

e-ISSN: 2319-8753 | p-ISSN: 2347-6710

TEDIA

International Journal of Innovative Research in SCIENCE | ENGINEERING | TECHNOLOGY

INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

Volume 11, Issue 6, June 2022

INTERNATIONAL STANDARD SERIAL NUMBER INDIA

Impact Factor: 8.118

9940 572 462

6381 907 438

 \bigcirc





| e-ISSN: 2319-8753, p-ISSN: 2320-6710| <u>www.ijirset.com</u> | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |

Using Machine Learning Algorithms to Forecast Huge Big Mart Sales

P Poornima, K Tharun Teja

Assistant Professor, Department of Computer Science and Engineering, Mahatma Gandhi Institute of Technology

Hyderabad, India

Department of Computer Science and Engineering, Mahatma Gandhi Institute of Technology

Hyderabad, India

ABSTRACT: Machine Learning is a category of algorithms that allows software applications to become more accurate in predicting outcomes without being explicitly programmed. The basic premise of machine learning is to build models and employ algorithms that can receive input data and use statistical analysis to predict an output while updating outputs as new data becomes available. These models can be applied in different areas and trained to match the expectations of management so that accurate steps can be taken to achieve the organization's target. The case of Big Mart, a one-stop-shopping centre, has been discussed to predict the sales of different types of items and for understanding the effects of different factors on the items' sales. Taking various aspects of a dataset collected for Big Mart, and the methodology followed for building a predictive model, results with high levels of accuracy are generated, and these observations can be employed to take decisions to improve sales.

KEYWORDS: Machine Learning, Random Forest, Linear Regression, Decision Tree.

I. INTRODUCTION

In today's modern world, huge shopping centers such as big malls and marts are recording data related to sales of items or products with their various dependent or independent factors as an important step to be helpful in prediction of future demands and inventory management. The dataset built with various dependent and independent variables is a composite form of item attributes, data gathered by means of customer, and also data related to inventory management in a data warehouse. The data is thereafter refined in order to get accurate predictions and gather new as well as interesting results that shed a new light on our knowledge with respect to the task's data. This can then further be used for forecasting future sales by means of employing machine learning algorithms such as the random forests and linear regression model.

1.1 Problem Statement

"To find out what role certain properties of an item play and how they affect their sales by understanding Big Mart sales." In order to help Big Mart achieve this goal, a predictive model can be built to find out for every store, the key factors that can increase their sales and what changes could be made to the product or store's characteristics

1.2 Existing System

The BigMart sales analysis can be done by many algorithms but presently working algorithms computational time is very high and the predicted accuracy of these algorithms is very less and these algorithms are not good in finding patterns. So, we are finding new ways to predicting the accuracy of bigmart sales data and finding patterns in the given data.

1.3 Proposed System

For building a model to predict accurate results the dataset of Big Mart sales undergoes several sequence of steps and in this work, we propose a model using Xgboost technique. Every step plays a vital role for building the proposed model. In our model we have used 2013 Big mart dataset.



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| www.ijirset.com | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

DOI:10.15680/IJIRSET.2022.1106182

After pre-processing and filling missing values, we used ensemble classifier using Decision trees, Linear regression, Ridge regression and Random forest regression. Both MSE and CVS are used as accuracy metrics for predicting the sales in Big Mart. From the accuracy metrics it was found that the model will predict best using minimum MSE and CVS. The details of the proposed method is explained in the following section.

II. LITERATURE SURVEY

- 1. A Forecast for Big Mart Sales Based on Random Forests and Multiple Linear Regression used Random Forest and Linear Regression for prediction analysis which gives less accuracy. To overcome this we can use XG boost Algorithm which will give more accuracy and will be more efficient.
- 2. Forecasting methods and applications contains Lack of Data and short life cycles. So some of the data like historical data, consumer-oriented markets face uncertain demands, can be prediction for accurate result.
- 3. Comparison of Different Machine Learning Algorithms for Multiple Regression on Black Friday Sales Data Used Neural Network for comparison of different algorithms. To overcome this Complex models like neural networks is used for comparison between different algorithms which is not efficient so we can use more simpler algorithm for prediction.
- 4. Prediction of retail sales of footwear using feed forward and recurrent neural networks used neural networks for prediction of sales. Using neural network for predicting of weekly retail sales, which is not efficient, So XG boost can work efficiently.

