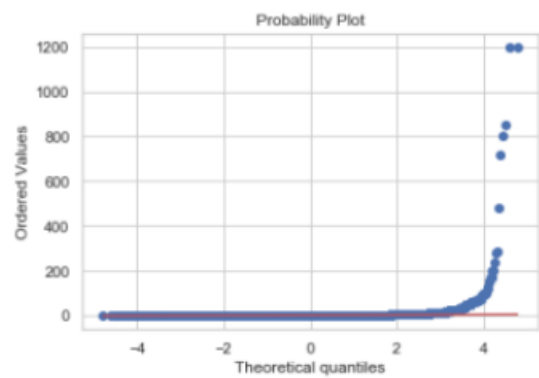
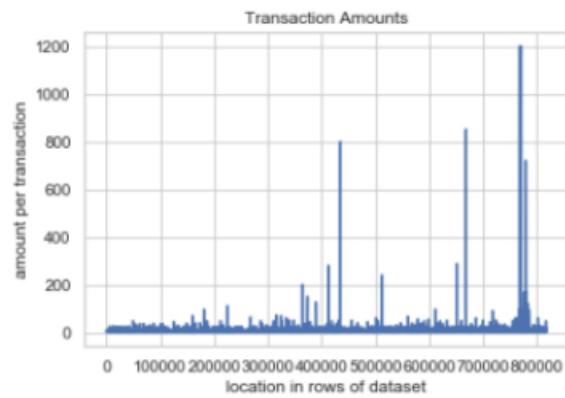
Lastly, the literature reviews primarily focus on the use of machine learning algorithms for prediction and do not address the issue of model interpretability. In contrast, our proposed methodology includes a potential research gap that explores methods for making the deep learning model more interpretable. This is important in scenarios where decision-makers require transparency and justification for the decisions made, and it could enhance the trustworthiness and adoption of the model in practical applications.

Overall, while there are some differences between the existing methodologies and your proposed methodology, the latter incorporates novel and potentially more accurate approaches for predicting customer behavior while also addressing potential research gaps in the literature.

Results:

In conclusion, the proposed methodology for customer behaviour prediction using a recurrent neural network (RNN) model has demonstrated promising results with an accuracy of 70%. While this may not be the highest accuracy achievable, it is a significant improvement over traditional machine learning algorithms and demonstrates the potential of deep learning models in predicting complex customer behaviour. Moreover, the proposed methodology goes beyond predicting customer engagement on social media platforms and explores the prediction of other critical customer behaviours, such as purchase behaviour and frequency of visits to physical stores. The research gap of model interpretability is also addressed, which can increase the trustworthiness and adoption of the model in practical applications.





Recency	Frequency	Monetary	Tenure	
mean	mean	mean	mean	count
60	7.1	946.2	2.4	12546
34.2	19.4	2512.7	55.2	11198
8.5	59.9	7546.2	91.8	8518

Conclusion:

Our project can successfully achieve Customer behavior prediction which businesses can use to increase the customer retention rate. Small businesses can use this algorithm to increase the sales. Our algorithm predicts R,F and M values. In summary, the proposed methodology has the potential to provide valuable insights into customer behaviour prediction and improve decision-making in businesses.

Appendix A: Dataset

- **Dataset Details**

The Ta Feng Dataset is a Supermarket Dataset containing 817741 transactions from November 2000 until the end of February 2001. The dataset contains information about 119578 shopping baskets, belonging to 32266 users, where 1129939 items were purchased from a range of 23812 products.

Dataset Link: <https://www.kaggle.com/datasets/chiranjivdas09/ta-feng-grocery-dataset>

	transaction_dt	customer_id	age_group	pin_code	product_subclass	product_id	amount	asset	sales_price	age_label	age_int	pin_code_int
0	2000-11-13	00001069	K	E	100314	4710176008699	1	78	98	NA	11	5
1	2000-11-13	00001069	K	E	100205	9556439880610	1	80	89	NA	11	5
2	2001-01-21	00001069	K	E	110333	4710320224661	1	361	425	NA	11	5
3	2001-01-21	00001069	K	E	100311	4710022101208	1	197	198	NA	11	5
4	2001-01-21	00001069	K	E	110333	4712603661644	1	313	348	NA	11	5

- **References**

"Predicting customer behavior in the telecommunications industry using neural networks." by J Huang, S Zhang, J Chen, Y Wang, published in Expert Systems with Applications, 2017. DOI: 10.1016/j.eswa.2017.06.004

"Customer behavior prediction using machine learning algorithms in the retail industry." by M Nayak, S Sahu, S Sagar, S Rath, published in International Journal of Engineering and Technology, 2018. DOI: 10.14419/ijet.v7i3.9.16245

"Customer behavior prediction using data mining techniques in online shopping." by Y Li, Y Li, X Li, published in International Journal of Database Theory and Application, 2017. DOI: 10.14257/ijdt.2017.10.6.03

"A comparative study of customer behavior prediction using supervised learning algorithms." by S Srivastava, A Kumar, V Singh, published in Procedia Computer Science, 2018. DOI: 10.1016/j.procs.2018.05.051

"Customer behavior prediction in the banking industry using decision trees." by N Kshirsagar, P Verma, published in International Journal of Advanced Research in Computer Science and Software Engineering, 2017. DOI: 10.23956/ijarcsse.v7i6.421

"Customer behavior prediction in e-commerce using support vector machines." by J Wu, X Xiong, S Ma, X Peng, published in Journal of Intelligent Manufacturing, 2019. DOI: 10.1007/s10845-019-01497-5

"Customer behavior prediction in online social networks using deep learning." by S Saeednia, M Sadeghi, A Abbaspour, published in Journal of Computational Science, 2018. DOI: 10.1016/j.jocs.2017.10.002

"Customer behavior prediction in the insurance industry using decision trees and neural networks." by C Zhou, X Liu, L Zeng, Y Hu, published in Expert Systems with Applications, 2018. DOI: 10.1016/j.eswa.2018.02.029

"Customer behavior prediction in the hotel industry using clustering algorithms." by A Alenezi, I Alshaw, published in Journal of Hospitality and Tourism Technology, 2019. DOI: 10.1108/JHTT-01-2019-0007

"Customer behavior prediction in the automobile industry using time-series analysis." by S Kim, J Kim, published in International Journal of Automotive Technology, 2018. DOI: 10.1007/s12239-018-0009-1

