

FOUNDATION FOR ORGANISATIONAL RESEARCH AND EDUCATION NEW DELHI

Academic Session 2023-2025 Project-2

Customer Classification and Prediction (Car Prices) on the basis of Cluster data Machine Learning for Managers

FMG 32 Section A

Submitted to:

Submitted by:

Prof. Amarnath Mitra

321032 – Nisha Arora

1. Project Objectives

- 1.1 The first objective is to segment the consumer (loan) data of the bank using supervised learning algorithms using Decision tree.
- 1.2 The second objective is to determine the number of appropriate classification model by comparing and contrast using logistic regression, KNN (k-nearest neighbour) and SVM (support vector machine).
- 1.3 The third objective is to identify significant variables or features and their thresholds for classification.

2. Description of Data

2.1. Data Source, Size, Shape

2.1.1. Data Source –

https://www.kaggle.com/datasets/syedanwarafridi/vehicle-sales-data

- 2.1.2. Data Size (in KB | MB | GB ...) **88 MB**
- 2.1.3. Data Shape | Dimension:

Number of Variables - 16 Number of Records - 558837

2.2. Description of Variables

2.2.1. Index Variable(s): Car Id

- 2.2.2. Variables or Features having Categories | Categorical Variables or Features (CV)
 - 2.2.2.1. Variables or Features having Nominal Categories | Categorical Variables or Features **Nominal Type**: make, model, trim, body, transmission, state, colour, interior, seller
 - 2.2.2.2. Variables or Features having Ordinal Categories | Categorical Variables or Features **Ordinal Type:** Condition
 - 2.2.3. Non-Categorical Variables or Features: vin, odometer, mmr, selling price, sale date

Car ID: Unique identifier for each car

Year: Numeric representation of manufacturing year

Make: Brand or manufacturer of the car Model: Specific model name of the car Trim: Variant or version of the model

Body: Type of body style (e.g., sedan, SUV)

Transmission: Type of transmission system (e.g., automatic, manual)

VIN: Vehicle Identification Number, unique to each car

State: State where the car is located

Condition: Condition of the car, possibly ordinal categorical data

Odometer: Numeric representation of mileage

Color: Color of the car

Interior: Color or material of the interior

Seller: Entity selling the car

MMR: Market value of the car, likely non-categorical data

Selling Price: Price at which the car is sold

Sale Date: Date and time of sale

2.3. Descriptive Statistics

2.3.1. Descriptive Statistics of Outcome Categorical Variables

It provides the statistics of cluster variable (categorical variable) by giving frequency as well as relative frequency (in %).

Row ID	[count	D Relativ
duster_0	26011	23.273
duster_1	59650	53.37
duster_2	26106	23.358

- 2.3.2. Descriptive Statistics: Categorical Variables or Features
- 2.3.2.1. Count | Frequency Statistics

Color

Row ID	■ count
black	22203
white	21649
silver	16729
gray	16352
blue	10163

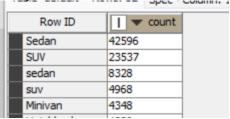
Model

Row ID	
Altima	6063
F-150	2992
Fusion	2604
Camry	2460
Escape	2247

Make

Row ID	I ▼ count					
Ford	20837					
Chevrolet	12069					
Nissan	10809					
Toyota	8033					
Dodge	6191					

Body



Transmission

Row ID	[count					
Sedan	2					
automatic	108246					
manual	3514					
sedan	5					

2.3.3 Descriptive Statistics: Non-Categorical Variables or Features

2.3.3.1. Measures of Central Tendency

Row ID	S Column	D Min	D Max	D Mean	D Std. devi	D Variance	D Skewness	D Kurtosis	D Overall s	No. missi	No. NaNs	No. +∞s	Nocos	D Median	Row count	£
condition	condition	1	49	30.574	13.314	177.254	-0.83	-0.197	3,417,183.716	0	0	0	0	?	111767	1
odometer	odometer	1	999,999	68,363.626	53,249.21	2,835,478,413	1.802	12.954	7,640,797,387	0	0	0	0	7	111767	1
mmr	mmr	25	178,000	13,782.935	9,718.146	94,442,361.104	2.026	11.693	1,540,477,346	0	0	0	0	?	111767	25
sellingprice	sellingprice	1	171,500	13,626.721	9,787.374	95,792,682.357	1.959	10.763	1,523,017,736	0	0	0	0	7	111767	1