III. SYSTEM DESIGN OF BIGMART SALES PREDICTION

3.1 Dataset Information

The data scientists at Big Mart have collected 2013 sales data for 1559 products across 10 stores in different cities. Also, certain attributes of each product and store have been defined. The aim is to build a predictive model and find out the sales of each product at a particular store.

Using this model, Big Mart will try to understand the properties of products and stores which play a key role in increasing sales.

Variable	Description
Item_Identifier	Unique product ID
Item_Weight	Weight of product
Item_Fat_Content	Whether the product is low fat or not
Item_Visibility	The % of total display area of all products in a store allocated to the particular product
Item_Type	The category to which the product belongs
Item_MRP	Maximum Retail Price (list price) of the product
Outlet_Identifier	Unique store ID
Outlet_Establishment_Year	The year in which store was established
Outlet_Size	The size of the store in terms of ground area covered
Outlet_Location_Type	The type of city in which the store is located
Outlet_Type	Whether the outlet is just a grocery store or some sort of supermarket
Item_Outlet_Sales	Sales of the product in the particulat store. This is the outcome variable to be predicted

Fig 3.1 Attributes of the dataset

The above figure 3.1 shows the attributes included in the dataset

3.2 Data Preprocessing

Data pre-processing is a data mining technique that involves transforming raw data into an understandable format. Real-world data is often incomplete, inconsistent, and/or lacking in certain behavior or trends, and is likely to contain many errors. Data pre-processing is a proven method of resolving such issues.

Three common data preprocessing steps are:

- Formatting
- Cleaning



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| <u>www.ijirset.com</u> | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

DOI:10.15680/IJIRSET.2022.1106182

• Sampling

3.3. Label Encoding

In machine learning, we usually deal with datasets that contain multiple labels in one or more than one column. These labels can be in the form of words or numbers. To make the data understandable or in human-readable form, the training data is often labelled in words.

Label Encoding refers to converting the labels into a numeric form so as to convert them into the machine-readable form. Machine learning algorithms can then decide in a better way how those labels must be operated. It is an important pre-processing step for the structured dataset in supervised learning.

3.4. Onehot Encoding

A one hot encoding is a representation of categorical variables as binary vectors. This first requires that the categorical values be mapped to integer values. Then, each integer value is represented as a binary vector that is all zero values except the index of the integer, which is marked with a 1.

3.5. Machine Learning Algorithms Used:

The data available is increasing day by day and such a huge amount of unprocessed data is needed to be analyzed precisely, as it can give very informative and finely pure gradient results as per current standard requirements. It is not wrong to say as with the evolution of Artificial Intelligence (AI) over the past two decades, Machine Learning (ML) is also on a fast pace for its evolution. ML is an important mainstay of IT sector and with that, a rather central, albeit usually hidden, part of our life. As the technology progresses, the analysis and understanding of data to give good results will also increase as the data is very useful in current aspects. In machine learning, one deals with both supervised and unsupervised types of tasks and generally a classification type problem accounts as a resource for knowledge discovery

1. Linear Regression Algorithm

Regression can be termed as a parametric technique which is used to predict a continuous or dependent variable on basis of a provided set of independent variables. This technique is said to be parametric as different assumptions are made on basis of data set.

$$Y = \beta o + \beta 1 X + \in (1)$$

Equation shown in eq.1 is used for simple linear regression. These parameters can be said as:

Y - Variable to be predicted

X - Variable(s) used for making a prediction

 β o - When X=0, it is termed as prediction value or can be referred to as intercept term

 β 1 - when there is a change in X by 1 unit it denotes change in Y. It can also be said as slope term

 \in -The difference between the predicted and actual values is represented by this parameter and represents the residual value. However efficiently the model is trained, tested, and validated, there is always a difference between actual and predicted values which is irreducible error thus we cannot rely completely on the predicted results by the learning algorithm. Alternative methods given by Dieterich can be used for comparing learning algorithms.

