

```
1 #%% md
2 # **ST1/ST1G Assignment 9 (Capstone Programming
Project)**
3 #%% md
4 ### This Project is based on E-commerce Price
Prediction available from Kaggle repository (https://www.kaggle.com/datasets/bhuwanesh340/ecommerce-price-prediction/data).
5
6
7 * It contains the details of 2452 products.
8 * Project task is to create a machine learning
model which can predict the price of products based
on its attributes.
9 * For solving this problem, I will approach the task
, with a step by step approach to create a data
analysis and prediction model based on (machine
learning/AI algorithms, regression algorithm for
example) available from different Python packages,
modules and classes.
10
11 #%% md
12 ### Step 1: Reading the data Reading the data with
python
13 #%%
14 # Supressing the warning messages
15 import warnings
16 warnings.filterwarnings('ignore')
17 #%%
18 # Reading the dataset
19 import pandas as pd
20 df=pd.read_csv('D:\PricePrediction\Train.csv',
encoding='latin')
21 print('Shape before deleting duplicate values:', df.
shape)
22
23 # Removing duplicate rows if any
24 df=df.drop_duplicates()
25 print('Shape After deleting duplicate values:', df.
shape)
26
```

```
27 # Printing sample data
28 # Start observing the Quantitative/Categorical/
  Qualitative variables
29 df.head(10)
30
31 #%% md
32 ## Key observations from Step 1 about Data
  Description
33
34
35 * This file contains 2452 product details from
  dataset.
36 * There are 8 attributes and they are outlined
  below.
37 * Product - Unique ID for each product
38 * Product_Brand
39 * Item_Category
40 * Subcategory_1
41 * Subcategory_2
42 * Item_Rating - customer ratings of products
43 * Date - The date the item was sold
44 * Selling_Price
45
46
47
48 #%% md
49 # Step 2 : Problem Statement Definition
50 * Creating a prediction model to predict the price
  of a product.
51 * Target Variable: SellingPrice
52 * Predictors/Features: Product/Product_Brand/
  Item_Category/Subcategory_1/Subcategory_2/Item_Rating
  /Date
53 #%% md
54 ## Step 3: Choosing the appropriate ML/AI Algorithm
  for Data Analysis.
55 * Based on the problem statement to create a
  supervised ML Regression model, as the target
  variable is **Continuous**.
56
57 #%% md
```

```
58 # Step 4: Looking at the class distribution (Target  
variable distribution to check if the data is  
balanced or skewed.  
59 #%%  
60 %matplotlib inline  
61 # Creating histogram as the Target variable is  
Continuous  
62  
63 df['Selling_Price'].hist(bins=20)  
64  
65 #%% md  
66 ## Observations from Step 4  
67 * The data distribution of the target variable is too  
skewed.  
68 * After removing the data with large deviation there  
still have sufficient number of rows for each type of  
values to learn from.  
69 #%% md  
70 ## Step 5: Basic Exploratory Data Analysis  
71  
72 #%%  
73 # Looking at sample rows in the data  
74 df.head()  
75 #%%  
76 # Looking at sample rows in the data  
77 df.tail()  
78 #%%  
79 # Summarise information of data  
80 df.info()  
81 #%%  
82 # Looking at the descriptive statistics of the data  
83 df.describe(include='all')  
84 #%%  
85 #Convert Date into categorical variable  
86 df['Date'] = pd.to_datetime(df['Date'],format='%d/%m  
/%Y')  
87  
88 # import the holidays list  
89 import holidays  
90  
91 df['Date'] = pd.to_datetime(df['Date'])
```

```

92
93 # Create a holidays object for Australia
94 aus_holidays = holidays.Australia(years=[df['Date'].dt.year.min(), df['Date'].dt.year.max()])
95
96 # Check each date in the Date column to see if it is
97 # a holiday
98 df['Is_Holiday'] = df['Date'].apply(lambda x: 1 if x
99 # in aus_holidays else 0)
100 print(df.head())
101 #%%
102 # Finding unique values for each column
103 # Typically if the number of unique values are < 20
104 # then the variable is likely to be a category
105 # otherwise continuous
106 df.nunique()
107 #%% md
108 ## Observations from Step 5 - Basic Exploratory Data
109 # Analysis
110 * Product - Just products' unique ID,
111 # waiting for delete.
112 * Product_Brand - Categorical. Selected.
113 * Item_Category - Categorical. Selected.
114 * Subcategory_1 - Categorical. Selected.
115 * Subcategory_2 - Categorical. Selected.
116 * Item_Rating - Continuous. Selected.
117 * Date - Transfer to Is_Holiday
118 * Selling_Price - Continuous. Selected. This is
119 # the Target variable.
120 * Is_Holiday - Categorical. Selected.
121 #%%
122 ## Step 7: Removing Unwanted columns
123 * If each product has a unique identifier, this
124 # column will not provide any generalizable insight
125 # into the selling price across different products, as
126 # it does not encode any repeatable pattern that can
127 # be learned by the model.
128 * Remove the Product column as it is likely to lead
129 # to overfitting, where the model learns noise
130 # specific to individual data points rather than

```

```
118 general trends.  
119 * Converting dates to categorical variables for  
research and machine learning. Considering that  
sales discounts are offered on holidays, there is an  
impact on the sales price.  
120 * Replace Date column with Is_Holidays.  
121 #%%  
122 #Convert Date into categorical variable  
123 df['Date'] = pd.to_datetime(df['Date'],format='%d/%m  
/%Y')  
124  
125 # import the holidays list  
126 import holidays  
127  
128 df['Date'] = pd.to_datetime(df['Date'])  
129  
130 # Create a holidays object for Australia  
131 aus_holidays = holidays.Australia(years=[df['Date'].  
dt.year.min(), df['Date'].dt.year.max()])  
132  
133 # Check each date in the Date column to see if it is  
a holiday  
134 df['Is_Holiday'] = df['Date'].apply(lambda x: 1 if x  
in aus_holidays else 0)  
135 print(df.head())  
136  
137 # Remove 'Product' and 'Date'.  