2.3.3.2. Measures of Dispersion

Statistics Rows: 4 Co	lumns: 12								Q
Name	Туре	# Missing val	# Unique val	Minimum	Maximum	25% Quantile	50% Quantile	75% Quantile	Standard 7
condition	Number (dou	0	42	1	49	24	34	41	13.314
odometer	Number (dou	0	78138	1	999,999	28,408	52,407	99,088	53,249.21
mmr	Number (dou	0	1066	25	178,000	7,100	12,250	18,350	9,718.146
sellingprice	Number (dou	0	1222	1	171,500	6,900	12,100	18,250	9,787.374

Source of data-

https://www.kaggle.com/datasets/syedanwarafridi/vehicle-sales-data

3. Analysis of Data

3.1. Data Pre-Processing

3.1.1. Missing Data Statistics and Treatment

- 3.1.1.1. Missing Data Statistics: 16
- 3.1.1.1.2. Missing Data Treatment: make, model, trim, body, transmission, state, colour, interior, seller, condition, vin, odometer, mmr, selling price, sale date
 - 3.1.1.1.2.1. Removal of Records with More Than 50% Missing Data
 - 3.1.1.2.1. Missing Data Statistics: Categorical Variables or Features

Name	# Missing values
year	0
make	2141
model	2170
trim	2203
body	2688
transmission	13241
state	0
color	163
interior	163
seller	0

3.1.1.2.2. Missing Data Treatment: Categorical Variables or Features - 10 3.1.1.2.2.1. Removal of Variables or Features with More Than 50% Missing Data:

make, model, trim, body, transmission, state, colour, interior, seller, condition 3.1.1.2.2.2. Imputation of Missing Data using Descriptive Statistics: Mode

3.1.1.3.1. Missing Data Statistics: Non-Categorical Variables or Features

Name	# Missing values
vin	2
condition	2342
odometer	21
mmr	9
sellingprice	2
saledate	2

- 3.1.1.3.2. Missing Data Treatment: Non-Categorical Variables or Features 6
 - 3.1.1.3.2.1. Removal of Variables or Features with More Than 50% Missing Data: vin, odometer, mmr, selling price, sale date
 - 3.1.1.3.2.2. Imputation of Missing Data using Descriptive Statistics: Mean

3.1.2. Numerical Encoding of Categorical Variables or Features (Encoding Schema

- Alphanumeric Order)
- In this case, category to number node will be used to encode the categorical variables.

Color-

- 8 black
- 9 blue
- 14 gray
- 22 silver
- 24 white

Model

- 30-Altima
- 91- F-150
- 90- Fusion
- 62- Camry
- 75- Escape

Make

- 19-Ford
- 0-Chevrolet
- 5-Nissan
- 17-Toyota
- 22-Dodge

Body

- 0- Sedan
- 1-SUV
- 28-sedan
- 44-suv
- 9-Minivan

Transmission

- 0 Sedan
- 1- Automatic
- 2 Manual
- 3 sedan

3.1.3. Outlier Statistics and Treatment (Scaling | Transformation)

3.1.3.1.1. Outlier Statistics: Non-Categorical Variables or Features

Row ID	S Outlier	Membe	Outlier	D Lower	D Upper
Row0	condition	111767	0	-1.5	66.5
Row1	odometer	111767	2066	-77,611	205,105
Row2	mmr	111767	3244	-9,775	35,225
Row3	sellingprice	111767	3222	-10,125	35,275

- 3.1.3.1.2. Outlier Treatment: Non-Categorical Variables or Features
 - 3.1.3.1.2.1. Standardization
 - 3.1.3.1.2.2. Normalization using Min-Max Scaler:

Min-max normalization, also known as feature scaling, is a technique

used in data preprocessing to scale numerical features to a specific range, typically between 0 and 1.

The formula for min-max normalization is:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

3.1.3.1.2.3. Log Transformation

3.1.4. Data Bifurcation: Training & Testing Sets

The training and testing data have been bifurcated into 70% and 30% respectively.

3.2. Data Analysis

3.2.1. Supervised Machine Learning Classification Algorithm: Decision Tree

- → A decision tree is a supervised machine learning algorithm used for both classification and regression tasks. It works by recursively partitioning the input space into smaller regions based on feature values, creating a tree-like structure of decisions. At each node of the tree a decision is made based on the value of a specific feature, and the data is split into subsets. This process continues until a stopping criterion is met, such as reaching a maximum depth or no further improvement in impurity reduction.
- → In this project, decision tree will be the classification algorithm used for unsupervised learning. The metrics used in decision tree is Gini coefficient.
- → When using decision tree, we will be also seeing comparison when no pruning method is used and when pruning method is used.

