

Predicting Big Mart Sales Using Machine Learning Models in KNIME

Overview:

In this project, machine learning workflows were built to forecast Big Mart's product sales by utilizing the visual interface of KNIME. The retail sales dataset was used to create models such as linear regression and random forest. Data import, exploration, preprocessing, modeling, testing, and result generation were all covered by the workflows, which were based on the precise actions taken during the project.

Training data set:

Item_Identifier	Item_Weight	Item_Fat_Content	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment_Year	Outlet_Size	Outlet_Location_Type	Outlet_Type	Item_Outlet_Sales
1	9.3	Low Fat	0.016047201	Dairy	249.8092	OUT049	1999	Medium	Tier 1	Supermarket Type1	3735.138
2	5.92	Regular	0.015278216	Soft Drinks	48.2692	OUT018	2009	Medium	Tier 3	Supermarket Type2	443.4228
3	17.5	Low Fat	0.016760075	Meat	141.618	OUT049	1999	Medium	Tier 1	Supermarket Type1	2097.27
4	19.2	Regular	0	Fruits and Vegetables	182.095	OUT010	1998	Tier 3		Grocery Store	732.38
5	8.93	Low Fat	0	Household	53.8614	OUT013	1987	High	Tier 3	Supermarket Type1	994.7052
6	10.395	Regular	0	Baking Goods	51.4008	OUT018	2009	Medium	Tier 3	Supermarket Type2	556.6088
7	13.65	Regular	0.012741089	Snack Foods	57.6588	OUT013	1987	High	Tier 3	Supermarket Type1	343.5528
8	12.7468857	Low Fat	0.012746857	Snack Foods	107.7622	OUT027	1985	Medium	Tier 3	Supermarket Type3	4022.7636
9	16.2	Regular	0.015687114	Frozen Foods	96.9726	OUT045	2002	Tier 2		Supermarket Type1	1076.5986
10	19.2	Regular	0.09444559	Frozen Foods	187.8214	OUT017	2007	Tier 2		Supermarket Type1	4710.535
11	11.8	Low Fat	0	Fruits and Vegetables	45.5402	OUT049	1999	Medium	Tier 1	Supermarket Type1	1516.0266
12	18.5	Regular	0.045463773	Dairy	144.1102	OUT046	1997	Small	Tier 1	Supermarket Type1	2187.153
13	15.1	Regular	0.1000135	Fruits and Vegetables	145.4786	OUT049	1999	Medium	Tier 1	Supermarket Type1	1589.2646
14	17.6	Regular	0.047257328	Snack Foods	119.6782	OUT046	1997	Small	Tier 1	Supermarket Type1	2145.2076
15	16.35	Low Fat	0.0680243	Fruits and Vegetables	196.4426	OUT013	1987	High	Tier 3	Supermarket Type1	1977.426
16	9	Regular	0.069088961	Breakfast	56.3614	OUT046	1997	Small	Tier 1	Supermarket Type1	1547.3192
17	11.8	Low Fat	0.008596051	Health and Hygiene	115.3452	OUT010	2009	Medium	Tier 3	Supermarket Type2	1621.8888
18	9	Regular	0.069196376	Breakfast	54.3614	OUT049	1999	Medium	Tier 1	Supermarket Type1	718.3982
19	13.2834	Low Fat	0.034237682	Hard Drinks	113.2834	OUT027	1985	Medium	Tier 3	Supermarket Type3	2303.668
20	13.35	Low Fat	0.10249212	Dairy	230.5352	OUT035	2004	Small	Tier 2	Supermarket Type1	2748.4224
21	18.85	Regular	0.138190277	Snack Foods	250.8724	OUT013	1987	High	Tier 3	Supermarket Type1	3775.086
22	14.4	Low Fat	0.035399923	Baking Goods	144.5444	OUT027	1985	Medium	Tier 3	Supermarket Type3	4064.0432
23	14.6	Low Fat	0.025698134	Household	196.5084	OUT035	2004	Small	Tier 2	Supermarket Type1	1587.2672
24	10.7555998	Low Fat	0.05755998	Baking Goods	107.6938	OUT019	1985	Small	Tier 1	Grocery Store	214.3876
25	13.85	Regular	0.025896485	Frozen Foods	165.021	OUT046	1997	Small	Tier 1	Supermarket Type1	4078.025
26	13	Low Fat	0.099887103	Household	45.906	OUT017	2007	Tier 2		Supermarket Type1	838.908
27	7.645	Regular	0.066693437	Snack Foods	42.3112	OUT035	2004	Small	Tier 2	Supermarket Type1	1065.28
28	11.65	Low Fat	0.019356132	Hard Drinks	39.1164	OUT013	1987	High	Tier 3	Supermarket Type1	308.9312
29	5.925	Regular	0.161466534	Dairy	45.5086	OUT010	1998	Tier 3		Grocery Store	178.4344
30	17.2221801	Regular	0.072221801	Canned	43.6454	OUT019	1985	Small	Tier 1	Grocery Store	125.8362
31	19.25	Low Fat	0.170348551	Dairy	55.7596	OUT010	1998	Tier 3		Grocery Store	163.7868
32	18.6	Low Fat	0.080829372	Health and Hygiene	96.4436	OUT018	2009	Medium	Tier 3	Supermarket Type2	2741.7644
33	18.7	Low Fat	0	Snack Foods	256.6672	OUT018	2009	Medium	Tier 3	Supermarket Type2	3068.0064
34	17.85	Low Fat	0	Breads	93.1436	OUT045	2002	Tier 2		Supermarket Type1	2174.5028
35	17.5	Low Fat	0.097904029	Soft Drinks	174.8738	OUT046	1997	Small	Tier 1	Supermarket Type1	2085.2856