138 df = df.drop('Product', axis=1)  
139 df = df.drop('Date', axis=1)  
140  
141 #%% md  
142 ## Step 8: Visual Exploratory Data Analysis  
143 * Based on the Basic Exploration Data Analysis in  
the previous step, we could spotted two categorical  
predictors in the data  
144 * Categorical Predictors:  
145  
146 * Product_Brand  
147 * Item_Category  
148 * Subcategory_1  
149 * Subcategory_2
```

```
150 * Is_Holiday
151
152 * Use bar charts to see how the data is distributed
    for these categorical columns.
153
154 #%%
155 # Plotting multiple bar charts at once for
    categorical variables
156 # Since there is no default function which can plot
    bar charts for multiple columns at once
157 # we are defining our own function for the same
158
159 import matplotlib.pyplot as plt
160 import seaborn as sns
161
162 # Item_Category
163 plt.figure(figsize=(12, 6))
164 sns.countplot(data=df, x='Item_Category', order=df['
    Item_Category'].value_counts().index[:50])
165 plt.xticks(rotation=90)
166 plt.title('Item Category Distribution')
167 plt.show()
168
169 # Subcategory_1
170 plt.figure(figsize=(12, 6))
171 sns.countplot(data=df, x='Subcategory_1', order=df['
    Subcategory_1'].value_counts().index[:20])
172 plt.xticks(rotation=90)
173 plt.title('Top 20 Subcategory 1 Distribution')
174 plt.show()
175
176 # Subcategory_2
177 plt.figure(figsize=(12, 6))
178 sns.countplot(data=df, x='Subcategory_2', order=df['
    Subcategory_2'].value_counts().index[:20])
179 plt.xticks(rotation=90)
180 plt.title('Top 20 Subcategory 2 Distribution')
181 plt.show()
182
183 #Is_Holiday
184 plt.figure(figsize=(12, 6))
```

```
185 sns.countplot(data=df, x='Is_Holiday', order=df['  
    Is_Holiday'].value_counts().index[:20])  
186 plt.xticks(rotation=90)  
187 plt.title('Is_Holiday Distribution')  
188 plt.show()  
189  
190 plt.figure(figsize=(12, 6))  
191 sns.countplot(data=df, x='Product_Brand', order=df['  
    Product_Brand'].value_counts().index[:50])  
192 plt.xticks(rotation=90)  
193 plt.title('ProductBrand Distribution')  
194 plt.show()  
195 #%% md  
196 ## Observations from Step 8 - Visual Exploratory  
    Data Analysis  
197 * Although the distribution of Product_Brand and  
    Item_Category bars charts is skewed, but they were  
    retained given the need for the predictive model to  
    include a variety of different categories of goods.  
198  
199 * Selected Categorical Variables (**Product_Brand,  
    Item_Category,Subcategory_1,Subcategory_2,Is_Holiday  
    **): Both the categorical variables are selected for  
    further analysis.  
200 #%% md  
201 ## Step 9: Now Visualize distribution of all the  
    Continuous Predictor variables in the data using  
    histograms  
202 * Based on the Basic Exploratory Data Analysis,  
    there are eleven continuous predictor variables '  
    CRIM', 'ZN', 'INDUS', 'NOX', 'RM', 'AGE', 'DIS', '  
    TAX', 'PTRATIO', 'B', and 'LSTAT'.  
203 #%%  
204 # Plotting histograms of multiple columns together  
205 df.hist(['Item_Rating', 'Selling_Price'], figsize=(  
    18,10))  
206 #%% md  
207 ## Observations from Step 9  
208 * The distribution of Selling_Price is too skewed,  
    with one piece of data dominating and under 6000  
    needing to be dealt with.
```

```
209
210 * Selected Continuous Variables:
211
212 * Selling_Price : Selected. Outliers seen beyond
   6000, need to treat them.
213 * Item_Rating : Selected. The distribution is good.
214 #%% md
215 ## Step 10: Outlier Analysis
216 * Data with too high a selling price but very low
   frequency can make the model very inaccurate, so
   these outliers are removed
217
218 #%%
219 #Removing outliers for 'Selling_Price'
220
221 df = df[df['Selling_Price'] <= 6000]
222
223 #%% md
224 Observation: Above result hence, deleting any value
   above 6000 with it.
225 #%% md
226 ## Step 11: Visualising Data Distribution after
   outlier removal
227
228 #%%
229 df.hist(['Selling_Price'], figsize=(18,5))
230 #%% md
231 ## Observation from Step 11
232 * The distribution has improved after the outlier
   treatment.
233 * There is still a tail but it is thick, that means
   there are many values in that range, hence, it is
   acceptable.
234 #%% md
235 ## Step 12: Missing Values Analysis
236 #%%
237 # Finding how many missing values are there for each
   column
238 df.isnull().sum()
239 #%% md
240 ## Observations from Step 12: Missing Value Analysis
```

```
241 * No missing values in this data!
242 * So no removal of any data samples(rows) is needed.
243
244 #%% md
245 ## Step 13: Feature Selection (Attribute Selection)
246
247 * **For this dataset, the Target variable is
    Continuous, hence following two scenarios will need
    attention**
248
249 * Continuous Target Variable Vs Continuous Predictor
250 * Continuous Target Variable Vs Categorical
    Predictor
251
252 #%% md
253 ## Relationship exploration: Continuous Vs
    Continuous -- Scatter Charts
254 * Using scatter plot and measure the strength of
    relation using a metric called pearson's correlation
    value.
255 #%%
256 ContinuousCols=['Item_Rating']
257
258 # Plotting scatter chart for each predictor vs the
    target variable
259 for predictor in ContinuousCols:
260     df.plot.scatter(x=predictor, y='Selling_Price',
        figsize=(10,5), title=predictor+" VS "+'
        Selling_Price')
261 #%% md
262 ## Scatter charts interpretation
263 * **No Trend**: I cannot see any clear increasing or
    decreasing trend. This means there is no
    correlation between the variables. Hence that
    predictor/feature may not be the best one for ML
    model building.
264
265 * I'll confirm this by looking at the correlation
    value in the next step.