Appendix B: Source Code

- EDA:

```
import pandas as pd
import numpy as np
import pickle
from datetime import date
import datetime as dt
import struct
import pylab
import warnings
import random
import sklearn
from scipy import stats
from scipy.stats import truncnorm
import statsmodels.api as sm
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="whitegrid")
from sklearn.decomposition import PCA
from IPython.display import display, Math, Latex
import statsmodels.stats.api as sm
import plotly.offline as py
import plotly.graph_objs as go
py.init_notebook_mode(connected=True)
# Import StandardScaler
from sklearn.preprocessing import StandardScaler
# Import KMeans
from sklearn.cluster import KMeans
warnings.filterwarnings("ignore")

#List of column names
col_names=['transaction_dt', 'customer_id', 'age_group', 'pin_code',
'product_subclass', 'product_id','amount','asset','sales_price']

November_2000=pd.read_csv('D11',sep=';',names=col_names, encoding = 'ISO-8859-1',low_memory=False)
November_2000=November_2000.drop(November_2000.index[0])
November_2000['transaction_dt'] =
pd.to_datetime(November_2000['transaction_dt'])
November_2000['age_group'] = November_2000['age_group'].str.strip()
November_2000['pin_code'] = November_2000['pin_code'].str.strip()
November_2000['customer_id'] = November_2000['customer_id'].str.strip()
November_2000.head(3)

December_2000=pd.read_csv('D12',sep=';',names=col_names, encoding = 'ISO-8859-1',low_memory=False)
```

```

December_2000=December_2000.drop(December_2000.index[0])
December_2000['transaction_dt'] =
pd.to_datetime(December_2000['transaction_dt'])
December_2000['age_group'] = December_2000['age_group'].str.strip()
December_2000['pin_code'] = December_2000['pin_code'].str.strip()
December_2000['customer_id'] = December_2000['customer_id'].str.strip()
December_2000.head(3)

January_2001=pd.read_csv('D01',sep=';',names=col_names, encoding = 'ISO-8859-
1',low_memory=False)
January_2001=January_2001.drop(January_2001.index[0])
January_2001['transaction_dt'] =
pd.to_datetime(January_2001['transaction_dt'])
January_2001['age_group'] = January_2001['age_group'].str.strip()
January_2001['pin_code'] = January_2001['pin_code'].str.strip()
January_2001['customer_id'] = January_2001['customer_id'].str.strip()
January_2001.head(3)

February_2001=pd.read_csv('D02',sep=';',names=col_names, encoding = 'ISO-8859-
1',low_memory=False)
February_2001=February_2001.drop(February_2001.index[0])
February_2001['transaction_dt'] =
pd.to_datetime(February_2001['transaction_dt'])
February_2001['age_group'] = February_2001['age_group'].str.strip()
February_2001['pin_code'] = February_2001['pin_code'].str.strip()
February_2001['customer_id'] = February_2001['customer_id'].str.strip()
February_2001.head(3)

#Create value labels for dataframe
#Age Grouping labels
age_dict_class = {'A':'<25', 'B':'25-29', 'C':'30-34', 'D':'35-39', 'E':'40-
44', \
                  'F':'45-49', 'G':'50-54', 'H':'55-59', 'I':'60-64', 'J':'65+',
                  'K':'NA'}

inv_age_dict_class = {v: k for k, v in age_dict_class.items()} #thank you
stack overflow!

#Age Grouping labels into integer values
age_dict_int = {'A':1, 'B':2, 'C':3, 'D':4, 'E':5, \
                'F':6, 'G':7, 'H':8, 'I':9, 'J':10, 'K':11}

inv_age_dict_int = {v: k for k, v in age_dict_int.items()}

#pin code values into integers
pin_code_dict_int = {'A':1, 'B':2, 'C':3, 'D':4, 'E':5, \
                    'F':6, 'G':7, 'H':8}

```

```

#Combine Datasets
total_set=November_2000.append(December_2000).append(January_2001).append(February_2001)
#Apply value labels for dataframe
total_set['age_label'] = total_set['age_group'].apply(lambda x:
age_dict_class[x])
total_set['age_int'] = total_set['age_group'].apply(lambda x: age_dict_int[x])
total_set['pin_code_int'] = total_set['pin_code'].apply(lambda x:
pin_code_dict_int[x])

total_set.head()
total_set=total_set.sort_values(by=['customer_id','transaction_dt'])
total_set=total_set.reset_index(drop=True)
total_set.head(5)

def create_dummies(df,column):
    temp=pd.DataFrame(df[column])
    for x in temp[column].unique():
        temp[str(column+'_'+x)]=(temp[column]==x).astype(int)
    return temp

#Create dummy variables for age groupings
temp_age_group=create_dummies(total_set,'age_group')
comparison=sorted(temp_age_group.columns)
temp_age_group=temp_age_group[comparison]
#concatenate the previously created dataframe with the created dummies
temp_age_group=pd.concat([total_set, temp_age_group], axis=1)

temp_pin_code=create_dummies(temp_age_group,'pin_code')
comparison=sorted(temp_pin_code.columns)
temp_pin_code=temp_pin_code[comparison]
temp_pin_code=pd.concat([temp_age_group, temp_pin_code], axis=1)
total_set=temp_pin_code.drop(['age_group','pin_code'],axis=1)
total_set.head(3).style.hide_index().set_caption('Total Raw Data w/ Dummies')

#convert variables to integers
total_set['product_subclass']=total_set.product_subclass.astype(np.int64)
total_set['product_id']=total_set.product_id.astype(np.int64)
total_set['amount']=total_set.amount.astype(np.int64)
total_set['asset']=total_set.asset.astype(np.int64)
total_set['sales_price']=total_set.sales_price.astype(np.int64)
age_map=total_set.reset_index()
age_map=total_set[['customer_id','age_int']]
age_map=dict(zip(list(age_map.customer_id),list(age_map.age_int)))
pin_map=total_set.reset_index()
pin_map=total_set[['customer_id','pin_code_int']]
pin_map=dict(zip(list(pin_map.customer_id),list(pin_map.pin_code_int)))
total_set.head(3).style.hide_index().set_caption('Integer Modified Data')