The Big Mart sales dataset provided contains retail transaction data for predicting product sales at specific store outlets. It has been split into a training and test set for developing and evaluating machine learning models.

source: <https://datahack.analyticsvidhya.com/contest/practice-problem-big-mart-sales-iii/#ProblemStatement>

The training dataset is available as a CSV file with 8523 rows and 12 columns. Each row represents a product sold at a particular store and the columns provide attributes of the item, outlet details, and the target sales variable to predict. The columns are:

Dega_INLS625_LAB04_KNIME_lightening_challenge

Item_Identifier (Numeric): Unique ID for each product

Item_Weight (Numeric): Weight of the product

Item_Fat_Content (Categorical): Whether the product is low fat or not.

Item_Visibility (Numeric): Display area percentage for product

Item_Type (Categorical): Category of the product

Item_MRP (Numeric): Maximum retail price of the product

Outlet_Identifier (Numeric): Unique ID for each outlet

Outlet_Establishment_Year (Numeric): Year when the outlet was established.

Outlet_Size (Categorical): Size of the outlet in sq. ft.

Outlet_Location_Type (Categorical): Type of city where the outlet is located.

Outlet_Type (Categorical): Grocery or supermarket outlet

Item_Outlet_Sales (Numeric): Sales of the product at the outlet. This is the target variable.

The test dataset contains 5681 rows and 11 columns. It has the same 11 predictor columns as training without the target variable.

This data can be used to develop a machine learning model to predict the sales for the test products at given outlets based on their attributes. The model's performance can then be evaluated by comparing predictions to actual sales.

Testing Data set:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
1	Item_Identifier	Item_Weight	Item_Fat_Content	Item_Visibility	Item_Type	Item_MRP	Outlet_Identifier	Outlet_Establishment_Year	Outlet_Size	Outlet_Location_Type	Outlet_Type									
2	FDW50	20.75	Low Fat	0.007564836	Snack Foods	107.8622	OUT049	1999	Medium	Tier 1	Supermarket Type1									
3	FDW14	8.3	reg	0.038427677	Dairy	87.3198	OUT017	2007		Tier 2	Supermarket Type1									
4	NCN55	34.6	Low Fat	0.099574908	Others	241.7538	OUT010	1998		Tier 3	Grocery Store									
5	FDQ58	7.315	Low Fat	0.015388393	Snack Foods	155.034	OUT017	2007		Tier 2	Supermarket Type1									
6	FDY38		Regular	0.118599314	Dairy	234.23	OUT027	1985	Medium	Tier 3	Supermarket Type3									
7	FDV56	9.8	Regular	0.063812206	Fruits and Vegetables	117.1492	OUT046	1997	Small	Tier 1	Supermarket Type1									
8	FDL48	19.35	Regular	0.082601517	Baking Goods	50.1034	OUT018	2009	Medium	Tier 3	Supermarket Type2									
9	FDC48		Low Fat	0.015782495	Baking Goods	81.0592	OUT027	1985	Medium	Tier 3	Supermarket Type3									
10	FDN33	6.305	Regular	0.123365446	Snack Foods	95.7436	OUT045	2002		Tier 2	Supermarket Type1									
11	FDA36	5.985	Low Fat	0.005698435	Baking Goods	186.8924	OUT017	2007		Tier 2	Supermarket Type1									
12	FDI44	16.6	Low Fat	0.103566075	Fruits and Vegetables	118.3466	OUT017	2007		Tier 2	Supermarket Type1									
13	FDQ56	6.59	Regular	0.10581147	Fruits and Vegetables	85.3908	OUT045	2002		Tier 2	Supermarket Type1									
14	NCC54		Low Fat	0.171079215	Health and Hygiene	240.4196	OUT019	1985	Small	Tier 1	Grocery Store									
15	FDU11	4.785	Low Fat	0.092737611	Breads	122.3098	OUT049	1999	Medium	Tier 1	Supermarket Type1									
16	DRJ59	16.75	LF	0.021206466	Hard Drinks	32.0298	OUT013	1987	High	Tier 3	Supermarket Type1									
17	FDN424	6.125	Regular	0.0794507	Baking Goods	151.6266	OUT049	1999	Medium	Tier 1	Supermarket Type1									
18	FDI57	19.85	Low Fat	0.05413521	Seafood	198.7768	OUT045	2002		Tier 2	Supermarket Type1									
19	DRJ12	17.85	Low Fat	0.037980963	Soft Drinks	192.2188	OUT018	2009	Medium	Tier 3	Supermarket Type2									
20	NCM42		Low Fat	0.028184344	Household	109.6912	OUT027	1985	Medium	Tier 3	Supermarket Type3									
21	FDH46	18.8	Low Fat	0.196897637	Snack Foods	193.7136	OUT010	1998		Tier 3	Grocery Store									
22	FDA31	7.1	Low Fat	0.109920138	Fruits and Vegetables	175.008	OUT013	1987	High	Tier 3	Supermarket Type1									
23	NCJ31	19.2	Low Fat	0.182619235	Others	239.9196	OUT035	2004	Small	Tier 2	Supermarket Type1									
24	FDG52	13.65	LF	0.065630844	Frozen Foods	47.7402	OUT046	1997	Small	Tier 1	Supermarket Type1									
25	NCL19		Low Fat	0.027447057	Others	142.347	OUT019	1985	Small	Tier 1	Grocery Store									
26	FDI10	19.2	Low Fat	0.035178915	Snack Foods	180.7318	OUT035	2004	Small	Tier 2	Supermarket Type1									
27	FDX22	6.785	Regular	0.038455125	Snack Foods	209.4928	OUT010	1998		Tier 3	Grocery Store									
28	NCF19	13	Low Fat	0.035102094	Household	47.6034	OUT035	2004	Small	Tier 2	Supermarket Type1									
29	NCE06	5.825	Low Fat	0.091485232	Household	161.3894	OUT046	1997	Small	Tier 1	Supermarket Type1									
30	DRJ27	13.8	Low Fat	0.058102469	Dairy	244.6802	OUT046	1997	Small	Tier 1	Supermarket Type1									
31	FDZ21	12.8	LF	0.022940249	Fruits and Vegetables	116.5492	OUT035	2004	Small	Tier 2	Supermarket Type1									
32	NCR42		Low Fat	0.06737681	Household	32.09	OUT019	1985	Small	Tier 1	Grocery Store									
33	FDX51	9.5	Regular	0.022148582	Meat	194.9452	OUT018	2009	Medium	Tier 3	Supermarket Type2									
34	NCR06	12.5	Low Fat	0.006792707	Household	42.4112	OUT018	2009	Medium	Tier 3	Supermarket Type2									
35	FDU31		Regular	0.024870035	Fruits and Vegetables	217.7508	OUT027	1985	Medium	Tier 3	Supermarket Type3									
36	FDU59	5.78	Low Fat	0.096931426	Breads	164.2552	OUT017	2007		Tier 2	Supermarket Type1									