2. Decision Tree:

Decision Tree is one of the most used, practical approaches for supervised learning. It can be used to solve both Regression and Classification tasks with the latter being put more into practical application.

It is a tree-structured classifier with three types of nodes. The *Root Node* is the initial node which represents the entire sample and may get split further into further nodes. The *Interior Nodes* represent the features of a data set and the branches represent the decision rules. Finally, the *Leaf Node* represent the outcome. This algorithm is very useful for solving decision-related problems.

With a particular data point, it is run completely through the entirely tree by answering *True/False* questions till it reaches the leaf node. The final prediction is the average of the value of the dependent variable in that particular leaf node. Through multiple iterations, the Tree is able to predict a proper value for the data point.

3. Random Forest Algorithm

Random forest algorithm is a very accurate algorithm to be used for predicting sales. It is easy to use and understand for the purpose of predicting results of machine learning tasks. In sales prediction, random forest classifier is used because it has decision tree like hyperparameters. The tree model is same as decision tool. Fig.5 shows the relation between decision trees and random forest. To solve regression tasks of prediction by virtue of random forest, the



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| www.ijirset.com | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |

sklearn.ensemble library's random forest regressor class is used. The key role is played by the parameter termed as n_estimators which also comes under random forest regressor. Random forest can be referred to as a meta-estimator used to fit upon numerous decision trees (based on classification) by taking the dataset's different sub-samples. min_samples_split is taken as the minimum number when splitting an internal node if integer number of minimum samples are considered. A split's quality is measured using mse (mean squared error), which can also be termed as feature selection criterion. This also means reduction in variance mae (mean absolute error), which is another criterion for feature selection. Maximum tree depth, measured in integer terms, if equals one, then all leaves are pure or pruning for better model fitting is done for all leaves less than min_samples_split samples.



Fig 3.5.1 Trees generation

The above figure 3.5.1 shows the relation between the Decision Tree and the Random Forest models and how the trees are generated

IV. IMPLEMENTATION

1. Importing libraries and reading the dataset

<pre>import pandas as pd</pre>
import numpy as np
<pre>import seaborn as sns</pre>
<pre>import matplotlib.pyplot as plt</pre>
import warnings
%matplotlib inline
<pre>warnings.filterwarnings('ignore')</pre>

Fig 4.1.1 Importing libraries

[n [8]:	df df.	= pd.read_cs	sv('Train.cs	sv')							
Dut[8]:		Item_Identifier	Item_Weight	Item_Fat_Content	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment_Year	Outlet_Size	Outlet_Location
	0	FDA15	9.30	Low Fat	0.016047	Dairy	249.8092	OUT049	1999	Medium	
	1	DRC01	5.92	Regular	0.019278	Soft Drinks	48.2692	OUT018	2009	Medium	
	2	FDN15	17.50	Low Fat	0.016760	Meat	141.6180	OUT049	1999	Medium	
	3	FDX07	19.20	Regular	0.000000	Fruits and Vegetables	182.0950	OUT010	1998	NaN	
	4	NCD19	8.93	Low Fat	0.000000	Household	53.8614	OUT013	1987	High	

Fig 4.1.2 Reading the dataset

1 [9]: #	# <i>stat</i> df.des	istical in j cribe()	Fo			
ut[9]:		Item_Weight	Item_Visibility	Item_MRP	Outlet_Establishment_Year	Item_Outlet_Sales
	count	7060.000000	8523.000000	8523.000000	8523.000000	8523.000000
	mean	12.857645	0.066132	140.992782	1997.831867	2181.288914
	std	4.643456	0.051598	62.275067	8.371760	1706.499616
	min	4.555000	0.000000	31.290000	1985.000000	33.290000
	25%	8.773750	0.026989	93.826500	1987.000000	834.247400
	50%	12.600000	0.053931	143.012800	1999.000000	1794.331000
	75%	16.850000	0.094585	185.643700	2004.000000	3101.296400
	max	21.350000	0.328391	266.888400	2009.000000	13086.964800