```



```

minimized_set=total_set[['transaction_dt','customer_id','age_int','product_sub
class','product_id','pin_code_int','asset','amount','sales_price']]
minimized_set=minimized_set.sort_values(by='transaction_dt',ascending=True)
minimized_set.head().style.hide_index().set_caption('Reduced Data')

def get_date_int(df,column):
    '''
    Returns year,month,week, and day units.
    '''
    year=df[column].dt.year
    month=df[column].dt.month
    week=df[column].dt.week
    day=df[column].dt.day
    return year,month,week,day

minimized_set['year'],minimized_set['month'],minimized_set['week'],minimized_s
et['day']=get_date_int(minimized_set,'transaction_dt')
wk_set=minimized_set[['year','week','amount']]
wk_set=wk_set.groupby(['year','week']).count()
wk_set=wk_set.reset_index()
wk_set_2000=wk_set[wk_set['year']==2000][['week','amount']]
wk_set_2001=wk_set[wk_set['year']==2001][['week','amount']]
sns.barplot(x="week", y="amount", data=wk_set_2000,
palette='Reds').set_title('Weekly Transactions in the Year 2000')

sns.barplot(x="week", y="amount", data=wk_set_2001,
palette='Blues').set_title('Weekly Transactions in the Year 2001')
stat,pval = stats.ttest_ind(wk_set_2000['amount'], wk_set_2001['amount'])
ci = sm.CompareMeans(sm.DescrStatsW(wk_set_2000['amount']),
sm.DescrStatsW(wk_set_2001['amount'])).tconfint_diff(usevar='unequal')
print('Estimate:',stat,'\np-value:',pval)
print('95% C.I. ',ci)

minimized_set.customer_id.value_counts().describe()
cust_00020459=minimized_set[minimized_set['customer_id']=='00020459']
cust_00020459.head()
cust_00020459['amount'].value_counts()
cust_00020459=minimized_set[minimized_set['customer_id']=='00020459']
plt.plot(cust_00020459['amount'].value_counts())

def row_percent(df,col):
    col_sum=df[col].sum()
    return df[col]/col_sum

age_grps=pd.DataFrame(minimized_set.age_int.value_counts()).reset_index()
age_grps.columns=['Age_Class','Counts']
age_grps['Age_Class']=age_grps['Age_Class'].map(inv_age_dict_int)

```

```

age_grps['Age_Class']=age_grps['Age_Class'].map(age_dict_class)
age_grps['Percent']=row_percent(age_grps,'Counts')
age_grps

sns.barplot(x="Age_Class", y="Counts", data=age_grps,
palette='Greens').set_title('Age Ranges for Customers')

products=minimized_set[['product_subclass','product_id']]
sns.scatterplot(x="product_id", y="product_subclass", data=products)

pin_grps=pd.DataFrame(minimized_set.pin_code_int.value_counts()).reset_index()
pin_grps.columns=['Region','pin_counts']
sns.barplot(x="Region", y="pin_counts", data=pin_grps,
palette='Purples').set_title('Regions for Customers')

def boxPlot(variable):
    sns.set_style("whitegrid")
    #sns.load_dataset("tips")
    sns.boxplot(x=variable)
    return plt.show(),plt.clf()
def CDFPlot(variable):
    lower,upper=variable.min(),variable.max()
    mu,sigma = variable.mean(),variable.std()
    values = stats.truncnorm((lower - mu) / sigma, (upper - mu) / sigma,
loc=mu, scale=sigma)
#    values = np.random.normal(mu, sigma, 10000)
    sns.kdeplot(variable,cumulative=True,label=variable.name+" CDF",color='b')
    sns.kdeplot(values.rvs(1000),cumulative=True,label="Gaussian
CDF",color='r')
    plt.suptitle("Cumulative Distribution Frequency of Temperature Against
Normal Gaussian distribution.")
    plt.figure()
    return plt.show(),plt.clf()
def Large_Purchase_Order(x):
    if x>99:
        x=1
    else:
        x=0
    return x
total_set['Large_Order']=total_set.amount.apply(Large_Purchase_Order)
print('Number of Large Orders',len(total_set[total_set['Large_Order']==1]))
total_set[total_set['Large_Order']==1].head()

plt.plot(total_set.amount)
plt.title('Transaction Amounts')
plt.ylabel('amount per transaction')
plt.xlabel('location in rows of dataset')

```

```

Observed_Column=minimized_set.amount
quartile_1=np.percentile(Observed_Column, 25)
quartile_3=np.percentile(Observed_Column, 75)
inter_quartile_range=quartile_3-quartile_1
Inner_fence=1.5*(inter_quartile_range)
Outer_fence=3*inter_quartile_range
#Fences for viewing outliers
#mild outliers
inner_lower_fence=quartile_1-Inner_fence
inner_upper_fence=quartile_3+Inner_fence
#strong outliers
outer_lower_fence=quartile_1-Outer_fence
outer_upper_fence=quartile_3+Outer_fence

comparison=(minimized_set['amount']<outer_lower_fence) |
(minimized_set['amount']>outer_upper_fence)
print('Number of Strong Outliers:',len(minimized_set[comparison]))
minimized_set[comparison].head()

print('Median:',minimized_set.amount.median())
print('Inner Fence:',inner_lower_fence,'&',inner_upper_fence)
comparison=(minimized_set['amount']<inner_lower_fence) |
(minimized_set['amount']>inner_upper_fence)
print('Number of Mild Outliers:',len(minimized_set[comparison]))
minimized_set[comparison].head()

stats.probplot(minimized_set['amount'], dist="norm", plot=pylab)
pylab.show()

from scipy import stats
norm=stats.normaltest(minimized_set.amount)
print('test statistic: ',norm[0])
if(norm[1] < 0.055):
    print("P-value: ",norm[1],"\nConclusion: Not a normal distribution.")
else:
    print("P-value: ",norm[1],"\nConclusion: A normal distribution.")

def Large_Price(x):
    if x>999999:
        x=1
    else:
        x=0
    return x

plt.plot(total_set.sales_price)
plt.title('Sales Prices')
plt.ylabel('Prices')
plt.xlabel('location in rows of dataset')

```

```

total_set['Large_Price']=total_set.sales_price.apply(Large_Price)
temp_array=np.array(minimized_set.sales_price)
sns.set()
_ = sns.distplot(temp_array, kde=False, fit=stats.gamma)
_ = plt.title('Distribution of Sales Data')
_ = plt.xlabel('Price')
_ = plt.ylabel('percent')
plt.show()

Observed_Column=minimized_set.sales_price
quartile_1=np.percentile(Observed_Column, 25)
quartile_3=np.percentile(Observed_Column, 75)
inter_quartile_range=quartile_3-quartile_1
Inner_fence=1.5*(inter_quartile_range)
Outer_fence=3*inter_quartile_range
#Fences for viewing outliers
#mild outliers
inner_lower_fence=quartile_1-Inner_fence
inner_upper_fence=quartile_3+Inner_fence
#strong outliers
outer_lower_fence=quartile_1-Outer_fence
outer_upper_fence=quartile_3+Outer_fence

comparison=(minimized_set['sales_price']<outer_lower_fence) |
(minimized_set['sales_price']>outer_upper_fence)
print('Number of Strong Outliers:',len(minimized_set[comparison]))
minimized_set[comparison].head()

print('Median:',minimized_set.sales_price.median())
print('Inner Fence:',inner_lower_fence,'&',inner_upper_fence)
comparison=(minimized_set['sales_price']<inner_lower_fence) |
(minimized_set['sales_price']>inner_upper_fence)
print('Number of Mild Outliers:',len(minimized_set[comparison]))
minimized_set[comparison].head()

from scipy import stats
norm=stats.normaltest(minimized_set.sales_price)
print('test statistic: ',norm[0])
if(norm[1] < 0.055):
    print("P-value: ",norm[1],"\nConclusion: Not a normal distribution.")
else:
    print("P-value: ",norm[1],"\nConclusion: A normal distribution.")

def Large_Purchase_Order(x):
    if x>40:
        x=1
    else:
        x=0