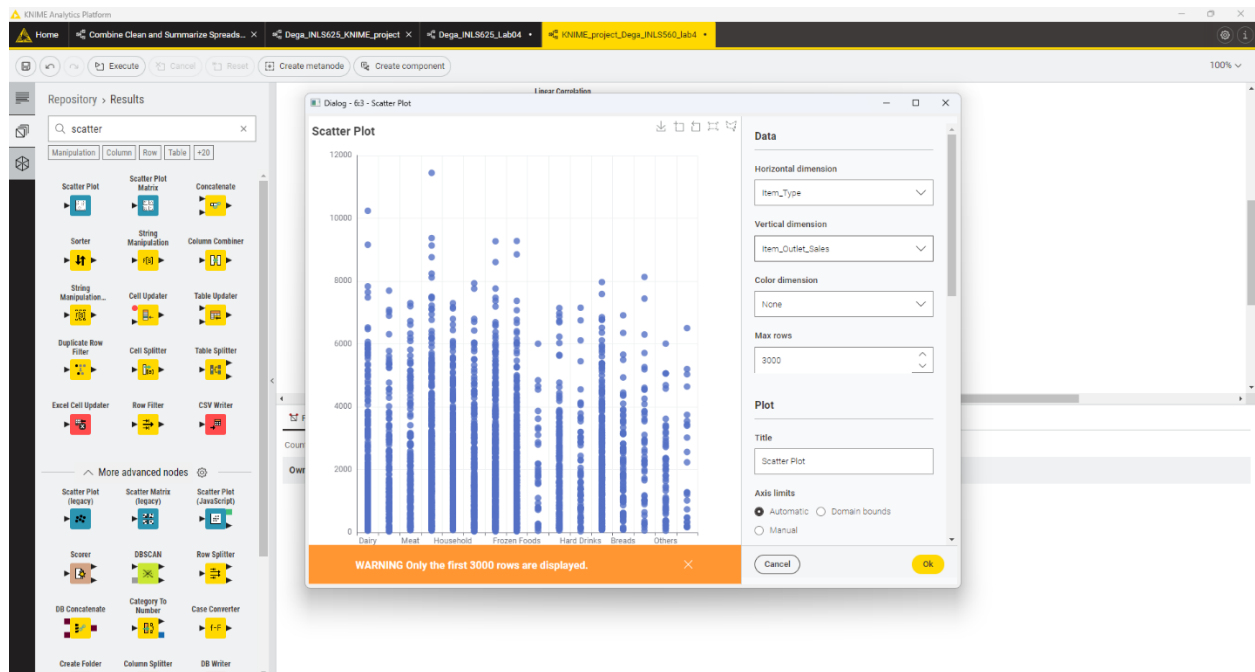
Dega_INLS625_LAB04_KNIME_lightening_challenge