266 #%% md
267 ## Step 14: Statistical Feature Selection (
```

```
267 Continuous Vs Continuous) using Correlation value
268 #%%
269 # Calculating correlation matrix
270 ContinuousCols=['Selling_Price', 'Item_Rating']
271
272 # Creating the correlation matrix
273 CorrelationData=df[ContinuousCols].corr()
274 CorrelationData
275 #%%
276 # Filtering only those columns where absolute
correlation > 0.5 with Target Variable
277 # reduce the 0.5 threshold if no variable is
selected
278 CorrelationData['Selling_Price'][abs(CorrelationData
['Selling_Price']) > 0.5 ]
279 #%% md
280 ## Observations from Step 14
281 * From the results we can see that there is no
correlation between Item_Rating and the target
variable. No Continuous columns selected.
282 #%% md
283 # Step 15: Relationship exploration: Categorical Vs
Continuous -- Box Plots
284 * Measure the strength of relation using Anova test.
285 #%%
286 # Box plots for continuous Target Variable "
Selling_Price" and Categorical predictors
287 CategoricalColsList=['Product_Brand','Item_Category'
, 'Subcategory_1', 'Subcategory_2', 'Is_Holiday']
288
289 import matplotlib.pyplot as plt
290 fig, PlotCanvas=plt.subplots(nrows=1, ncols=len(
CategoricalColsList), figsize=(18,5))
291
292 # Creating box plots for each continuous predictor
against the Target Variable "MEDV"
293 for PredictorCol , i in zip(CategoricalColsList,
range(len(CategoricalColsList))):
294     df.boxplot(column='Selling_Price', by=
PredictorCol, figsize=(5,5), vert=True, ax=
PlotCanvas[i])
```

```
295 %% md
296 ##Observations from Step 15: Box-Plots
   interpretation
297 * The distribution is different for each category( the boxes are not in same line). It hints that these variables might be correlated with Selling_Price.
298
299 * For this data, all the categorical predictors looks correlated with the Target variable.
300
301 We confirm this by looking at the results of ANOVA test below
302 %% md
303 ## Step 16: Statistical Feature Selection ( Categorical Vs Continuous) using ANOVA test
304 %%
305 # Defining a function to find the statistical relationship with all the categorical variables
306 def FunctionAnova(inpData, TargetVariable,
   CategoricalPredictorList):
307     from scipy.stats import f_oneway
308
309     # Creating an empty list of final selected predictors
310     SelectedPredictors=[]
311
312     print('##### ANOVA Results ##### \n')
313     for predictor in CategoricalPredictorList:
314         CategoryGroupLists=inpData.groupby(predictor)
315         [TargetVariable].apply(list)
316         AnovaResults = f_oneway(*CategoryGroupLists)
317
318         # If the ANOVA P-Value is <0.05, that means we reject H0
319         if (AnovaResults[1] < 0.05):
320             print(predictor, 'is correlated with',
321                   TargetVariable, '| P-Value:', AnovaResults[1])
322             SelectedPredictors.append(predictor)
323         else:
324             print(predictor, 'is NOT correlated with',
325                   TargetVariable, '| P-Value:', AnovaResults[1])
```

```
323
324     return(SelectedPredictors)
325 #%%
326 #Calling the function to check which categorical
variables are correlated with target
327 CategoricalPredictorList=['Product_Brand', '
Subcategory_1', 'Subcategory_2', 'Item_Category', '
Is_Holiday']
328 FunctionAnova(inpData=df,
329                 TargetVariable='Selling_Price',
330                 CategoricalPredictorList=
CategoricalPredictorList)
331 #%% md
332 ##Observations from Step 16
333 * Final selected Categorical columns:
334
335 'Product_Brand', 'Subcategory_1', 'Subcategory_2', '
Item_Category', 'Is_Holiday'
336
337 #%% md
338 ## Selecting final Predictors/Features for building
Machine Learning/AI model.
339 * Select the final features/predictors/columns for
machine learning model building as:
340 * **'Product_Brand', 'Subcategory_1', 'Subcategory_2
', 'Item_Category', 'Is_Holiday**'
341
342 #%%
343 SelectedColumns=['Product_Brand', 'Subcategory_1', '
Subcategory_2', 'Item_Category', 'Is_Holiday']
344
345 # Selecting final columns
346 DataForML=df[SelectedColumns]
347 DataForML.head()
348
349 #%%
350 # Saving this final data subset for reference during
deployment
351 DataForML.to_pickle('DataForML.pkl')
352 #%% md
353 ## Step 17: Data Pre-processing for Machine Learning
```

```
353 Model Building or Model Development
354 Considering the large number of categorical
variables in the predicted values, it is necessary
to encode the data using one-hot coding for easy use
in machine learning.
355
356 #%% md
357 ## Converting the nominal variable to numeric using
get_dummies()
358 #%%
359 # Treating all the nominal variables at once using
dummy variables
360 DataForML_Numeric=pd.get_dummies(DataForML)
361
362 # Adding Target Variable to the data
363 DataForML_Numeric['Selling_Price']=df['Selling_Price']
364
365 # Printing sample rows
366 DataForML_Numeric.head()
367 #%% md
368 ## Step 18: Machine Learning Model Development:
369 #%%
370 # Printing all the column names for reference
371 DataForML_Numeric.columns
372 #%%
373 #Separate Target Variable and Predictor Variables
374 import numpy as np
375 import pandas as pd
376 from sklearn.preprocessing import MinMaxScaler
377 from sklearn.model_selection import train_test_split
378 from sklearn.linear_model import LinearRegression
379
380 TargetVariable = 'Selling_Price'
381 Predictors = [column for column in DataForML_Numeric
.columns if column != TargetVariable]
382
383
384 # Extraction of features and target variables
385 X=DataForML_Numeric[Predictors].values
386 y=DataForML_Numeric[TargetVariable].values
```

```
387
388 # Split the data into training(70%) and testing set(30%).
389 from sklearn.model_selection import train_test_split
390 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=428)
391
392 #%% md
393 ## Step 19: Standardization/Normalization of data
394 * I'm using KNN, so this step is necessary.
395 #%%
396 ### Standardization of data ####
397 from sklearn.preprocessing import StandardScaler,
    MinMaxScaler
398
399 # Using MinMAX normalization is better
400 PredictorScaler=MinMaxScaler()
401
402 # Storing the fit object for later reference
403 PredictorScalerFit=PredictorScaler.fit(X)
404
405 # Generating the standardized values of X
406 X=PredictorScalerFit.transform(X)
407
408 # Split the data into training and testing set
409 from sklearn.model_selection import train_test_split
410 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
411 #%%
412 # Sanity check for the sampled data
413 print(X_train.shape)
414 print(y_train.shape)
415 print(X_test.shape)
416 print(y_test.shape)
417 #%% md
418 ## Step 20: Multiple Linear Regression Algorithm For ML/AI model building
419 #%% md
420 # Decision Tree Regressor
421 #%%
422 # Decision Trees (Multiple if-else statements!)