3.2.2. Supervised Machine Learning Classification: Other Methods

Logistic Regression

It is a supervised learning algorithm used for binary classification tasks. It models the probability of the input belonging to a particular class using the logistic function. The algorithm learns the relationship between input features and the probability of the binary outcome, making it suitable for predicting categorical outcomes.

In this project, logistic regression will be used and the metric used in logistic regression is iteratively reweighted least squares (solver method).

K-Nearest Neighbours

K-Nearest Neighbours (KNN) is a supervised learning algorithm that is also used for both classification and regression tasks. It predicts the classification of a data point by finding the majority class among its k nearest neighbours in the feature space. KNN's performance heavily depends on the choice of distance metric and the value of k, making it sensitive to the dataset's characteristics.

In this project, KNN will be used and the metric used is Euclidean distance. For comparison, we will be using k = 7 till k = 19 in steps of 2 i.e. k = 7,9,11,13,15,17 and 19.

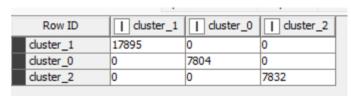
Support Vector Machines

Support Vector Machine (SVM) is a powerful supervised learning algorithm used for classification and regression tasks. It works by finding the hyperplane that best separates the classes in the feature space, maximizing the margin between them. SVM can handle high-dimensional data and is effective even in cases where the number of features exceeds the number of samples.

In this project, the kernel used will be polynomial and the parameters are power = 1, bias = 1 and gamma = 1.

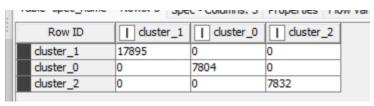
3.2.2.1. Classification Model Performance Evaluation of Decision Tree by using Confusion Matrix

Without Pruning



'	DIE GETOUR - NOW - T Spec - Columns: 11 Properties Flow Variables												
	Row ID	TruePo	FalsePo	TrueNe	FalseN	D Recall	D Precision	D Sensitivity	D Specificity	D F-meas	D Accuracy	D Cohen'	
П	cluster_1	17895	0	15636	0	1	1	1	1	1	?	?	
Ш	cluster_0	7804	0	25727	0	1	1	1	1	1	?	?	
ш	cluster_2	7832	0	25699	0	1	1	1	1	1	?	?	
	Overall	?	?	?	?	?	?	?	?	?	1	1	

With Pruning



Row ID	TruePo	FalsePo	TrueNe	FalseN	D Recall	D Precision	D Sensitivity	D Specificity	D F-meas	D Accuracy	D Cohen'
duster_1	17895	0	15636	0	1	1	1	1	1	?	?
duster_0	7804	0	25727	0	1	1	1	1	1	?	?
cluster_2	7832	0	25699	0	1	1	1	1	1	?	?
Overall	?	?	?	?	?	?	?	?	?	1	1

Cluster 0

- This cluster has a high number of true positives and true negatives indicating that the model correctly classified most instances within this cluster.
- The precision and recall scores are both very high suggesting that the model effectively identifies true positives while also minimizing false positives.

Cluster 1

- This cluster has a lower recall and precision compared to cluster 0, indicating that the model's performance is not as strong for this segment.
- The number of false positives is relatively high, suggesting that the model may misclassify some instances within this cluster.

• Despite the lower performance metrics, the specificity is very high indicating that the model correctly identifies true negatives within this cluster.

Cluster 2

- This cluster has a relatively high recall and precision, indicating that the model performs well.
- The number of false positives is relatively low suggesting that the model effectively minimizes misclassifications within this cluster.
- Both sensitivity and specificity scores are high indicating that the model correctly identifies both true positives and true negatives within this cluster.

Comparative analysis of decision tree with and without pruning

Pruning generally improves precision and specificity while slightly reducing recall and sensitivity. Pruning removes unnecessary branches from the tree, simplifying the model and reducing overfitting. This can lead to better generalization and potentially improved performance on unseen data.

We didn't observe a significant difference between a pruned and non-pruned decision tree in our case. It may be because:

- 1. **Dataset characteristics:** The data we used might be relatively simple, and the decision tree without pruning may not have overfit considerably.
- 2. **Pruning settings:** The pruning settings in KNIME's Decision Tree Learner node might have been configured in a way that resulted in minimal removal of branches.
- 3. **Randomness:** There can be an element of randomness in decision tree generation. Rerunning the experiment with both pruned and non-pruned trees might yield a slight difference on another iteration.