Interactive View Statistics

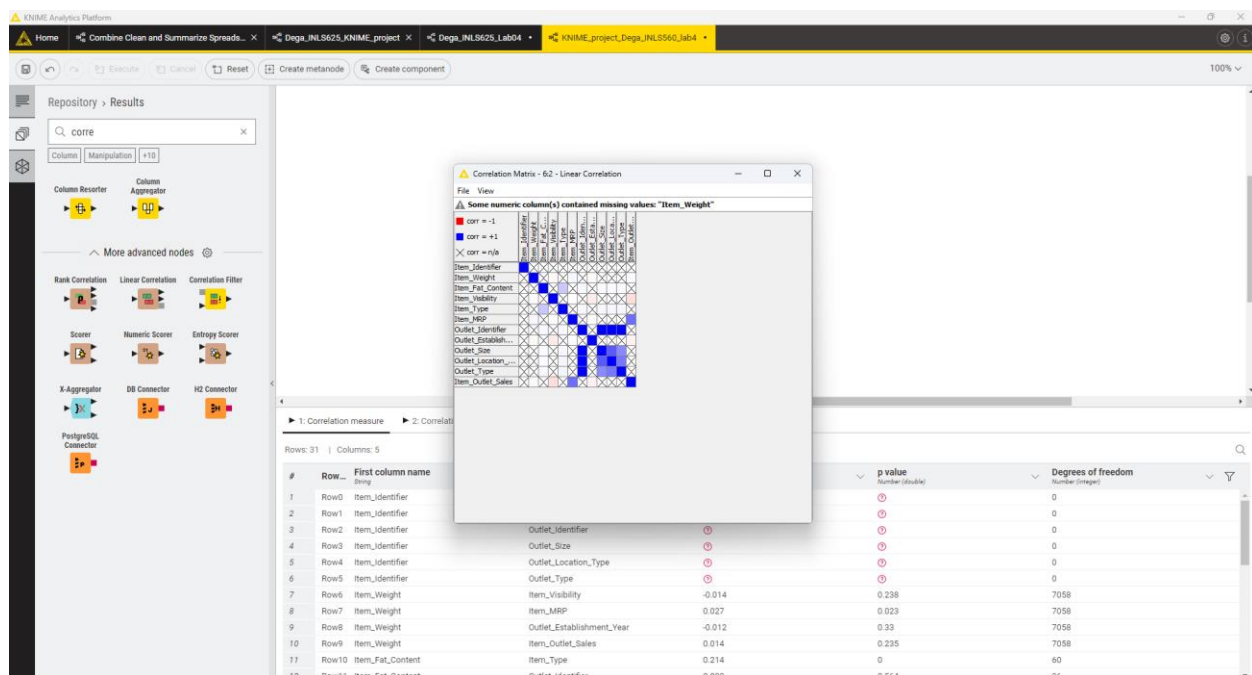
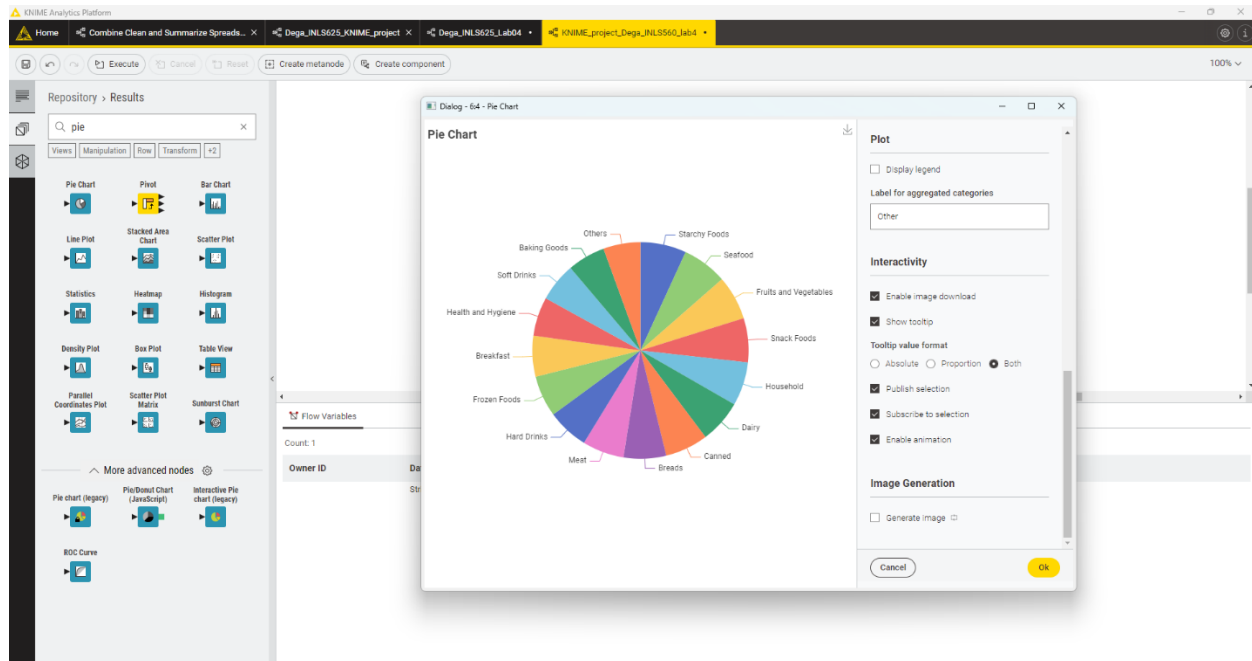
Rows: 12 | Columns: 14