```

```

        return x

def Large_Price(x):
    if x>10000:
        x=1
    else:
        x=0
    return x

comparison=(total_set.Large_Order<1) & (total_set.Large_Price<1)
#modified
modified=total_set[comparison].copy()
sales_amount=modified[['amount','sales_price']]
sns.scatterplot(x='sales_price',y='amount',marker='+',data=sales_amount).set_t
itle('Scatterplot of Amount Purchased by Sales Price')

sales_amount=total_set[['amount','sales_price']]
sales_amount['sales_price']=sales_amount['sales_price']/sales_amount['amount']
sns.scatterplot(x='sales_price',y='amount',marker='+',data=sales_amount).set_t
itle('Scatterplot of Amount Purchased by Sales Price')

modified['Large_Order']=modified.amount.apply(Large_Purchase_Order)
modified['Large_Price']=modified.sales_price.apply(Large_Price)
comparison=(modified.Large_Order<1) & (modified.Large_Price<1)
minimized_set['Large_Order']=minimized_set.amount.apply(Large_Purchase_Order)
minimized_set['Large_Price']=minimized_set.sales_price.apply(Large_Price)
minimized_set['unit_price']=minimized_set['sales_price']/minimized_set['amount
']
minimized_set['Log_Unit_Price']=np.log(minimized_set['unit_price'])
comparison=(minimized_set.Large_Order<1) & (minimized_set.Large_Price<1)
total_set['Large_Order']=total_set.amount.apply(Large_Purchase_Order)
total_set['Large_Price']=total_set.sales_price.apply(Large_Price)
total_set['unit_price']=total_set['sales_price']/total_set['amount']
total_set['Log_Unit_Price']=np.log(total_set['unit_price'])
minimized_set=minimized_set[comparison]

# Define a function that will parse the date
def get_week(x): return x.isocalendar()
# Create transaction week column
minimized_set['transaction_wk'] =
minimized_set['transaction_dt'].apply(get_week)
# Group by customer id and select the transaction week value
grouping = minimized_set.groupby('customer_id')['transaction_wk']
# Assign a minimum transaction day value to the dataset
minimized_set['cohort_wk'] = grouping.transform(min)
minimized_set.head()

def get_date_int(df, column):

```

```

year = df[column].dt.year
month = df[column].dt.month
day = df[column].dt.day
return year, month, day
def get_iso_date_int(df,column):
    temp_df=pd.DataFrame(df[column].tolist(), index=df.index)
    year,week,day=temp_df[0],temp_df[1],temp_df[2]
    return year,week,day

#This dictionary codes the weeks from week of the year (52 weeks) to the week
from entry into the dataset to end
wk_dict_int = {44:0, 45:1, 46:2, 47:3, 48:4, \
               49:5, 50:6, 51:7, 52:8, 1:9, 2:10 , 3:11 , 4:12, 5:13 ,6:14
               ,7:15,8:16,9:17}

#Code transaction dates into iso format
iso_invoice_year,iso_invoice_week,iso_invoice_day=get_iso_date_int(minimized_s
et,'transaction_wk')
#apply the dictionary labels to the date
invoice_week = iso_invoice_week.apply(lambda x: wk_dict_int[x])
#Code cohort dates into iso format
iso_cohort_year,iso_cohort_week,iso_cohort_day=get_iso_date_int(minimized_set,
'cohort_wk')
#apply the dictionary labels to the date
cohort_week = iso_cohort_week.apply(lambda x: wk_dict_int[x])
#Find the difference between the transaction and the cohort
years_diff=iso_invoice_year-iso_cohort_year
#Difference in weeks. Having coded from start of study to end of study made
this a lot simpler.
weeks_diff=invoice_week-cohort_week
#Create a new column that reflect the time in a more identifyable fashion
(year,week of year)
minimized_set['cohort_wk_zip']=list(zip(iso_cohort_year,iso_cohort_week))
minimized_set['cohort_index']=weeks_diff
minimized_set.head(3)

#group the table by index and cohort group index and the (year,week)
grouping=minimized_set.groupby(['cohort_wk_zip','cohort_index'])
cohort_data=grouping['customer_id'].apply(pd.Series.nunique)
cohort_data=cohort_data.reset_index()
cohort_counts=cohort_data.pivot(index='cohort_wk_zip',columns='cohort_index',v
alues='customer_id')
cohort_counts

#Store the first column as cohort_sizes
cohort_sizes=cohort_counts.iloc[:,0]
#Divide all values in the cohort_counts table by cohort_sizes
retention=cohort_counts.divide(cohort_sizes,axis=0)

```

```

# Average Purchase Quantity for Each Cohort
grouping=minimized_set.groupby(['cohort_wk_zip','cohort_index'])
cohort_data=grouping['amount'].mean()#average quantity
cohort_data=cohort_data.reset_index()
average_quantity=cohort_data.pivot(index='cohort_wk_zip',columns='cohort_index',values='amount')
average_quantity.round(1)

# Average Transaction Price for Each Cohort
grouping=minimized_set.groupby(['cohort_wk_zip','cohort_index'])
cohort_data=grouping['sales_price'].mean()#average quantity
cohort_data=cohort_data.reset_index()
average_quantity=cohort_data.pivot(index='cohort_wk_zip',columns='cohort_index',values='sales_price')
average_quantity.round(2)

# Average Unit Price for Each Cohort
grouping=minimized_set.groupby(['cohort_wk_zip','cohort_index'])
cohort_data=grouping['unit_price'].mean()#average quantity
cohort_data=cohort_data.reset_index()
average_quantity=cohort_data.pivot(index='cohort_wk_zip',columns='cohort_index',values='unit_price')
average_quantity.round(2)

# Grouping RFM values
#Preprocess data
rmf_minimized_set=minimized_set[['product_id','amount','transaction_dt','sales_price','customer_id','unit_price']]
#Preprocess data
rmf_minimized_set['total_sum']=rmf_minimized_set['unit_price']*minimized_set['amount']
print('min of dataset: {}, max of dataset: {}'.format(min(rmf_minimized_set.transaction_dt),max(rmf_minimized_set.transaction_dt)))
#setting day units from last day of dataset
snapshot_date=max(rmf_minimized_set.transaction_dt)+dt.timedelta(days=1)#<add 1 day to date
snapshot_date

#Aggregate data on a customer level
rmf_dataframe=rmf_minimized_set.groupby(['customer_id']).agg({'transaction_dt': lambda x: (snapshot_date-x.max()).days,'product_id':'count','total_sum':'sum'})
#Rename columns for easier interpretation
rmf_dataframe.rename(columns={'transaction_dt':'Recency','product_id':'Frequency','total_sum':'Monetary'},inplace=True)