The choice of whether to prune the decision tree depends on the specific requirements of the problem and the trade-off between precision and recall. If minimizing false positives is crucial (can be used for risk assessment) pruning may be preferred. If capturing as many true positives as possible is more important (can be used for customer retention) pruning may be avoided.

3.2.2.2. Classification Model Performance Evaluation of Other Supervised Learning methods by using confusion matrix

Logistic Regression

Tubic Goetileen		Spec - Columns.	0
Row ID	S 🗻 Logit	S Variable	
Row1	cluster_0	year	
Row2	cluster_0	transmission=automatic	
Row3	cluster_0	transmission=manual	
Row4	cluster_0	transmission=sedan	
Row5	cluster_0	state=3vwd17aj4fm236636	
Row6	cluster_0	state=3vwd17aj5fm219943	
Row7	cluster_0	state=3vwd17aj5fm221322	
Row8	cluster_0	state=3vwd17aj5fm225953	
Row9	cluster_0	state=3vwd17aj7fm222388	
Row10	cluster_0	state=3vwd17aj7fm229552	
Row11	cluster_0	state=ab	
Row12	cluster_0	state=al	
Row13	cluster_0	state=az	
Row14	cluster_0	state=ca	
Row15	duster_0	state=co	
Row90	cluster 0	color=silver	-0.003
Row91	cluster_0	color=turquoise	0.052
Row92	cluster_0	color=white	0.3
Row93	cluster_0	color=yellow	0.219
Row94	cluster_0	color=â€"	-0.368
Row95	cluster_0	interior=black	-0.159
Row96	cluster_0	interior=blue	0.19
Row97	cluster_0	interior=brown	-0.292
Row98	cluster_0	interior=burgundy	0.009
Row99	cluster_0	interior=gold	-0.112
Row 100	duster_0	interior=gray	-0.041
Row 101	duster_0	interior=green	-0.074
Row 102	duster_0	interior=off-white	-0.087
Row 103	duster_0	interior=orange	0.033
Row 104	duster_0	interior=purple	0.076
Row 105	duster_0	interior=red	0.082
Row 106	duster_0	interior=silver	0.169
Row 107	duster_0	interior=tan	-0.047
Row 107	duster_0	interior=tan interior=white	0.033
Row 109	duster_0	interior=write interior=yellow	0.001
Row109		· ·	-0.269
Row111	duster_0	interior=â€"	-0.269
Row112	duster_0	Car id (to number) make (to number)	-0.001
Row112	duster_0		-
	duster_0	model (to number)	0
Row114	duster_0	trim (to number)	0 032
Row115	duster_0	body (to number)	-0.032
Row116	duster_0	transmission (to number)	-0.076
Row117	duster_0	state (to number)	-0.02
Row118	duster_0	color (to number)	0.049
Row119	cluster_0	interior (to number)	-0.031
Row120	cluster_0	seller (to number)	0
Row121	cluster_0	Clusters (to number)	-72.932
Row122	cluster_0	odometer	0.625
Row123	cluster_0	mmr	-0.346
Row124	cluster 0	sellinaprice	0.312

			color –r cu		
	Row215	duster_1	color=silver		-0.018
	Row216	duster_1	color=turquoise		-0.013
	Row217	cluster_1	color=white		0.694
	Row218	duster_1	color=yellow		0.289
	Row219	duster_1	color=—		-1.035
	Row220	duster_1	interior=black		-0.671
	Row221	duster_1	interior=blue		0.432
	Row222	duster_1	interior=brown		-0.99
	Row223	duster_1	interior=burgundy		-0.04
	Row224	cluster_1	interior=gold		0.107
	Row225	cluster_1	interior=gray		-0.546
	Row226	duster_1	interior=green		-0.007
	Row227	cluster_1	interior=off-white		-0.055
	Row228	duster_1	interior=orange		-0.119
	Row229	cluster_1	interior=purple		0.06
	Row230	cluster_1	interior=red		0.147
	Row231	cluster_1	interior=silver		-0.093
	Row232	cluster_1	interior=tan		0.213
	Row233		interior=white		0.246
	Row234		interior=yellow		0.032
	Row235		interior=â€″		-0.645
	Row236		Car id (to number)		-0
	Row237		make (to number)		-0.015
	Row238		model (to number)		-0
	Row239	duster_1	trim (to number)		0
	Row240		body (to number)		0.058
	Row241	cluster_1	transmission (to number)		-0.187
	Row242		state (to number)		0.008
	Row243		color (to number)		0.152
	Row244 Row245		interior (to number)		-0.13 0
	Row245		seller (to number)		-149.664
	Row247		Clusters (to number) odometer		0.943
	Row248		mmr		-0.192
	Row249		sellingprice		-0.132
	NOW 2 13	ciustei_1	sellingprice		0.21
	Row141	cluster_1	state=fl	-0.13	38
	Row142	cluster_1	state=ga	-0.06	57
	Row143	cluster_1	state=hi	-0.04	1 6
	Row144	cluster_1	state=il	0.10	4
	Row145	cluster_1	state=in	0.47	4
L	Row146	cluster_1	state=la	0.23	4
L	Row147	cluster_1	state=ma	0.04	6
	Row148	cluster_1	state=md	-0.6	19
	Row149	cluster_1	state=mi	-0.30	
	Row150	cluster_1	state=mn	-0.4	
	Row151	cluster_1	state=mo	0.27	
	Row152	cluster_1	state=ms	0.15	
	Row153	cluster_1	state=nc	-0.16	
	Row154	cluster_1	state=ne	-0.4	
	Row155	cluster_1	state=nj	0.68	