Statistics

Name	Type	# Missing values	# Unique values	Minimum	Maximum	25% Quantile	50% Quantile (M...	75% Quantile	Mean	Mean Absolute ...	Standard Deviat...	Sum	10 most commo...
Item_Identifier	String	0	1559	0	0	0	0	0	0	0	0	0	FD033 (10, 0.12%)...
Item_Weight	Number (double)	1463	415	4.555	21.35	8.771	12.6	16.85	12.858	4.046	4.643	90,774.975	12.15 (86, 1.22%)...
Item_Fat_Content	String	0	5	0	0	0	0	0	0	0	0	0	Low Fat (5089, 59...
Item_Visibility	Number (double)	0	7880	0	0.328	0.027	0.054	0.095	0.066	0.041	0.052	563.643	0 (526, 6.17%), 0.0...
Item_Type	String	0	16	0	0	0	0	0	0	0	0	0	Fruits and Vegeta...
Item_MRP	Number (double)	0	5938	31.29	266.888	93.809	143.013	185.661	140.993	52.539	62.275	1,201,681.481	172.042 (7, 0.08%)...
Outlet_Identifier	String	0	10	0	0	0	0	0	0	0	0	0	OUT027 (935, 10.9...
Outlet_Establishm...	Number (integer)	0	9	1,985	2,009	1,987	1,999	2,004	1,997.832	6.956	8.372	17,027,521	1,985 (1463, 17.17...
Outlet_Size	String	2410	3	0	0	0	0	0	0	0	0	0	Medium (2793, 45...
Outlet_Location_T...	String	0	3	0	0	0	0	0	0	0	0	0	Tier 3 (3350, 39.3...
Outlet_Type	String	0	4	0	0	0	0	0	0	0	0	0	Supermarket Type...
Item_Outlet_Sales	Number (double)	0	3493	33.29	13,086.965	833.582	1,794.331	3,101.296	2,181.289	1,346.571	1,706.5	18,591,125.41	958.752 (17, 0.2%)...



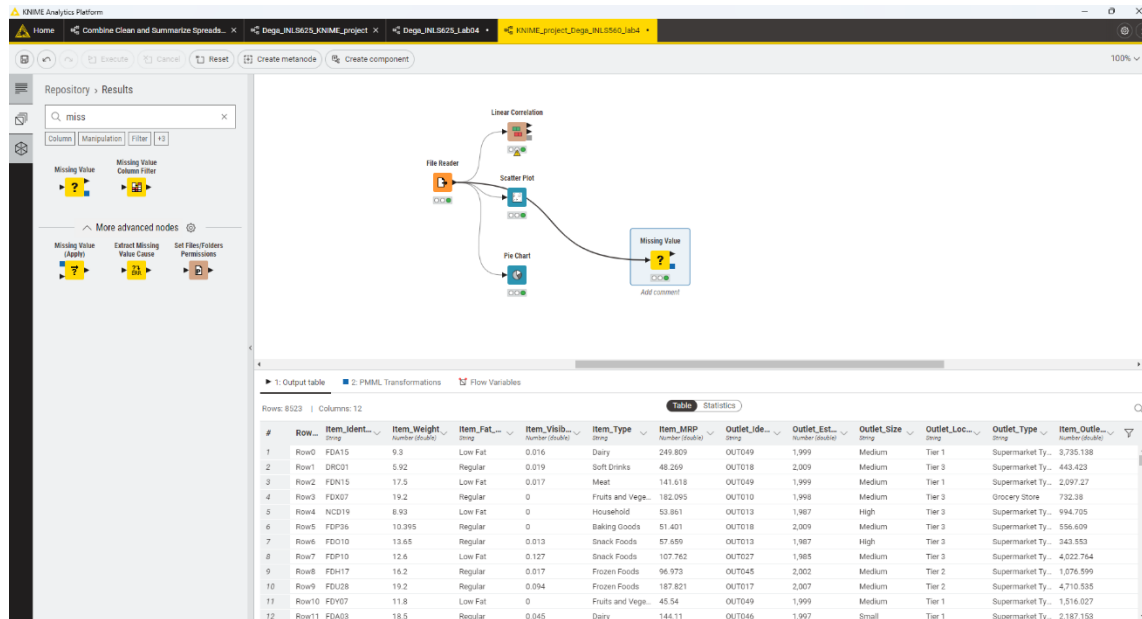
Dega_INLS625_LAB04_KNIME_lightening_challenge



Data Preprocessing

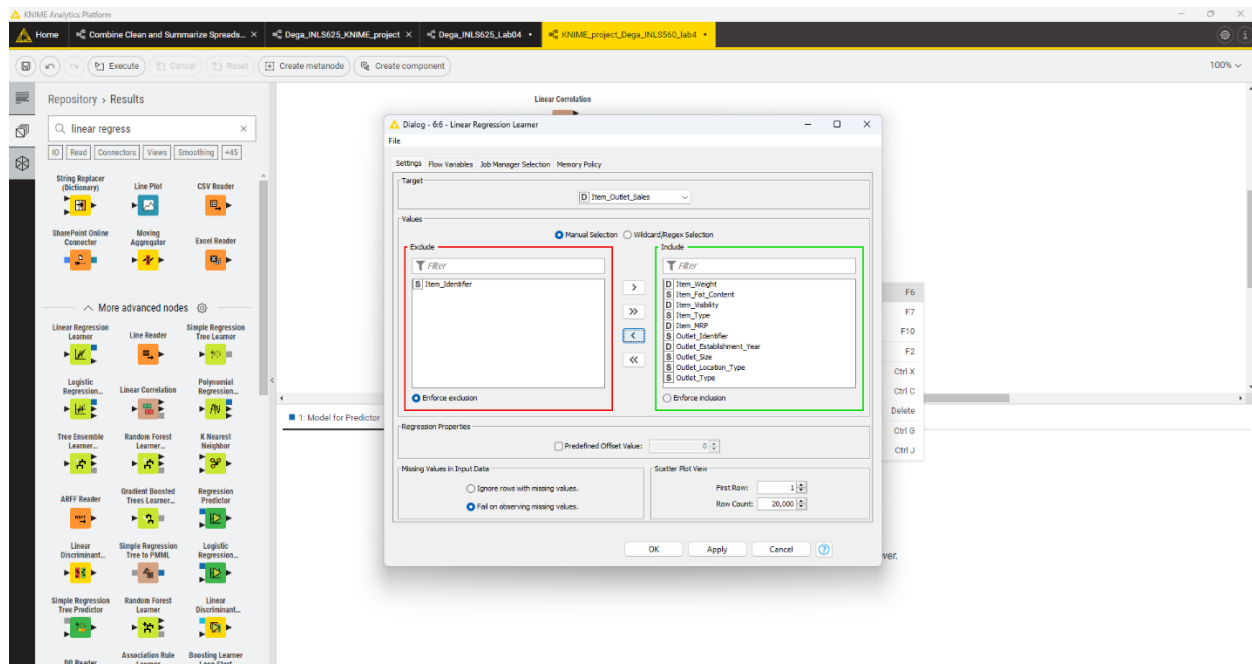
The raw training data had issues like missing values and non-numeric columns which were addressed before modeling. The Missing Value node helped identify and replace missing values in numerical columns with a median. These steps output clean, normalized training data ready for ML modeling.