Fig 4.1.3. Getting the statistical info of the dataset



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| <u>www.ijirset.com</u> | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

DOI:10.15680/LJIRSET.2022.1106182

2. Exploring the dataset

In [10]: # datatype of attributes
df.info()

<pre><class 'pandas.core.frame.dataframe'=""> RangeIndex: 8523 entries, 0 to 8522 Data columns (total 12 columns):</class></pre>									
#	Column	Non-Null Count	Dtype						
0	Item Identifier	8523 non-null	object						
1	Item_Weight	7060 non-null	float64						
2	Item_Fat_Content	8523 non-null	object						
3	Item_Visibility	8523 non-null	float64						
4	Item_Type	8523 non-null	object						
5	Item_MRP	8523 non-null	float64						
6	Outlet_Identifier	8523 non-null	object						
7	Outlet_Establishment_Year	8523 non-null	int64						
8	Outlet_Size	6113 non-null	object						
9	Outlet_Location_Type	8523 non-null	object						
10	Outlet_Type	8523 non-null	object						
11	Item_Outlet_Sales	8523 non-null	float64						
dtyp	dtypes: float64(4), int64(1), object(7)								
memory usage: 799.2+ KB									

Fig 4.2.1 Datatype of attributes

In [11]:	<pre># check unique values in dataset df.apply(lambda x: len(x.unique()))</pre>						
Out[11]:	Item Identifier	1559					
	Item Weight	416					
	Item Fat Content	5					
	Item Visibility	7880					
	Item Type	16					
	Item_MRP	5938					
	Outlet Identifier	10					
	Outlet Establishment Year	9					
	Outlet_Size	4					
	Outlet_Location_Type	3					
	Outlet Type	4					
	Item_Outlet_Sales	3493					
	dtype: int64						

Fig 4.2.2 checking unique values

3. Data preprocessing

In [12]:	<pre># check for null values df.isnull().sum()</pre>		
Out[12]:	Item Identifier	0	
	Item Weight	1463	
	Item Fat Content	0	
	Item Visibility	0	
	Item Type	0	
	Item MRP	0	
	Outlet Identifier	0	
	Outlet Establishment Year	0	
	Outlet Size	2410	
	Outlet Location Type	0	
	Outlet Type	0	
	Item_Outlet_Sales dtype: int64	0	

Fig 4.3.1 Checking for null values

In [13]:	<pre># check for categorical attributes cat_col = [] if df.dtypes.index: if df.dtypes[x] == 'object': cat_col.append(x) cat_col</pre>
Out[13]:	<pre>['Item Jdentifier', 'Item Jac Content', 'Item Jype', 'Outlet_Identifier', 'Outlet_Size', 'Outlet_Size', 'Outlet_Size', 'Outlet_Size', 'Outlet_Type']</pre>
In [14]:	<pre>cat_col.remove('Item_Identifier') cat_col.remove('Outlet_Identifier') cat_col</pre>
Out[14]:	['Item_Fat_Content', 'Item_Type', 'Outlet_Size', 'Outlet_Cocation_Type', 'Outlet_Type']





| e-ISSN: 2319-8753, p-ISSN: 2320-6710| <u>www.ijirset.com</u> | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