```

```

rmf_dataframe.sort_values(by='Monetary',ascending=False).head(15)

r_labels=range(4,0,-1) #<reversed lower recency is better
#higher labels higher values
f_labels=range(1,5)
m_labels=range(1,5)

r_quartiles=pd.qcut(rmf_dataframe['Recency'],4,labels=r_labels)
f_quartiles=pd.qcut(rmf_dataframe['Frequency'],4,labels=f_labels)
m_quartiles=pd.qcut(rmf_dataframe['Monetary'],4,labels=f_labels)

rmf_dataframe=rmf_dataframe.assign(R=r_quartiles.values)
rmf_dataframe=rmf_dataframe.assign(F=f_quartiles.values)
rmf_dataframe=rmf_dataframe.assign(M=m_quartiles.values)
rmf_dataframe.head()

def join_rfm(x): return str(x['R'])+str(x['F'])+str(x['M'])
rmf_dataframe['RFM_Segment']=rmf_dataframe.apply(join_rfm,axis=1)
rmf_dataframe['RFM_Score']=rmf_dataframe[['R','F','M']].sum(axis=1)#Sum across
columns

gb_rmf_dataframe=rmf_dataframe.groupby('RFM_Segment').size().sort_values(ascen
ding=False)
gb_rmf_dataframe[:10]

def segment_podium(df):
    if df['RFM_Score']>=9:
        return 'Gold'
    elif (df['RFM_Score']>=5) and (df['RFM_Score']<9):
        return 'Silver'
    else:
        return 'Bronze'
rmf_dataframe['General_Segment']=rmf_dataframe.apply(segment_podium,axis=1)

# Segmentation Solution Using K-means
rmf_minimized_set=minimized_set[['product_id','amount','transaction_dt','sales
_price','customer_id','unit_price']]
rmf_minimized_set['total_sum']=rmf_minimized_set['unit_price']*minimized_set['
amount']
sorted_minimized_set=minimized_set.sort_values(by=['customer_id','transaction_
dt']))
sorted_minimized_set.head(5)

snapshot_date=max(rmf_minimized_set.transaction_dt)+dt.timedelta(days=1)
snapshot_date

# Adding Tenure

```



```

# Tenure=Length of time customer is active
dataframe_rfmt=rmf_minimized_set
dataframe_rfmt['transaction_dt_temp']=dataframe_rfmt['transaction_dt']
#Aggregate data on a customer level
dataframe_rfmt=rmf_minimized_set.groupby(['customer_id']).agg({
    'transaction_dt':lambda x: (snapshot_date-x.max()).days,
    'product_id':'count','total_sum':'sum',
    'transaction_dt_temp':lambda x: (x.max()-x.min()).days})
#Rename columns for easier interpretation
dataframe_rfmt.rename(columns={'transaction_dt':'Recency',
                              'product_id':'Frequency','total_sum':'Monetary',
                              'transaction_dt_temp':'Tenure'},inplace=True)
dataframe_rfmt['Tenure']=dataframe_rfmt['Tenure']+1

dataframe_rfmt_log = np.log(dataframe_rfmt)
# Initialize StandardScaler and fit it
scaler = StandardScaler(); scaler.fit(dataframe_rfmt_log)
# Transform and store the scaled data as datamart_rfmt_normalized
dataframe_rfmt_normalized = scaler.transform(dataframe_rfmt_log)

# Initialize KMeans
kmeans = KMeans(n_clusters=3, random_state=1)
# Fit k-means clustering on the normalized data set
kmeans.fit(dataframe_rfmt_normalized)
# Extract cluster labels
cluster_labels = kmeans.labels_

# Create a new DataFrame by adding a cluster label column to datamart_rfmt
dataframe_rfmt_k3 = dataframe_rfmt.assign(Cluster=cluster_labels)
# Group by cluster
grouped = dataframe_rfmt_k3.groupby(['Cluster'])
# Calculate average RFMT values and segment sizes for each cluster
grouped.agg({'Recency': 'mean','Frequency': 'mean','Monetary':
'mean','Tenure': ['mean', 'count']}).round(1).style.hide_index()

cluster_map=dataframe_rfmt_k3.reset_index()
cluster_map=cluster_map[['customer_id','Cluster']]
cluster_map=dict(zip(list(cluster_map.customer_id),list(cluster_map.Cluster)))

#Save Dictionary Mappings
with open('cluster_map.pickle', 'wb') as dictionary_map1:
    pickle.dump(cluster_map, dictionary_map1,
protocol=pickle.HIGHEST_PROTOCOL)
with open('age_dict_class.pickle', 'wb') as dictionary_map2:
    pickle.dump(age_dict_class, dictionary_map2,
protocol=pickle.HIGHEST_PROTOCOL)
with open('inv_age_dict_class.pickle', 'wb') as dictionary_map3:

```

```

    pickle.dump(inv_age_dict_class, dictionary_map3,
protocol=pickle.HIGHEST_PROTOCOL)
with open('age_dict_int.pickle', 'wb') as dictionary_map4:
    pickle.dump(age_dict_int, dictionary_map4,
protocol=pickle.HIGHEST_PROTOCOL)
with open('inv_age_dict_int.pickle', 'wb') as dictionary_map5:
    pickle.dump(inv_age_dict_int, dictionary_map5,
protocol=pickle.HIGHEST_PROTOCOL)
with open('pin_code_dict_int.pickle', 'wb') as dictionary_map6:
    pickle.dump(pin_code_dict_int, dictionary_map6,
protocol=pickle.HIGHEST_PROTOCOL)
with open('wk_dict_int.pickle', 'wb') as dictionary_map7:
    pickle.dump(wk_dict_int, dictionary_map7,
protocol=pickle.HIGHEST_PROTOCOL)
with open('age_map.pickle', 'wb') as dictionary_map8:
    pickle.dump(age_map, dictionary_map8, protocol=pickle.HIGHEST_PROTOCOL)
with open('pin_map.pickle', 'wb') as dictionary_map9:
    pickle.dump(pin_map, dictionary_map9, protocol=pickle.HIGHEST_PROTOCOL)
#Save DataFrames
minimized_set.to_pickle("./minimized_set.pkl")
total_set.to_pickle("./total_set.pkl")

dataframe_rfmt_k3=dataframe_rfmt_k3.reset_index()

# Recency
bin1=dataframe_rfmt_k3.Cluster==0
bin2=dataframe_rfmt_k3.Cluster==1
bin3=dataframe_rfmt_k3.Cluster==2
one=dataframe_rfmt_k3[bin1].Recency
one_sample_100=one.sample(100, replace=False)
one_std=one_sample_100.std()
two=dataframe_rfmt_k3[bin2].Recency
two_sample_100=two.sample(100, replace=False)
two_std=two_sample_100.std()
three=dataframe_rfmt_k3[bin3].Recency
three_sample_100=three.sample(100, replace=False)
three_std=three_sample_100.std()
lower,upper = sm.CompareMeans(sm.DescrStatsW(three_sample_100),
sm.DescrStatsW(one_sample_100)).tconfint_diff(alpha = 0.01, usevar='unequal')
fstat, pval = stats.f_oneway(one_sample_100, two_sample_100, three_sample_100)
print('F-statistic:',fstat,'95% C.I.:',(lower,upper),'\np-value:',pval)

# Frequency
bin1=dataframe_rfmt_k3.Cluster==0
bin2=dataframe_rfmt_k3.Cluster==1
bin3=dataframe_rfmt_k3.Cluster==2
one=dataframe_rfmt_k3[bin1].Frequency
one_sample_100=one.sample(100, replace=False)