Dega_INLS625_LAB04_KNIME_lightening_challenge

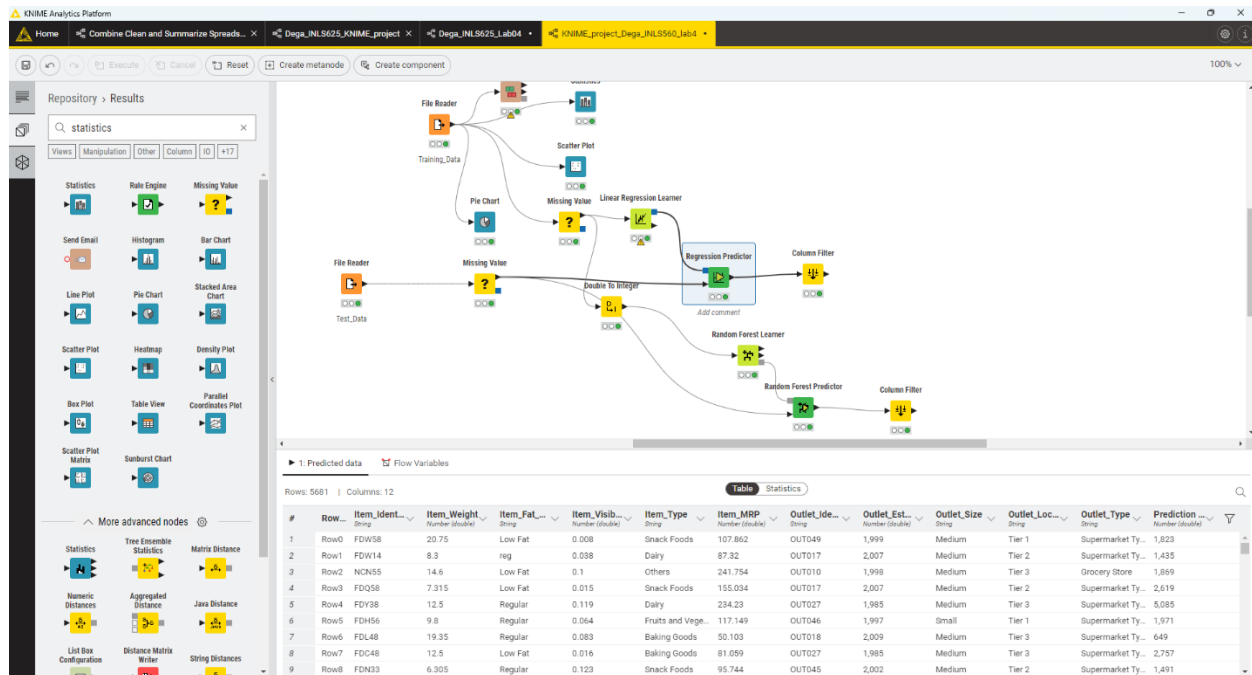


Linear Regression Model

A linear regression model was built to predict the Item_Outlet_Sales. The Linear Regression Learner node was added, and the cleaned training data was connected to its input port. In the node configuration, Item_Outlet_Sales was set as the target column while numeric columns like Item_MRP, Outlet_Establishment_Year, etc. were chosen as predictors. The node was executed to train the regression model.