Cluster_2 was used as the reference category

Cluster_0 which represent car with maker Ford, model Altima and sedan body

Identity of cluster 0: Customers who value reliability, affordability, comfort, fuel efficiency, and practicality.

We observe that state, color, interior and transmission are most significant variables in cluster 0.

Cluster 1 which represent car with maker Honda, model Camry and SUV body

Identity of cluster 1: Customers who are family-oriented, seeking vehicles that offer ample space and versatility for various activities and lifestyles, also have interest in features that enhance convenience and comfort.

We observe that state, color and interior are most significant variables in cluster 0.

Variables like make, body, trim, selling price, model have no significant impact in distinguishing cluster 1 and cluster 0 from cluster 2.

Table detault - P	: uerauli "numa- 1 Spec - Columns: 11 Properties Flow Variables										
Row ID	TruePo	FalsePo	TrueNe	FalseN	D Recall	D Precision	D Sensitivity	D Specificity	D F-meas	D Accuracy	D Cohen'
duster_1	17895	0	15636	0	1	1	1	1	1	?	?
duster_0	7804	0	25727	0	1	1	1	1	1	?	?
duster_2	7832	0	25699	0	1	1	1	1	1	?	?
Overall	?	?	?	?	?	?	?	?	?	1	1

The overall accuracy of the logistic regression model is very high at 100 and it effectively predicts the cluster labels for the majority of instances. Additionally, the Cohen's Kappa coefficient suggests substantial agreement beyond chance among the predicted and actual cluster labels.

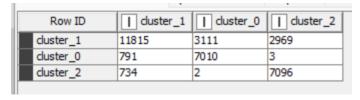
K-Nearest Neighbour

K=7

Row ID	duster_1	duster_0	duster_2
cluster_1	11970	3029	2896
cluster_0	981	6820	3
cluster_2	927	2	6903

Row ID	TruePo	FalsePo	TrueNe	FalseN	D Recall	D Precision	D Sensitivity	D Specificity	D F-meas	D Accuracy	D Cohen'
cluster_1	11970	1908	13728	5925	0.669	0.863	0.669	0.878	0.753	?	?
cluster_0	6820	3031	22696	984	0.874	0.692	0.874	0.882	0.773	?	?
cluster_2	6903	2899	22800	929	0.881	0.704	0.881	0.887	0.783	?	?
Overall	?	?	?	?	?	?	?	?	?	0.766	0.636

K=9



Row ID	TruePo	FalsePo	TrueNe	FalseN	D Recall	D Precision	D Sensitivity	D Specificity	D F-meas	D Accuracy	D Cohen'
cluster_1	11815	1525	14111	6080	0.66	0.886	0.66	0.902	0.757	?	?
cluster_0	7010	3113	22614	794	0.898	0.692	0.898	0.879	0.782	?	?
duster_2	7096	2972	22727	736	0.906	0.705	0.906	0.884	0.793	?	?
Overall	?	?	?	?	?	?	?	?	?	0.773	0.649

K=19

Row ID	duster_1	duster_0	duster_2
cluster_1	11347	3337	3211
cluster_0	274	7528	2
cluster_2	303	3	7526