DOI:10.15680/LJIRSET.2022.1106182

<pre>[15]: # print the categor for col in cat_col: print(col) print(df[col],v</pre>	ical columns	add to circ
print()		outlet_Size
Item Fat Content Low Fat 5089 Regular 2889 LF 316 reg 117		Medium 2793 Small 2388 High 932 Name: Outlet_Size, dtype: int64
low fat 112		
Name: Item_Fat_Cont	ent, dtype: int64	Outlet Location Type
Item_Type		dicte_beddan_jpe
Fruits and Vegetabl	es 1232	Tier 3 3350
Snack Foods	1200	Tier 2 2785
Household	910	
Frozen Foods	856	lier 1 2388
Dairy	682	Name: Outlet Location Type dtype; int64
Canned	649	halle, outret_coutron_type, utype, intou
Baking Goods	648	
Health and Hygiene	520	Outlat Type
SOTU DELINKS	445	oucceripe
Recorde	423	Supermarket Type1 5577
Hard Drinks	214	Grocery Store 1083
others	169	
Starchy Eoods	148	Supermarket Type3 935
Breakfast	110	Supermarket Type2 928
Seafood	64	Supermanice Typez 525
Name: Item Type, dt	vpe: int64	Name: Outlet Type, dtype: int64



[10]:						
	Item_Identifier	æm_Weight				
	DRA12	11.600				
	DRA24	19.350				
	DRA59	8.270				
	DRB01	7.390				
	DRB13	6.115				
	S					
	NCZ30	6.590				
	NCZ41	19.850				
	NCZ42	10.500				
	NCZ53	9.600				
	NCZ54	14.650				
	1555 rows × 1 co	lumns				
[17]:	<pre>miss_bool = df[miss_bool</pre>	'Item_Weight'].ismull()			
[17]:	0 False					
	1 False					
	Z False					
	4 False					
	8518 False					
	8519 False					
	8520 False					
	8522 False					
	Name: Item_Weig	ght, Length: 8	523, dtype: bool			
	<pre>for i, item in enumerate(df['Item_Identifier']): if miss bool[i]: if item in item_weight_mean:</pre>					
[18]:	for i, item in if miss boo if item	enumerate(df[bl[i]: i in item_weig	'Item_Identifier']): ht_mean:			





Fig 4.3.6 Filling missing values for Item visibility

e-ISSN: 2319-8753, p-ISSN: 2320-6710 www.ijirset.com | Impact Factor: 8.118

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |

[26]: # d d	f comb f['It f['It	ine it em_Fat em_Fat	em fat co _Content' _Content'	ontent] = df['Item_Fa].value_counts(t_Content'].rep)	place({'Lf	':'Low Fat', 'r	eg':'Regu	ilar', 'low fa	t':'Low Fat'	'})
[26]: L R	ow Fa	t S r 3 Ttem F	517 006 at Conter	nt. dtype: int64							
				Fig	g 4.3.7 Coi	mbinin	g item fat c	conten	t		
n [27]:	df['N df['N	ew tte ew_tte	n Type'] n_Type']	= df['Iten Ident	ifier'].apply(1	anbda x:)	([:2])				
it[27]:	0	FD									
	1	DR									
	2	FD									
	4	NC									
	8518	FD									
	8520	NC									
	8521	FD									
	8522	DR	ton Turn	Longth, 8522	tunni abiast						
	wane:	NEM_I	cen_type,	cengen: 8523, 6	type: object						
[28]:	df["N df["N	ew_Ite ew_Ite	n_Type'] n_Type'].	<pre>= df['New_Item_T value_counts()</pre>	ype'].map({'FD'	:'Food', '	'NC': 'Non-Consuma	ble', 'DF	(':'Drinks'})		
11 28 1:	Food		61	25							
100	Non-C	onsuna	ble 15	99							
	Drink	s	7	99							
	Name:	New_1	ten type,	dtype: int64							
1 [20]:	df In	eldfi'	Now Them	Type 'les 'Non-Con	sumable' 'ttom	Eat Conto	ont'l = 'Mon Edit	de'			
. [6].	df['I	tem_Fa	t_Content].value_counts()	- ac_contr	ine j = mar care				
t1291:	Low E	at	2010								
arter li	Regul	ar	3006								
	Non-D	dible	1599								
	Name:	Iten_	Fat_Conte	nt, dtype: int64	ŧ.						
					Fig 4.	3.8Nev	v item type				
					U		21				
In	[30]:	# crea	te small v	alues for establis	shment year	nt Year'l					
				J - Lors and o							
In	[31]:	df['Ou	tlet_Years	1							
Out	t[31]:	0	14								
		2	14								
		3	26								
		8518									
		8519	11								
		8520 8521	9								
		8522	16 Outlet Mar	ing Longthy 0522	dtumos inted						
		Name:	Outlet_Yea	rs, Length: 8523,	dtype: int64						
In	[32]:	df.hea	d()								
Out	t[32]:	em_Type	Item_MRP	Outlet_Identifier Outle	et_Establishment_Year	Outlet_Size	Outlet_Location_Type	Outlet_Type	Item_Outlet_Sales	New_Item_Type	Outlet_Years
		Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1	3735.1380	Food	14
		off Drinke	49 2692	OUT018	2000	Medium	Tion 2	Supermarket	443 4220	Drinke	
		en erniks	40.2002	001010	2009	Medium	ner a	Type2	440.4228	Diriks	
		Meat	141.6180	OUT049	1999	Medium	Tier 1	Supermarket Type1	2097.2700	Food	14
		ruits and	182 0950	OUT010	1998	Small	Tier 3	Grocery	732 3800	Food	15