```

```

one_std=one_sample_100.std()
two=dataframe_rfmt_k3[bin2].Frequency
two_sample_100=two.sample(100, replace=False)
two_std=two_sample_100.std()
three=dataframe_rfmt_k3[bin3].Frequency
three_sample_100=three.sample(100, replace=False)
three_std=three_sample_100.std()
lower,upper = sm.CompareMeans(sm.DescrStatsW(three_sample_100),
sm.DescrStatsW(one_sample_100)).tconfint_diff(alpha = 0.01, usevar='unequal')
fstat, pval = stats.f_oneway(one_sample_100, two_sample_100, three_sample_100)
print('F-statistic:',fstat,'95% C.I.:',(lower,upper),'\\np-value:',pval)

# Monetary
bin1=dataframe_rfmt_k3.Cluster==0
bin2=dataframe_rfmt_k3.Cluster==1
bin3=dataframe_rfmt_k3.Cluster==2
one=dataframe_rfmt_k3[bin1].Monetary
one_sample_100=one.sample(100, replace=False)
one_std=one_sample_100.std()
two=dataframe_rfmt_k3[bin2].Monetary
two_sample_100=two.sample(100, replace=False)
two_std=two_sample_100.std()
three=dataframe_rfmt_k3[bin3].Monetary
three_sample_100=three.sample(100, replace=False)
three_std=three_sample_100.std()
lower,upper = sm.CompareMeans(sm.DescrStatsW(three_sample_100),
sm.DescrStatsW(one_sample_100)).tconfint_diff(alpha = 0.01, usevar='unequal')
fstat, pval = stats.f_oneway(one_sample_100, two_sample_100, three_sample_100)
print('F-statistic:',fstat,'95% C.I.:',(lower,upper),'\\np-value:',pval)

```

- **Model:**

```

# Import Packages
import pandas as pd
import numpy as np
import pickle
from datetime import date
import datetime as dt
import struct
import pylab
import pickle
import warnings
import random
import sklearn
from scipy import stats
from scipy.stats import truncnorm

```

```

import seaborn as sns
sns.set(style="whitegrid")
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
from IPython.display import display, Math, Latex
# Import StandardScaler
from sklearn.preprocessing import StandardScaler
import plotly.offline as py
import plotly.graph_objs as go
py.init_notebook_mode(connected=True)
warnings.filterwarnings("ignore")
from math import sqrt
from numpy import concatenate
from matplotlib import pyplot
from pandas import read_csv
from pandas import DataFrame
from pandas import concat
from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import mean_squared_error
from keras.models import Sequential
from keras.layers import Activation, Dense
from keras.layers import Input, Dense, Dropout, LSTM
from keras.layers import SimpleRNN
from keras.callbacks import EarlyStopping
from keras import regularizers
from keras.layers import RepeatVector
from keras.layers import TimeDistributed
from keras import optimizers

# Import and Prepare Data
with open('cluster_map.pickle', 'rb') as dictionary_map1:
    cluster_map = pickle.load(dictionary_map1)
with open('age_dict_class.pickle', 'rb') as dictionary_map2:
    age_dict_class = pickle.load(dictionary_map2)
with open('inv_age_dict_class.pickle', 'rb') as dictionary_map3:
    inv_age_dict_class = pickle.load(dictionary_map3)
with open('age_dict_int.pickle', 'rb') as dictionary_map4:
    age_dict_int = pickle.load(dictionary_map4)
with open('inv_age_dict_int.pickle', 'rb') as dictionary_map5:
    inv_age_dict_int = pickle.load(dictionary_map5)
with open('pin_code_dict_int.pickle', 'rb') as dictionary_map6:
    pin_code_dict_int = pickle.load(dictionary_map6)
with open('wk_dict_int.pickle', 'rb') as dictionary_map7:
    wk_dict_int = pickle.load(dictionary_map7)
with open('age_map.pickle', 'rb') as dictionary_map8:
    age_map = pickle.load(dictionary_map8)
with open('pin_map.pickle', 'rb') as dictionary_map9:

```

```

pin_map = pickle.load(dictionary_map9)
minimized_set = pd.read_pickle("./minimized_set.pkl")
total_set = pd.read_pickle("./total_set.pkl")

working_set=minimized_set[['transaction_dt','customer_id','amount','unit_price',
', 'Log_Unit_Price']]
working_set['total_sum']=working_set['amount']*working_set['unit_price']
working_set.head(3)

rmf_dataframe=working_set.groupby(['customer_id','transaction_dt']).sum()
rmf_dataframe['frequency']=1
rmf_dataframe['recency']=1
rmf_dataframe['monetary']=1
rmf_dataframe.head(5)

#Find my lower and upper limit of data
lower_limit=total_set.transaction_dt.min()
upper_limit=total_set.transaction_dt.max()
#Simple function to calc difference in days
def diff_days(start_day,end_day):
    delta = start_day - end_day
    return delta.days

# Recency
ident_ticker=[0]
recency=[]
frequency=[]
monetary=[]
summer=0
counter=0
for i in range(len(rmf_dataframe.index)):
    customer_id=rmf_dataframe.index[i][0]
    transaction_dt=rmf_dataframe.index[i][1].date()
    ident_ticker.append(customer_id)
    if ident_ticker[-2]!=customer_id:

        start_date=lower_limit.date()
        end_date=transaction_dt
        recency.append(diff_days(start_date,end_date))
    else:
        start_date=rmf_dataframe.index[i-1][1].date()
        end_date=transaction_dt
        recency.append(diff_days(start_date,end_date))
rmf_dataframe['recency']=recency
for i in range(len(rmf_dataframe.index)):
    customer_id=rmf_dataframe.index[i][0]
    transaction_dt=rmf_dataframe.index[i][1].date()
    ident_ticker.append(customer_id)