```

```
423 from sklearn.tree import DecisionTreeRegressor
424 RegModel = DecisionTreeRegressor(max_depth=5,
425     criterion='friedman_mse')
426 # Good Range of Max_depth = 2 to 20
427
428 # Printing all the parameters of Decision Tree
429 print(RegModel)
430
431 # Creating the model on Training Data
432 DT=RegModel.fit(X_train,y_train)
433 prediction=DT.predict(X_test)
434
435 # Measuring Goodness of fit in Training data
436 print('R2 Value:',metrics.r2_score(y_train, DT.
437     predict(X_train)))
438
439 # Plotting the feature importance for Top 10 most
440 # important columns
441 %matplotlib inline
442 feature_importances = pd.Series(DT.
443     feature_importances_, index=Predictors)
444 feature_importances.nlargest(10).plot(kind='barh')
445
446 ###### Model Validation and Accuracy
447 Calculations #####
448
449 print('\n##### Model Validation and Accuracy
500 Calculations #####')
501
502 # Printing some sample values of prediction
503 TestingDataResults=pd.DataFrame(data=X_test, columns
504     =Predictors)
505 TestingDataResults[TargetVariable]=y_test
506 TestingDataResults[['Predicted'+TargetVariable]]=np.
507     round(prediction)
508
509 # Printing sample prediction values
510 print(TestingDataResults.head())
511
512 # Calculating the error for each row
513 TestingDataResults['APE']=100 * ((abs(
```

```
456 TestingDataResults['Selling_Price']-
    TestingDataResults['PredictedSelling_Price'])/
    TestingDataResults['Selling_Price'])
457
458 MAPE=np.mean(TestingDataResults['APE'])
459 MedianMAPE=np.median(TestingDataResults['APE'])
460
461 Accuracy =100 - MAPE
462 MedianAccuracy=100- MedianMAPE
463 print('Mean Accuracy on test data:', Accuracy) # Can
    be negative sometimes due to outlier
464 print('Median Accuracy on test data:',
    MedianAccuracy)
465
466 # Defining a custom function to calculate accuracy
467 # Make sure there are no zeros in the Target
    variable
468 def Accuracy_Score(orig,pred):
469     MAPE = np.mean(100 * (np.abs(orig-pred)/orig))
470     #print('#'*70, 'Accuracy:', 100-MAPE)
471     return(100-MAPE)
472
473 # Custom Scoring MAPE calculation
474 from sklearn.metrics import make_scorer
475 custom_Scoring=make_scorer(Accuracy_Score,
    greater_is_better=True)
476
477 # Importing cross validation function from sklearn
478 from sklearn.model_selection import cross_val_score
479
480 # Running 10-Fold Cross validation on a given
    algorithm
481 # Passing full data X and y because the K-fold will
    split the data and automatically choose train/test
482 Accuracy_Values=cross_val_score(RegModel, X , y, cv=
    10, scoring=custom_Scoring)
483 print('\nAccuracy values for 10-fold Cross
    Validation:\n',Accuracy_Values)
484 print('\nFinal Average Accuracy of the model:', 
    round(Accuracy_Values.mean(),2))
485 %% md
```

```
486 # Plotting/Visualising the Decision Tree
487 #%%
488
489 # Load libraries
490 from IPython.display import Image
491 from sklearn import tree
492 import pydotplus
493
494 # Create DOT data
495 dot_data = tree.export_graphviz(RegModel, out_file=
    None,
        feature_names=
    Predictors, class_names=TargetVariable)
497
498 # printing the rules
499 print(dot_data)
500
501 # Draw graph
502 graph = pydotplus.graph_from_dot_data(dot_data)
503
504 # Show graph
505 Image(graph.create_png(), width=2000,height=2000)
506 # Double click on the graph to zoom in
507 #%% md
508 # Random Forest Regressor
509 #%%
510 # Random Forest (Bagging of multiple Decision Trees)
511 from sklearn.ensemble import RandomForestRegressor
512 RegModel = RandomForestRegressor(max_depth=4,
    n_estimators=400,criterion='friedman_mse')
513 # Good range for max_depth: 2-10 and n_estimators:
    100-1000
514
515 # Printing all the parameters of Random Forest
516 print(RegModel)
517
518 # Creating the model on Training Data
519 RF=RegModel.fit(X_train,y_train)
520 prediction=RF.predict(X_test)
521
522 from sklearn import metrics
```

```
523 # Measuring Goodness of fit in Training data
524 print('R2 Value:',metrics.r2_score(y_train, RF.
    predict(X_train)))
525
526 # Plotting the feature importance for Top 10 most
    important columns
527 %matplotlib inline
528 feature_importances = pd.Series(RF.
    feature_importances_, index=Predictors)
529 feature_importances.nlargest(10).plot(kind='barh')
530
531 ##### Model Validation and Accuracy
##### Calculations #####
532 print('\n##### Model Validation and Accuracy
    Calculations #####')
533
534 # Printing some sample values of prediction
535 TestingDataResults=pd.DataFrame(data=X_test, columns
    =Predictors)