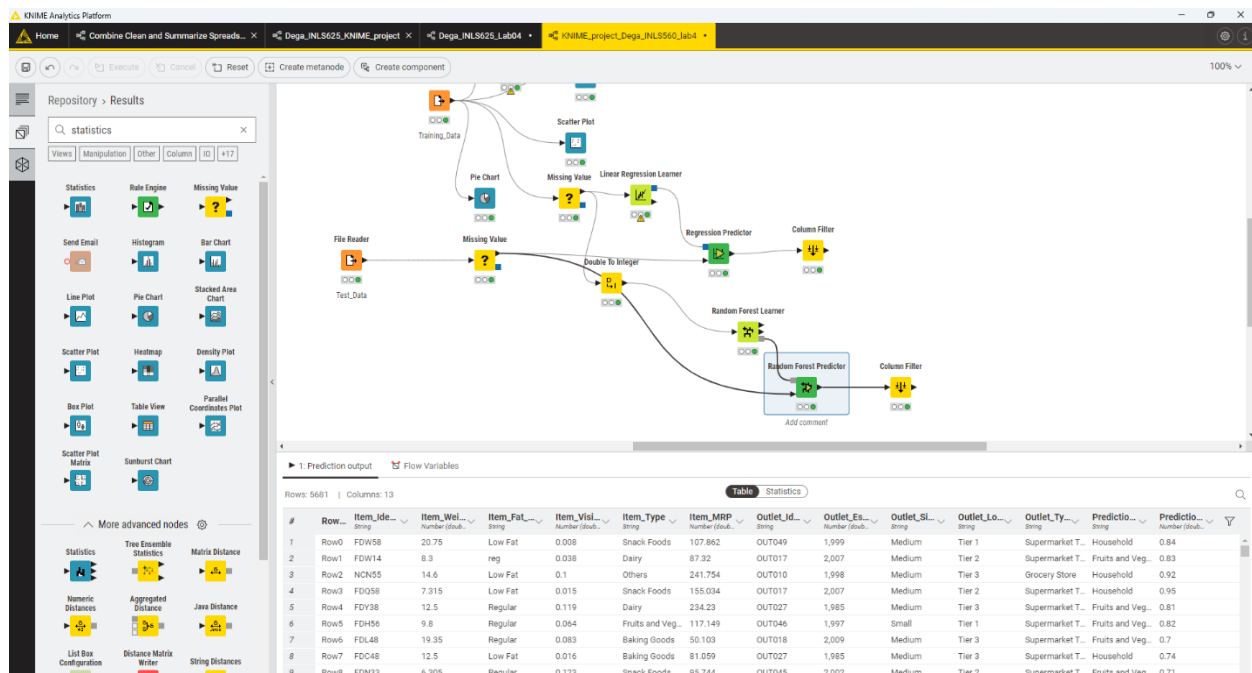
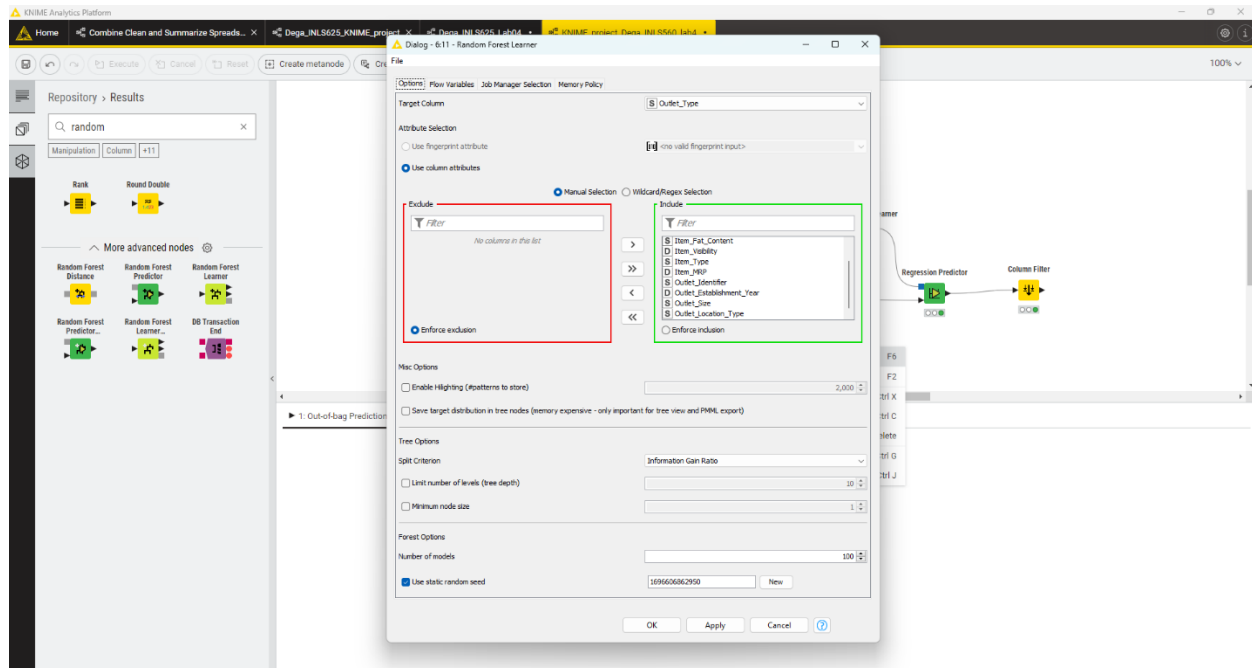
The trained model was connected to a Linear Regression Predictor node. The clean test data was connected to the predictor's second input port. After executing this workflow, the predicted sales values were compared to actuals to evaluate performance on the test set.



Random Forest Model

A random forest classifier was built with Outlet_type as the target variable. The Random Forest Learner node was configured to use the preprocessed training data. Outlet_type was selected as the target column while columns like item_fat_content, outlet_size, etc. were chosen as predictors after converting Double to an integer using the Double to integer node.