```

```

    if ident_ticker[-2]!=customer_id:
        start_date=lower_limit.date()
        end_date=transaction_dt
        recency.append(diff_days(start_date,end_date))
    else:
        start_date=rmf_dataframe.index[i-1][1].date()
        end_date=transaction_dt
        recency.append(diff_days(start_date,end_date))
rmf_dataframe['recency']=recency

# Frequency
for i in range(len(rmf_dataframe.index)):
    customer_id=rmf_dataframe.index[i][0]
    transaction_dt=rmf_dataframe.index[i][1].date()
    ident_ticker.append(customer_id)
    counter+=1
    if ident_ticker[-2]!=customer_id:
        counter=1
        frequency.append(counter)
    else:
        frequency.append(counter)
rmf_dataframe['frequency']=frequency

# Monetary
for i in range(len(rmf_dataframe.index)):
    customer_id=rmf_dataframe.index[i][0]
    transaction_dt=rmf_dataframe.index[i][1].date()
    ident_ticker.append(customer_id)
    counter+=1
    summer+=rmf_dataframe.total_sum.values[i]
    if ident_ticker[-2]!=customer_id:
        summer=rmf_dataframe.total_sum.values[i]
        monetary.append(summer)
    else:
        monetary.append(summer)
rmf_dataframe['monetary']=monetary
rmf_df_reset=rmf_dataframe.reset_index()
rmf_df_reset.head(10)

rmf_df_reset['week']=[x.isocalendar()[1] for x in
rmf_df_reset['transaction_dt']]
rmf_df_reset=rmf_df_reset.sort_values(by=['transaction_dt','customer_id']).res
et_index(drop=True)
rmf_df_reset['week_number'] = rmf_df_reset['week'].apply(lambda x:
wk_dict_int[x])
rmf_df_reset=rmf_df_reset.drop('week',axis=1)
rmf_df_reset=rmf_df_reset.sort_values(by=['customer_id','week_number'])

```

```

rmf_df_reset['Cluster'] = rmf_df_reset['customer_id'].apply(lambda x:
cluster_map[x])
rmf_df_reset['age_int'] = rmf_df_reset['customer_id'].apply(lambda x:
age_map[x])
rmf_df_reset['pin_code_int'] = rmf_df_reset['customer_id'].apply(lambda x:
pin_map[x])
final_data=rmf_df_reset[['transaction_dt','customer_id','week_number','amount'
,'total_sum','frequency','recency','monetary','Cluster','age_int','pin_code_in
t','unit_price','Log_Unit_Price']]
comparison=(final_data['frequency']>1)
print('Avg Frequency:',round(final_data[comparison].frequency.mean()))
print('Avg Week Number:',round(final_data[comparison].week_number.mean()))
round(final_data[comparison].week_number.value_counts())

final_data.to_pickle('test_set.pkl')
#month leading up to the final weeks
comparison=((final_data['week_number']<=15) & (final_data['frequency']>1))
X_set=final_data[comparison]
#last weeks of data and has had at least 1 previous transaction
comparison=((final_data['week_number']==16) | (final_data['week_number']==17)
& (final_data['frequency']>1))
y_set=final_data[comparison]
y_id_list=sorted(list(y_set.customer_id.unique()))#unique customer ids from
final week
X_set=X_set.loc[X_set['customer_id'].isin(y_id_list)]#have to be represented
in final week
list_of_indexes = [np.argmin(g['transaction_dt']) for l, g in
X_set.groupby('customer_id')]
X_set=X_set.ix[list_of_indexes]#keep only latest transaction from previous
month purchases
X_set=X_set.sort_values(by=['customer_id','transaction_dt'])#sort by id
X_set=X_set.reset_index(drop=True)
print('number of unique customer
ids:',len(set(list(X_set.customer_id))),'\nlength of the dataset:',len(X_set))
X_set.head()

X_set_list=sorted(list(X_set.customer_id.unique()))#get list of unique ids
from those represented in the last month
y_set=y_set.loc[y_set['customer_id'].isin(X_set_list)]#return only individuals
who made a purchase in the prior month
y_set=y_set.sort_values(by=['customer_id','transaction_dt'])
list_of_indexes = [np.argmin(g['transaction_dt']) for l, g in
y_set.groupby('customer_id')]
y_set=y_set.ix[list_of_indexes]
y_set=y_set.reset_index(drop=True)
print('number of unique customer
ids:',len(set(list(y_set.customer_id))),'\nlength of the dataset:',len(y_set))
y_set.head()

```

```

X_set_df=X_set[['customer_id', 'week_number', 'amount', 'total_sum',
'Cluster', 'age_int', 'pin_code_int', 'unit_price', 'Log_Unit_Price', 'frequency',
'recency', 'monetary']]#reduce Xset
y_set_df=y_set[['customer_id', 'frequency', 'recency', 'monetary']]#reduce yset
X_set_df.to_pickle('X_set_df.pkl')
y_set_df.to_pickle('y_set_df.pkl')
with open('X_set_df.pkl', 'rb') as fp:
    X_set_df = pickle.load(fp)
with open('y_set_df.pkl', 'rb') as fp:
    y_set_df = pickle.load(fp)
X_set_df.head(3)

y_set_df.head(3)

# Full Model
def return_pair(df,indexer=0):
    data=df.iloc[indexer]
    x=[]
    part1=list(map(int,list(data.iloc[0])))
    part2=list(data.iloc[1:])
    x.extend(part1)
    x.extend(part2)
    x=np.array(x)
    return x
X_set=np.array([return_pair(X_set_df,i) for i in range(0,len(X_set_df))])
X_set[0]

y_set=np.array([return_pair(y_set_df,i) for i in range(0,len(y_set_df))])
y_set[0]

n_orig_size=len(X_set)
n_train_percent = round(n_orig_size*.2)
print('Original size:',n_orig_size,'\n20 Percent of size:',n_train_percent)

train_X=X_set[n_train_percent:,:]
train_y=y_set[n_train_percent:,:]
test_X=X_set[:n_train_percent, :]
test_y=y_set[:n_train_percent, :]
print(train_X.shape,test_X.shape)
print(train_y.shape,test_y.shape)

# Normalize the X,y sets
# Both sets have values on different scales.
scaler_y = MinMaxScaler(feature_range=(0, 1))
train_y = scaler_y.fit_transform(train_y)
test_y = scaler_y.fit_transform(test_y)
scaler_x = MinMaxScaler(feature_range=(0, 1))