536 TestingDataResults[TargetVariable]=y_test
537 TestingDataResults[('Predicted'+TargetVariable)]=np.
    round(prediction)
538
539 # Printing sample prediction values
540 print(TestingDataResults.head())
541
542 # Calculating the error for each row
543 TestingDataResults['APE']=100 * ((abs(
544     TestingDataResults['Selling_Price']-
        TestingDataResults['PredictedSelling_Price'])/
    TestingDataResults['Selling_Price']))
545
546 MAPE=np.mean(TestingDataResults['APE'])
547 MedianMAPE=np.median(TestingDataResults['APE'])
548
549 Accuracy =100 - MAPE
550 MedianAccuracy=100- MedianMAPE
551 print('Mean Accuracy on test data:', Accuracy) # Can
        be negative sometimes due to outlier
552 print('Median Accuracy on test data:',
    MedianAccuracy)
```

```
553
554
555 # Defining a custom function to calculate accuracy
556 # Make sure there are no zeros in the Target
      variable if you are using MAPE
557 def Accuracy_Score(orig,pred):
558     MAPE = np.mean(100 * (np.abs(orig-pred)/orig))
559     #print('#'*70, 'Accuracy:', 100-MAPE)
560     return(100-MAPE)
561
562 # Custom Scoring MAPE calculation
563 from sklearn.metrics import make_scorer
564 custom_Scoring=make_scorer(Accuracy_Score,
      greater_is_better=True)
565
566 # Importing cross validation function from sklearn
567 from sklearn.model_selection import cross_val_score
568
569 # Running 10-Fold Cross validation on a given
      algorithm
570 # Passing full data X and y because the K-fold will
      split the data and automatically choose train/test
571 Accuracy_Values=cross_val_score(RegModel, X , y, cv=
      10, scoring=custom_Scoring)
572 print('\nAccuracy values for 10-fold Cross
      Validation:\n',Accuracy_Values)
573 print('\nFinal Average Accuracy of the model:',
      round(Accuracy_Values.mean(),2))
574 #%% md
575
576 ## Step 21: AdaBoost Algorithm For ML/AI model
      building
577 #%%
578 # Adaboost (Boosting of multiple Decision Trees)
579 from sklearn.ensemble import AdaBoostRegressor
580 from sklearn.tree import DecisionTreeRegressor
581
582 # Choosing Decision Tree with 6 level as the weak
      learner
583 DTR=DecisionTreeRegressor(max_depth=3)
584 RegModel = AdaBoostRegressor(n_estimators=100 ,
```

```
584 learning_rate=0.04)
585
586 # Printing all the parameters of Adaboost
587 print(RegModel)
588
589 # Creating the model on Training Data
590 AB=RegModel.fit(X_train,y_train)
591 prediction=AB.predict(X_test)
592
593 from sklearn import metrics
594 # Measuring Goodness of fit in Training data
595 print('R2 Value:',metrics.r2_score(y_train, AB.
    predict(X_train)))
596
597 # Plotting the feature importance for Top 10 most
    important columns
598 %matplotlib inline
599 feature_importances = pd.Series(AB.
    feature_importances_, index=Predictors)
600 feature_importances.nlargest(10).plot(kind='barh')
601
602 ##### Model Validation and Accuracy
    Calculations #####
603 print('\n##### Model Validation and Accuracy
    Calculations #####')
604
605 # Printing some sample values of prediction
606 TestingDataResults=pd.DataFrame(data=X_test, columns
    =Predictors)
607 TestingDataResults[TargetVariable]=y_test
608 TestingDataResults[('Predicted'+TargetVariable)]=np.
    round(prediction)
609
610 # Printing sample prediction values
611 print(TestingDataResults.head())
612
613 # Calculating the error for each row
614 TestingDataResults['APE']=100 * ((abs(
615     TestingDataResults['Selling_Price']-
        TestingDataResults['PredictedSelling_Price']))/
    TestingDataResults['Selling_Price'])
```

```
616
617 MAPE=np.mean(TestingDataResults['APE'])
618 MedianMAPE=np.median(TestingDataResults['APE'])
619
620 Accuracy =100 - MAPE
621 MedianAccuracy=100- MedianMAPE
622 print('Mean Accuracy on test data:', Accuracy) # Can
   be negative sometimes due to outlier
623 print('Median Accuracy on test data:',
   MedianAccuracy)
624
625
626 # Defining a custom function to calculate accuracy
627 # Make sure there are no zeros in the Target
   variable if you are using MAPE
628 def Accuracy_Score(orig,pred):
629     MAPE = np.mean(100 * (np.abs(orig-pred)/orig))
630     #print('#'*70,'Accuracy:', 100-MAPE)
631     return(100-MAPE)
632
633 # Custom Scoring MAPE calculation
634 from sklearn.metrics import make_scorer
635 custom_Scoring=make_scorer(Accuracy_Score,
   greater_is_better=True)
636
637 # Importing cross validation function from sklearn
638 from sklearn.model_selection import cross_val_score
639
640 # Running 10-Fold Cross validation on a given
   algorithm
641 # Passing full data X and y because the K-fold will
   split the data and automatically choose train/test
642 Accuracy_Values=cross_val_score(RegModel, X , y, cv=
   10, scoring=custom_Scoring)
643 print('\nAccuracy values for 10-fold Cross
   Validation:\n',Accuracy_Values)
644 print('\nFinal Average Accuracy of the model:', round(Accuracy_Values.mean(),2))
645 #%% md
646 # XGBoost Regressor
647 #%%
```

```
648 # Xtreme Gradient Boosting (XGBoost)
649 from xgboost import XGBRegressor
650 RegModel=XGBRegressor(max_depth=2,
651                         learning_rate=0.1,
652                         n_estimators=1000,
653                         objective='reg:linear',
654                         booster='gbtree')
655
656 # Printing all the parameters of XGBoost
657 print(RegModel)
658
659 # Creating the model on Training Data
660 XGB=RegModel.fit(X_train,y_train)
661 prediction=XGB.predict(X_test)
662
663 from sklearn import metrics
664 # Measuring Goodness of fit in Training data
665 print('R2 Value:',metrics.r2_score(y_train, XGB.
666     predict(X_train)))
667 # Plotting the feature importance for Top 10 most
668 # important columns
669 %matplotlib inline
670 feature_importances = pd.Series(XGB.