Dega_INLS625_LAB04_KNIME_lightening_challenge

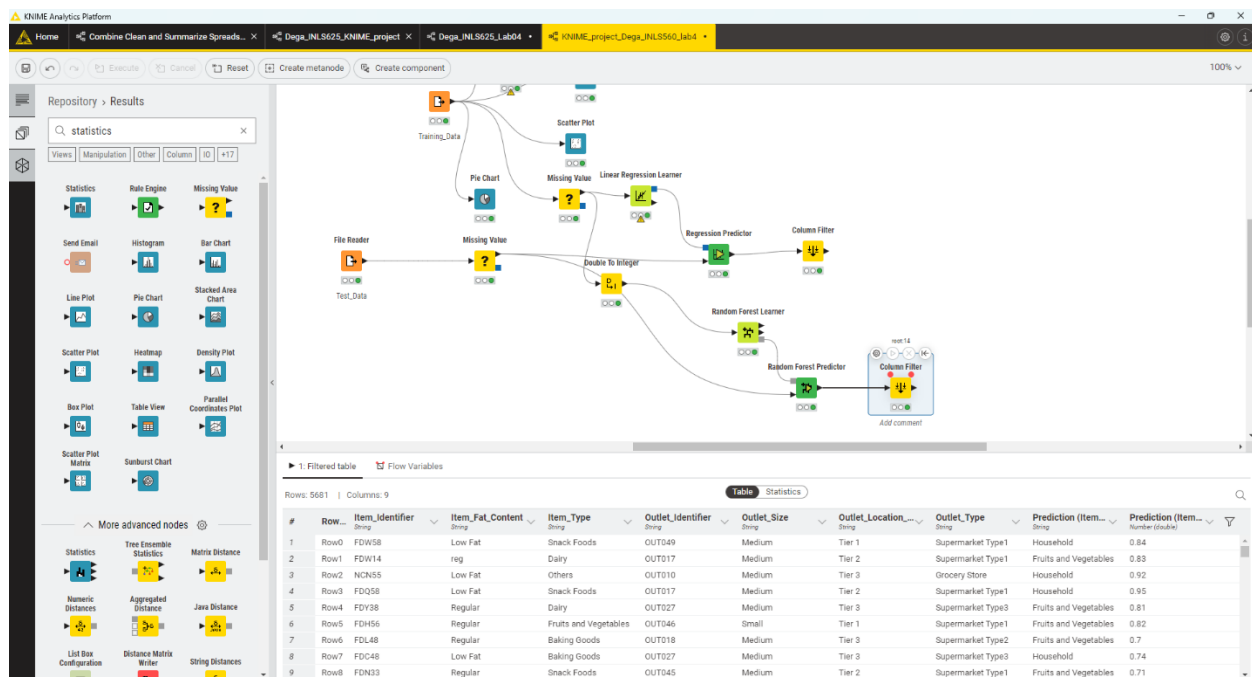
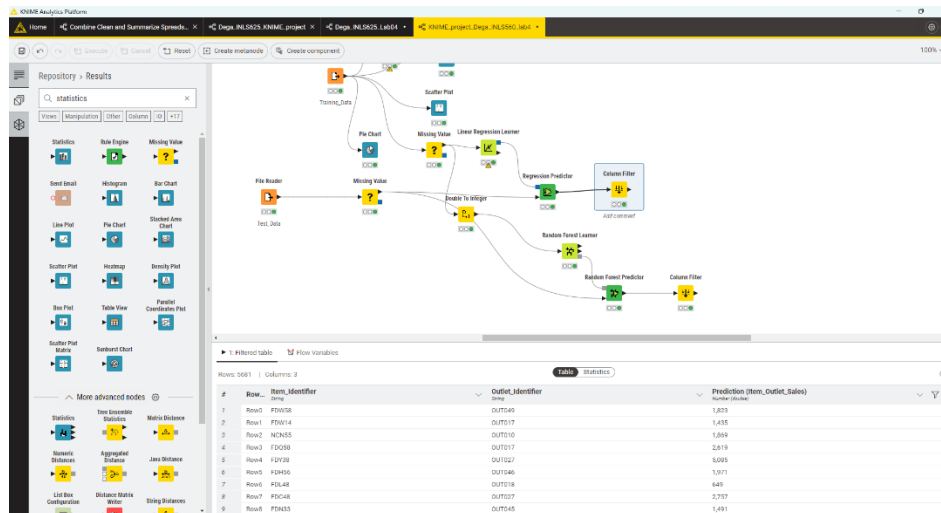


The trained random forest model was connected to the Random Forest Predictor node. The preprocessed test data was fed to the second input port. The predicted item types were compared to actual values to assess classification accuracy on the test data.

Result Analysis

Dega_INLS625_LAB04_KNIME_lightening_challenge

The Column Filter node was used to output only the required columns from the predictors. Summary statistics and scoring metrics were computed to evaluate model performance. The predictions were reasonably accurate and can be further improved by tuning models and adding more relevant features.



Conclusion:

This project provided hands-on experience with KNIME's capabilities for building end-to-end ML workflows visually without coding. The linear regression model performed decently for sales prediction. The random forest model also had reasonable accuracy in classifying item types. The workflows can be extended by adding more data, features, and tuning techniques.

List of other project sites that you searched and explored before choosing your project (at least 5):

1. Healthcare fraud Source: <https://www.kaggle.com/datasets/tamilsel/healthcare-providers-data>
2. Medical Practitioner and Prescription Behavior-National Survey of Physician Characteristics and Prescription Practices in the USA
<https://www.kaggle.com/datasets/tubmak/dataset>
3. Titanic - ML Project <https://hub.knime.com/augustinejoseph/spaces/Public/latest/Titanic%20-%20ML%20Project~TsA5hTL4pnFIfDHg>
4. ECG Heartbeat Categorization Dataset <https://www.kaggle.com/datasets/shayanfazeli/heartbeat>
https://hub.knime.com/knime/spaces/Digital%20Healthcare/latest/ECG%20Arrhythmia%20Detection/ecg_cnn_mit~bWv0UtH6sTLAgcnn
5. Triage Score Prediction
<https://hub.knime.com/knime/spaces/Digital%20Healthcare/latest/Triage%20Score%20Prediction/Triage%20Score%20Prediction~-wdfGA8UecAI42L8>