Fig 4.3.9 Creating small values for establishment year

High

1987

OUT013

53.8614

V. DATA VISUALISATION





Tier 3 Supermarket Type1 994.7052

Con



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| www.ijirset.com | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |











Fig 5.4 Item outlet sales plot



Fig 5.5 Log transformation plot



Fig 5.6 Item fat content plot



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| <u>www.ijirset.com</u> | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |



Fig 5.8 Outlet establishment year plot

200



In [43]: sns.countplot(df['Outlet_Location_Type'])
Out[43]: <AxesSubplot:xlabel='Outlet_Location_Type', ylabel='count'>



e-ISSN: 2319-8753, p-ISSN: 2320-6710 www.ijirset.com | Impact Factor: 8.118



Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |







6.1 Preprocessing for categorical attributes

Label encoding											
	In [47]:	<pre>from sklearn.preprocessing import LabelEncoder le = LabelEncoder() df['outlet'] = Lefit_transform(df['outlet_Identifier']) cat_col = ['Item_rat_Content', 'Item_type', 'outlet_size', 'Outlet_tocation_type', 'Outlet_type', 'New_Item_type'] for col in cat_col: df[col] = Ie.fit_transform(df[col])</pre>									
		Onehot	Encod	ling							
	In [48]; Out[48]:	df = pd.get df.head()	_dumies	(df, columns=['item_	Fat_Content', 'Out]	et_Size', 'Outlet_L	ocation_Type'	, 'Outlet_Type',	, 'New_Item_	(ype'])	
		Outlet_Years	Outlet	Outlet_Location_Type_0	Outlet_Location_Type_1	Outlet_Location_Type_2	Outlet_Type_0	Outlet_Type_1 Out	tlet_Type_2 O	utlet_Type_3 N	
		14	9	1	0	C	0	1	0	0	
		4	3	0	0	1	0	0	1	0	
		14	9	1	0	0	0	1	0	0	
		15	0	0	0	× 1	8 8	0	0	0	
		28	1	0	0	1	0	1	0	0	
									-		

Fig 6.1 label, Onehot encoding



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| www.ijirset.com | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |

6.2 Model training and splitting the dataset

In [49]:	<pre>X = df.drop(columns=['Outlet_Establishment_Year', 'Item_Identifier', 'Outlet_Identifier', 'Item_Outlet_Sales']) y = df['Item_Outlet_Sales']</pre>
	Model Training
In [62]:	<pre>from sklearn.model_selection import cross_val_score from sklearn.metrics import mean_squared_error,confusion_matrix # train the model model.fit(x, y) # predict the training set pred = model.predict(x) # perform cross-validation cv_score = ross_val_score(model, x, y, scoring='neg_mean_squared_error', cv=5) cv_score = np.abs(np.mean(cv_score)) print("Model Report") print("Model Report") print("Model Report") print("Model Report") </pre>

Fig 6.2 Dataset split and model training

6.3 Linear regression



Fig 6.3 Linear Regression

The above figure 6.3 shows prediction of linear regression and taken title as Model Coefficient where we have got the MSE = 0.2880787833 and CV score = 0.28920602795. MSE and CV score is very minimal so that this model trained and predicted efficiently. The model normalizes the cofficients so that when x increases, y decreases or y increases, x decreases. So it is proven the basic models also predict efficiently.