671     feature_importances_, index=Predictors)
672 feature_importances.nlargest(10).plot(kind='barh')
673 #####
674 print('\n##### Model Validation and Accuracy
675 Calculations #####')
676
677 # Printing some sample values of prediction
678 TestingDataResults=pd.DataFrame(data=X_test, columns
679     =Predictors)
680 TestingDataResults[TargetVariable]=y_test
681 TestingDataResults[('Predicted'+TargetVariable)]=np.
682     round(prediction)
683
684 # Printing sample prediction values
685 print(TestingDataResults.head())
686
687 # Printing sample prediction values
688 print(TestingDataResults.head())
689
690 # Printing sample prediction values
691 print(TestingDataResults.head())
692
693 # Printing sample prediction values
694 print(TestingDataResults.head())
695
696 # Printing sample prediction values
697 print(TestingDataResults.head())
698
699 # Printing sample prediction values
700 print(TestingDataResults.head())
701
702 # Printing sample prediction values
703 print(TestingDataResults.head())
704
705 # Printing sample prediction values
706 print(TestingDataResults.head())
707
708 # Printing sample prediction values
709 print(TestingDataResults.head())
710
711 # Printing sample prediction values
712 print(TestingDataResults.head())
713
714 # Printing sample prediction values
715 print(TestingDataResults.head())
716
717 # Printing sample prediction values
718 print(TestingDataResults.head())
719
720 # Printing sample prediction values
721 print(TestingDataResults.head())
722
723 # Printing sample prediction values
724 print(TestingDataResults.head())
725
726 # Printing sample prediction values
727 print(TestingDataResults.head())
728
729 # Printing sample prediction values
730 print(TestingDataResults.head())
731
732 # Printing sample prediction values
733 print(TestingDataResults.head())
734
735 # Printing sample prediction values
736 print(TestingDataResults.head())
737
738 # Printing sample prediction values
739 print(TestingDataResults.head())
740
741 # Printing sample prediction values
742 print(TestingDataResults.head())
743
744 # Printing sample prediction values
745 print(TestingDataResults.head())
746
747 # Printing sample prediction values
748 print(TestingDataResults.head())
749
750 # Printing sample prediction values
751 print(TestingDataResults.head())
752
753 # Printing sample prediction values
754 print(TestingDataResults.head())
755
756 # Printing sample prediction values
757 print(TestingDataResults.head())
758
759 # Printing sample prediction values
760 print(TestingDataResults.head())
761
762 # Printing sample prediction values
763 print(TestingDataResults.head())
764
765 # Printing sample prediction values
766 print(TestingDataResults.head())
767
768 # Printing sample prediction values
769 print(TestingDataResults.head())
770
771 # Printing sample prediction values
772 print(TestingDataResults.head())
773
774 # Printing sample prediction values
775 print(TestingDataResults.head())
776
777 # Printing sample prediction values
778 print(TestingDataResults.head())
779
780 # Printing sample prediction values
781 print(TestingDataResults.head())
782
783 # Printing sample prediction values
784 print(TestingDataResults.head())
785
786 # Printing sample prediction values
787 print(TestingDataResults.head())
788
789 # Printing sample prediction values
790 print(TestingDataResults.head())
791
792 # Printing sample prediction values
793 print(TestingDataResults.head())
794
795 # Printing sample prediction values
796 print(TestingDataResults.head())
797
798 # Printing sample prediction values
799 print(TestingDataResults.head())
800
801 # Printing sample prediction values
802 print(TestingDataResults.head())
803
804 # Printing sample prediction values
805 print(TestingDataResults.head())
806
807 # Printing sample prediction values
808 print(TestingDataResults.head())
809
810 # Printing sample prediction values
811 print(TestingDataResults.head())
812
813 # Printing sample prediction values
814 print(TestingDataResults.head())
815
816 # Printing sample prediction values
817 print(TestingDataResults.head())
818
819 # Printing sample prediction values
820 print(TestingDataResults.head())
821
822 # Printing sample prediction values
823 print(TestingDataResults.head())
824
825 # Printing sample prediction values
826 print(TestingDataResults.head())
827
828 # Printing sample prediction values
829 print(TestingDataResults.head())
830
831 # Printing sample prediction values
832 print(TestingDataResults.head())
833
834 # Printing sample prediction values
835 print(TestingDataResults.head())
836
837 # Printing sample prediction values
838 print(TestingDataResults.head())
839
840 # Printing sample prediction values
841 print(TestingDataResults.head())
842
843 # Printing sample prediction values
844 print(TestingDataResults.head())
845
846 # Printing sample prediction values
847 print(TestingDataResults.head())
848
849 # Printing sample prediction values
850 print(TestingDataResults.head())
851
852 # Printing sample prediction values
853 print(TestingDataResults.head())
854
855 # Printing sample prediction values
856 print(TestingDataResults.head())
857
858 # Printing sample prediction values
859 print(TestingDataResults.head())
860
861 # Printing sample prediction values
862 print(TestingDataResults.head())
863
864 # Printing sample prediction values
865 print(TestingDataResults.head())
866
867 # Printing sample prediction values
868 print(TestingDataResults.head())
869
870 # Printing sample prediction values
871 print(TestingDataResults.head())
872
873 # Printing sample prediction values
874 print(TestingDataResults.head())
875
876 # Printing sample prediction values
877 print(TestingDataResults.head())
878
879 # Printing sample prediction values
880 print(TestingDataResults.head())
881
```

```
682 # Calculating the error for each row
683 TestingDataResults['APE']=100 * ((abs(
684     TestingDataResults['Selling_Price']-
TestingDataResults['PredictedSelling_Price']))/
TestingDataResults['Selling_Price']))
685
686
687 MAPE=np.mean(TestingDataResults['APE'])
688 MedianMAPE=np.median(TestingDataResults['APE'])
689
690 Accuracy =100 - MAPE
691 MedianAccuracy=100- MedianMAPE
692 print('Mean Accuracy on test data:', Accuracy) # Can
    be negative sometimes due to outlier
693 print('Median Accuracy on test data:',
    MedianAccuracy)
694
695
696 # Defining a custom function to calculate accuracy
697 # Make sure there are no zeros in the Target
    variable if you are using MAPE
698 def Accuracy_Score(orig,pred):
699     MAPE = np.mean(100 * (np.abs(orig-pred)/orig))
700     #print('#'*70, 'Accuracy:', 100-MAPE)
701     return(100-MAPE)
702
703 # Custom Scoring MAPE calculation
704 from sklearn.metrics import make_scorer
705 custom_Scoring=make_scorer(Accuracy_Score,
    greater_is_better=True)
706
707 # Importing cross validation function from sklearn
708 from sklearn.model_selection import cross_val_score
709
710 # Running 10-Fold Cross validation on a given
    algorithm
711 # Passing full data X and y because the K-fold will
    split the data and automatically choose train/test
712 Accuracy_Values=cross_val_score(RegModel, X , y, cv=
    10, scoring=custom_Scoring)
713 print('\nAccuracy values for 10-fold Cross
```

```
713 Validation:\n',Accuracy_Values)
714 print('\nFinal Average Accuracy of the model:', round(Accuracy_Values.mean(),2))
715 #%% md
716 #Plotting a single Decision tree out of XGBoost
717 #%%
718 #Plotting a single Decision tree out of XGBoost
719 from xgboost import plot_tree
720 import matplotlib.pyplot as plt
721 fig, ax = plt.subplots(figsize=(20, 8))
722 plot_tree(XGB, num_trees=10, ax=ax)
723 #%% md
724 # K-Nearest Neighbor(KNN)
725 #%%
726 #KNN
727 # K-Nearest Neighbor(KNN)
728 from sklearn.neighbors import KNeighborsRegressor
729 RegModel = KNeighborsRegressor(n_neighbors=3)
730
731 # Printing all the parameters of KNN
732 print(RegModel)
733
734 # Creating the model on Training Data
735 KNN=RegModel.fit(X_train,y_train)
736 prediction=KNN.predict(X_test)
737
738 from sklearn import metrics
739 # Measuring Goodness of fit in Training data
740 print('R2 Value:',metrics.r2_score(y_train, KNN.predict(X_train)))
741
742 # Plotting the feature importance for Top 10 most
    important columns
743 # The variable importance chart is not available for
    KNN
744
745 ##########
##########
746 print('\n##### Model Validation and Accuracy
    Calculations #####')
747
```

```
748 # Printing some sample values of prediction
749 TestingDataResults=pd.DataFrame(data=X_test, columns
=Predictors)