```



```

rescaled_x_train = scaler_x.fit_transform(train_X)
rescaled_x_test = scaler_x.fit_transform(test_X)
# reshape input to be 3D [samples, timesteps, features]
train_X = rescaled_x_train.reshape((rescaled_x_train.shape[0], 1,
rescaled_x_train.shape[1]))
test_X = rescaled_x_test.reshape((rescaled_x_test.shape[0], 1,
rescaled_x_test.shape[1]))
print('train_X',train_X.shape,'\ntrain_y', train_y.shape,'\ntest_X',
test_X.shape,'\ntest_y', test_y.shape)

epochs = 1000
learning_rate = 0.001
model = Sequential()
model.add(SimpleRNN(19,activation='relu',input_shape=(train_X.shape[1],
train_X.shape[2])))
model.add(Dense(250,kernel_initializer='random_uniform',
activity_regularizer=regularizers.l1(0.0001)))
model.add(Dense(11))
model.compile(loss='mse', optimizer='adam', metrics=['mae', 'acc'])
model.summary()

history = model.fit(train_X, train_y, epochs=epochs, batch_size=120,
validation_data=(test_X, test_y), validation_split=0.3, verbose=1,
shuffle=True)
file_out=captured.stdout
print(file_out[-20000:])

pyplot.plot(history.history['loss'], label='train')
pyplot.plot(history.history['val_loss'], label='test')
pyplot.legend()
pyplot.show()

yhat = model.predict(test_X)
yhat.shape

test_X = test_X.reshape(test_X.shape[0], test_X.shape[2])
test_X.shape

inv_yhat = concatenate((scaler_x.inverse_transform(test_X),
scaler_y.inverse_transform(yhat)),axis=1)
inv_y =
concatenate((scaler_x.inverse_transform(test_X),scaler_y.inverse_transform(test_y)), axis=1)
rmse = sqrt(mean_squared_error(inv_y[0], inv_yhat[0]))
print('Test RMSE (Prediction): %.3f' % rmse)

indexer=10

```

```

transformed_id=[''.join(list(map(str,list(map(int,list(scaler_x.inverse_transform(test_X)[indexer][:-11])))))))]
week_number=[scaler_x.inverse_transform(test_X)[indexer][-11]]
amount=[scaler_x.inverse_transform(test_X)[indexer][-10]]
total_sum=[scaler_x.inverse_transform(test_X)[indexer][-9]]

cluster=[scaler_x.inverse_transform(test_X)[indexer][-8]]
age=[scaler_x.inverse_transform(test_X)[indexer][-7]]
pin=[scaler_x.inverse_transform(test_X)[indexer][-6]]
unit_price=[scaler_x.inverse_transform(test_X)[indexer][-5]]
log_unit_price=[scaler_x.inverse_transform(test_X)[indexer][-4]]
prediction=list(scaler_y.inverse_transform(yhat)[:,-3:][indexer])

input_set=pd.DataFrame((transformed_id+week_number+amount+total_sum+cluster+pin+age+unit_price+log_unit_price)).T.copy()
input_set.columns=['ID','Week_number','Amount','Total_sum','Cluster','Age_group','Pin_code','Unit_price','Log_unit_price']
input_set

# Training R,F,M Values
outputID=input_set['ID'][0]
comparison=X_set_df['customer_id']==outputID
X_set_df[comparison][['frequency','recency','monetary']]

# Prediction
output_set=pd.DataFrame((prediction)).T.copy()
output_set.columns=['frequency','recency','monetary']
output_set

# True value
comparison=y_set_df['customer_id']==outputID
y_set_df[comparison][['frequency','recency','monetary']]

# Reduced Model
X_set_df=X_set_df[['customer_id','frequency','recency','monetary']]#reduce Xset
y_set_df=y_set_df[['customer_id','frequency','recency','monetary']]#reduce yset
X_set_df.head(3)

y_set_df.head(3)

X_set=np.array([return_pair(X_set_df,i) for i in range(0,len(X_set_df))])
X_set[0]

y_set=np.array([return_pair(y_set_df,i) for i in range(0,len(y_set_df))])
y_set[0]

```

```

n_orig_size=len(X_set)
n_train_percent = round(n_orig_size*.2)
print('Original size:',n_orig_size,'\n20 Percent of size:',n_train_percent)

train_X=X_set[n_train_percent:,:]
train_y=y_set[n_train_percent:,:]
test_X=X_set[:n_train_percent, :]
test_y=y_set[:n_train_percent, :]

scaler_y = MinMaxScaler(feature_range=(0, 1))
train_y = scaler_y.fit_transform(train_y)
test_y = scaler_y.fit_transform(test_y)
scaler_x = MinMaxScaler(feature_range=(0, 1))
rescaled_x_train = scaler_x.fit_transform(train_X)
rescaled_x_test = scaler_x.fit_transform(test_X)
train_X = rescaled_x_train.reshape((rescaled_x_train.shape[0], 1,
rescaled_x_train.shape[1]))
test_X = rescaled_x_test.reshape((rescaled_x_test.shape[0], 1,
rescaled_x_test.shape[1]))
print('train_X',train_X.shape,'\ntrain_y', train_y.shape,'\ntest_X',
test_X.shape,'\ntest_y', test_y.shape)

epochs = 1000
learning_rate = 0.001
model = Sequential()
model.add(SimpleRNN(11,activation='relu',input_shape=(train_X.shape[1],
train_X.shape[2])))
model.add(Dense(250,kernel_initializer='random_uniform',
activity_regularizer=regularizers.l1(0.0001)))
model.add(Dense(11))
model.compile(loss='mse', optimizer='adam', metrics=['mae', 'acc'])
model.summary()

history = model.fit(train_X, train_y, epochs=epochs, batch_size=120,
validation_data=(test_X, test_y), validation_split=0.3, verbose=1,
shuffle=True)
file_out=captured.stdout
print(file_out[-20000:])

# plot history
pyplot.plot(history.history['loss'], label='train')
pyplot.plot(history.history['val_loss'], label='test')
pyplot.legend()
pyplot.show()

yhat = model.predict(test_X)
test_X = test_X.reshape(test_X.shape[0], test_X.shape[2])

```

```

inv_yhat = concatenate((scaler_x.inverse_transform(test_X),
scaler_y.inverse_transform(yhat)),axis=1)
inv_y =
concatenate((scaler_x.inverse_transform(test_X),scaler_y.inverse_transform(test_y)), axis=1)
# calculate RMSE
rmse = sqrt(mean_squared_error(inv_y[0], inv_yhat[0]))
print('Test RMSE (Prediction): %.3f' % rmse)

indexer=10
transformed_id=[''.join(list(map(str,list(map(int,list(scaler_x.inverse_transform(test_X)[indexer][: -3])))))))]
prediction=list(scaler_y.inverse_transform(yhat)[:, -3:][indexer])

# Prediction
output=pd.DataFrame((transformed_id+prediction)).T.copy()
output.columns=['ID', 'Recency', 'Frequency', 'Monetary']
output

# Training
outputID=output['ID'][0]
comparison=X_set_df['customer_id']==outputID
X_set_df[comparison]

# True value
comparison=y_set_df['customer_id']==outputID
y_set_df[comparison]

```