6.4 Decision Tree Regressor



Fig 6.4 Decision Tree Regressor



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| <u>www.ijirset.com</u> | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |

The above figure 6.4 shows prediction of linear regression and taken title as Feature Importance where we have got the MSE = 5.55340306385e-34 and CV score = 0.5750919329. MSE is minimal and CV score is low so that this model trained efficiently and it predicted well.

6.5 Random Forest Regressor



Fig 6.5 Random Forest Regressor

The above figure 5.5 shows prediction of linear regression and taken title as Feature Importance where we have got the MSE = 0.0422380173 and CV score = 0.3102969540. MSE is minimal and CV score is low so that this model trained efficiently and also it predicted well.

VII. CONCLUSION AND FUTURE SCOPE

The basics of machine learning and the associated data processing and modeling algorithms are described and followed by their application for the task of sales prediction in Big Mart shopping centers at different locations. On implementation, the prediction results show the correlation among different attributes considered and how a particular location of medium size recorded the highest sales, suggesting that other shopping locations should follow similar patterns for improved sales.

Multiple instances parameters and various factors can be used to make this sales prediction more innovative and successful. Accuracy, which plays a key role in prediction-based systems, can be significantly increased as the number of parameters used are increased. Also, a look into how the sub-models work can lead to increase in productivity of system. The paper can be further collaborated in a web-based application or in any device supported with an in-built intelligence by virtue of Internet of Things (IoT), to be more feasible for use. Various stakeholders concerned with sales information can also provide more inputs to help in hypothesis generation and more instances can be taken into consideration such that more precise results that are closer to real world situations are generated. When combined with effective data mining methods and properties, the traditional means could be seen to make a higher and positive effect on the overall development of corporation's tasks on the whole.

REFERENCES

- [1] Cerrada, M., & Aguilar, J. (2008). Reinforcement learning in system identification.
- [2] Downey, A. B. (2011). Think stats. "O'Reilly Media, Inc.".
- [3] Géron, A. (2019). Hands-On Machine Learning with ScikitLearn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems. O'Reilly Media.
- [4] Learning, M. (1994). Neural and Statistical Classification. Editors D. Mitchie et. al, 350.
- [5] MacKay, D. J., & Mac Kay, D. J. (2003). Information theory, inference and learning algorithms. Cambridge university press.
- [6] Mitchell, T. M. (1999). Machine learning and data mining. Communications of the ACM, 42(11), 30-36.
- Cambridge University, UK, 32, 34.
- [7] Saltz, J. S., & Stanton, J. M. (2017). An introduction to data science. Sage Publications.
- [8] Shashua, A. (2009). Introduction to machine learning: Class notes 67577. arXiv preprint arXiv:0904.3664.



| e-ISSN: 2319-8753, p-ISSN: 2320-6710| <u>www.ijirset.com</u> | Impact Factor: 8.118|

Volume 11, Issue 6, June 2022

| DOI:10.15680/IJIRSET.2022.1106182 |

[9] Smola, A., & Vishwanathan, S. V. N. (2008). Introduction to machine learning.

- [10] Quinlan, J. R. (2014). C4. 5: programs for machine learning. Elsevier. In Reinforcement Learning. IntechOpen.
- [11] Welling, M. (2011). A first encounter with Machine Learning. Irvine, CA.: University of California, 12.





Impact Factor: 8.118





INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY





www.ijirset.com