750 TestingDataResults[TargetVariable]=y_test
751 TestingDataResults[('Predicted'+TargetVariable)]=np.
round(prediction)
752
753 # Printing sample prediction values
754 print(TestingDataResults.head())
755
756 # Calculating the error for each row
757 TestingDataResults['APE']=100 * ((abs(
758     TestingDataResults['Selling_Price']-
    TestingDataResults['PredictedSelling_Price']))/
    TestingDataResults['Selling_Price'])
759
760 MAPE=np.mean(TestingDataResults['APE'])
761 MedianMAPE=np.median(TestingDataResults['APE'])
762
763 Accuracy =100 - MAPE
764 MedianAccuracy=100- MedianMAPE
765 print('Mean Accuracy on test data:', Accuracy) # Can
    be negative sometimes due to outlier
766 print('Median Accuracy on test data:',
    MedianAccuracy)
767
768 # Defining a custom function to calculate accuracy
769 # Make sure there are no zeros in the Target
    variable if you are using MAPE
770 def Accuracy_Score(orig,pred):
771     MAPE = np.mean(100 * (np.abs(orig-pred)/orig))
772     #print('#'*70,'Accuracy:', 100-MAPE)
773     return(100-MAPE)
774
775 # Custom Scoring MAPE calculation
776 from sklearn.metrics import make_scorer
777 custom_Scoring=make_scorer(Accuracy_Score,
    greater_is_better=True)
778
779 # Importing cross validation function from sklearn
780 from sklearn.model_selection import cross_val_score
```

```
781
782 # Running 10-Fold Cross validation on a given
    algorithm
783 # Passing full data X and y because the K-fold will
    split the data and automatically choose train/test
784 Accuracy_Values=cross_val_score(RegModel, X , y, cv=
    10, scoring=custom_Scoring)
785 print('\nAccuracy values for 10-fold Cross
    Validation:\n',Accuracy_Values)
786 print('\nFinal Average Accuracy of the model:', round(Accuracy_Values.mean(),2))
787 #%% md
788 # Step 21: Model Deployment
789 * I am choosing XGboost as the final model it has
    the highest accuracy(35.78).
790
791 * For this dataset, the most important predictor
    variables are 'Product_Brand','Subcategory_1', '
    Subcategory_2', 'Item_Category'.
792 #%%
793 # Separate Target Variable and Predictor Variables
794 TargetVariable='Selling_Price'
795
796 # Selecting the final set of predictors for the
    deployment
797 # Based on the variable importance charts of
    multiple algorithms above
798 Predictors=[column for column in DataForML_Numeric.
    columns if column != TargetVariable]
799
800 X=DataForML_Numeric[Predictors].values
801 y=DataForML_Numeric[TargetVariable].values
802
803 ### Standardization of data ####
804 from sklearn.preprocessing import StandardScaler,
    MinMaxScaler
805 # Choose either standardization or Normalization
806 # On this data Min Max Normalization produced better
    results
807
808 # Choose between standardization and MinMax
```

```
808 normalization
809 #PredictorScaler=StandardScaler()
810 PredictorScaler=MinMaxScaler()
811
812 # Storing the fit object for later reference
813 PredictorScalerFit=PredictorScaler.fit(X)
814
815 # Generating the standardized values of X
816 X=PredictorScalerFit.transform(X)
817
818 print(X.shape)
819 print(y.shape)
820 #%% md
821 # Cross validating the final model accuracy with
     less predictors
822 #%%
823 # Importing cross validation function from sklearn
824 from sklearn.model_selection import cross_val_score
825
826 # choose from different tunable hyper parameters
827 from xgboost import XGBRegressor
828 RegModel=XGBRegressor(max_depth=2,
829                         learning_rate=0.1,
830                         n_estimators=1000,
831                         objective='reg:linear',
832                         booster='gbtree')
833
834 # Running 10-Fold Cross validation on a given
     algorithm
835 # Passing full data X and y because the K-fold will
     split the data and automatically choose train/test
836 Accuracy_Values=cross_val_score(RegModel, X , y, cv=
     10, scoring=custom_Scoring)
837 print('\nAccuracy values for 10-fold Cross
     Validation:\n',Accuracy_Values)
838 print('\nFinal Average Accuracy of the model:', round(Accuracy_Values.mean(),2))
839 #%% md
840 # Step 22: Retraining the final model using 100%
     data
841 #%%
```

```
842 # Training the model on 100% Data available
843 Final_XGB_Model=RegModel.fit(X,y)
844 #%% md
845 # Step 23: Save the model as a serialized file which
     can be stored anywhere
846 #%%
847 import pickle
848 import os
849
850 # Saving the Python objects as serialized files can
     be done using pickle library
851 # Save the Final model
852 with open('Final_XGB_Model.pkl', 'wb') as
     fileWriteStream:
853     pickle.dump(Final_XGB_Model, fileWriteStream)
854     # Close the filestream!
855     fileWriteStream.close()
856
857 print('pickle file of Predictive Model is saved at
     Location:',os.getcwd())
858 #%% md
859 # Step 24: Create a python function
860 #%%
861 from re import IGNORECASE
862
863 def FunctionPredictResult(InputData):
864     import pandas as pd
865     Num_Inputs=InputData.shape[0]
866
867     # Appending the new data with the Training data
868     DataForML=pd.read_pickle('DataForML.pkl')
869     #InputData=InputData.append(DataForML,
     ignore_index=True)
870     InputData = pd.concat([InputData, DataForML],
     ignore_index=True)
871
872     # Generating dummy variables for rest of the
     nominal variables
873     InputData=pd.get_dummies(InputData)
874
875     # Maintaining the same order of columns as it
```

```
875 was during the model training
876     Predictors=[column for column in
877         DataForML_Numeric.columns if column != TargetVariable]
878
879     # Generating the input values to the model
880     X=InputData[Predictors].values[0:Num_Inputs]
881
882     # Generating the standardized values of X since
883     # it was done while model training also
884     X=PredictorScalerFit.transform(X)
885
886     # Loading the Function from pickle file
887     import pickle
888     with open('Final_XGB_Model.pkl', 'rb') as
889         fileReadStream:
890             PredictionModel=pickle.load(fileReadStream)
891             # Don't forget to close the filestream!
892             fileReadStream.close()
893
894     # Generating Predictions
895     Prediction=PredictionModel.predict(X)
896     PredictionResult=pd.DataFrame(Prediction,
897         columns=['Prediction'])
898     return(PredictionResult)
899
900 #%% md
901 # Step 25: Calling the function for some new data
902 #%%
903 # Load the test dataframe to test the function
904 test_data = pd.read_csv('D:\PricePrediction\Test.csv')
905 NewSampleData=pd.DataFrame(data=test_data,columns=['Product_Brand','Item_Category','Subcategory_1','Subcategory_2'])
906 print(NewSampleData)
907 # Passing the loaded DataFrame
908 FunctionPredictResult(InputData=NewSampleData)
909
910 #%% md
911 # Conclusion
912 * The Function FunctionPredictResult() can be used
```

```
907 to produce the predictions for one or more new cases  
at a time.  
908 * Hence, it can be scheduled using a batch job or  
cron job to run every night and generate predictions  
for all the products in the E-commerce platform.  
909 #%% md  
910 # Desktop App deployment: Tkinter package  
911  
912 #%%  
913 import tkinter as tk  
914 from tkinter import messagebox  
915 from tkinter import ttk  
916 import pandas as pd  
917 from sklearn.model_selection import train_test_split  
918 from xgboost import XGBRegressor  
919  
920 class SellingPricePredictionApp:  
921     def __init__(self, master):  
922         self.master = master  
923         self.master.title('Price Prediction')  
924  
925         # Read dataset for training and applying one  
-hot encoding  
926         self.data = pd.read_csv('Train.csv')  
927  
928         # Applying one-hot encoding  
929         self.data = pd.get_dummies(self.data)  
930  
931         # List to potentially store slider widgets  
932         self.sliders = []  
933  
934         # Separating features (X) and target  
variable (y)  
935         self.X = self.data.drop('Selling_Price',  
axis=1).values  
936         self.y = self.data['Selling_Price'].values  
937  
938         # Splitting data into training and testing  
sets  
939         self.X_train, self.X_test, self.y_train,  
self.y_test = train_test_split(self.X, self.y,
```

```
939 test_size=0.2, random_state=42)
940
941     # Initialize and train the XGBoost
942     regression model
943     self.model = XGBRegressor()
944     self.model.fit(self.X_train, self.y_train)
945
946     # Create GUI widgets for user interaction
947     self.create_widgets()
948
949     def create_widgets(self):
950         raw_data = pd.read_csv('Train.csv')
951         dropdowns = {}
952
953         # To track rows for placing widgets
954         row_idx = 0
955
956         # Loop through each categorical column to
957         # create dropdown menus
958         for column in ['Product_Brand', 'Item_Category', 'Subcategory_1', 'Subcategory_2']:
959             label = tk.Label(self.master, text=
960             column + ': ')
961             label.grid(row=row_idx, column=0)
962             options = list(raw_data[column].unique
963             ())
964             # Sort options alphabetically
965             options = sorted(options)
966             # for storing dropdown selection
967             selected_value = tk.StringVar(self.
968             master)
969             # Default to first option
970             selected_value.set(options[0])
971             dropdown = ttk.Combobox(self.master,
972             textvariable=selected_value, values=options)
973             dropdown.grid(row=row_idx, column=1)
974             dropdowns[column] = selected_value
975             row_idx += 1
976
977             # Button for price prediction
978             predict_button = tk.Button(self.master, text
```

```
972 ='Predict Price', command=lambda: self.  
    predict_price(dropdowns))  
973         predict_button.grid(row=row_idx, columnspan  
=2)  
974  
975     def predict_price(self, dropdowns):  
976         # Convert dropdown selections to model  
        input format  
977         input_data = {col: val.get() for col, val  
in dropdowns.items()}  
978  
979         # Create a new DataFrame and convert to one  
-hot encoding format  
980         input_df = pd.DataFrame([input_data])  
981         input_df = pd.get_dummies(input_df)  
982         # Fill missing columns with 0  
983         for col in self.data.columns:  
984             if col not in input_df.columns:  
985                 input_df[col] = 0  
986  
987         # Ensure the correct feature order  
988         input_df = input_df[self.data.drop('  
Selling_Price', axis=1).columns]  
989  
990         # Predict and Print results  
991         price = self.model.predict(input_df)  
992         messagebox.showinfo('Predicted Price', f'  
The predicted selling price is ${price[0]:.2f}')  
993  
994  
995 if __name__ == '__main__':  
996     root = tk.Tk()  
997     app = SellingPricePredictionApp(root)  
998     root.mainloop()  
999  
1000 %% md  
1001 # END OF PROGRAMMING PROJECT
```