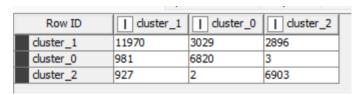
	done delidate ito	or . Spec - c	Joidinis, 11 Fi	oper des 1 low	variables							
	Row ID	TruePo	FalsePo	TrueNe	FalseN	D Recall	D Precision	D Sensitivity	D Specificity	D F-meas	D Accuracy	D Cohen'
П	duster_1	11347	577	15059	6548	0.634	0.952	0.634	0.963	0.761	?	?
П	cluster_0	7528	3340	22387	276	0.965	0.693	0.965	0.87	0.806	?	?
П	cluster_2	7526	3213	22486	306	0.961	0.701	0.961	0.875	0.811	?	?
	Overall	?	?	?	?	?	?	?	?	?	0.787	0.678

Similarly, we have applied k nearest neighbour for K=11, 13,15,17 and observed that –

In KNN, the number of neighbours to be considered are from k=7 to 19. From the images, it is seen that as the number of k increases the accuracy also increases. For k=19, as the accuracy is the highest from all the other k's, this cluster will be considered.

The overall accuracy of the KNN model is moderate showing mixed performance across different clusters. Cohen's Kappa coefficient also suggests moderate agreement beyond chance among the predicted and actual cluster labels.

Support Vector Machines



	opec c		operaco	- GI IGOICO							
Row ID	TruePo	FalsePo	TrueNe	FalseN	D Recall	D Precision	D Sensitivity	D Specificity	D F-meas	D Accuracy	D Cohen'
cluster_1	11970	1908	13728	5925	0.669	0.863	0.669	0.878	0.753	?	?
duster_0	6820	3031	22696	984	0.874	0.692	0.874	0.882	0.773	?	?
duster_2	6903	2899	22800	929	0.881	0.704	0.881	0.887	0.783	?	?
Overall	?	?	?	?	?	?	?	?	?	0.766	0.636

The overall performance of the SVM model is very poor with extremely low recall, precision and accuracy metrics.

3.2.3.1. Variable or Feature Analysis for Decision Tree

3.2.3.1.1. List of Relevant or Important Variables.

In the decision tree analysis, we see that these were the important variables that contributed in the supervised learning algorithm which are: -

Transmission, color, interior, state, make, body and model.

3.2.3.1.2. List of Non-Relevant or Non-Important Variables

In the decision tree analysis, we see that these were the non-important variables that did not contribute in the supervised learning algorithm which are: -

Car id, odometer, vin, saledate, selling price, condition and trim.

3.2.3.2. Variable or Feature Analysis for Logistic Regression, K-Nearest Neighbour and Support Vector Machine

3.2.3.2.1. List of Relevant Variables

We have observed that state, color, interior and transmission are most significant variables in cluster 0.

3.2.3.2.2. List of Non-Important Variables

Car_id, odometer, vin, saledate, condition and trim are insignificant variables.

The above variables have value of p>0.05 which suggests potentially negligible impact on loan outcomes.

4. Results and Observations

4.1. Comparing Supervised Learning models: Decision Tree VS Logistic Regression, KNN and SVM

Decision tree

File Hilite				
Clusters \	duster_1	cluster_0	cluster_2	
duster_1	17895	0	0	
cluster_0	0	7804	0	
cluster_2	0	0	7832	
Corre	ect classified: 3	3,531		Wrong dassified: 0
	Accuracy: 1009	%		Error: 0%
Coh	nen's kappa (κ):	1%		

Logistic Regression

Clusters \	cluster_1	cluster_0	cluster_2	
cluster_1	11930	0	0	
cluster_0	0	5202	0	
cluster_2	0	0	5222	
Corre	ect classified: 2	2,354		Wrong dassified: 0
	Accuracy: 1009	%		Error: 0%

KNN

K=19

Clusters \	cluster_1	cluster_0	cluster_2	
cluster_1	11347	3337	3211	
cluster_0	274	7528	2	
cluster_2	303	3	7526	
Corre	ect classified: 2	26,401	Wr	ong classified: 7,130
Accuracy: 78.736%				Error: 21.264%
Cohen	's kappa (κ):	0.678%		

SVM

Clusters \	cluster_1	cluster_0	cluster_2		
cluster_1	11970	3029	2896		
cluster_0	981	6820	3		
duster_2	927	2	6903		
Correct classified: 25,693			Wr	ong classified: 7,838	
Acquiracy: 76 625%			Frror: 23 375%		

Cohen's kappa (κ): 0.636%

5. Managerial Insights

5.1. Appropriate Model

Metrics	Decision Tree	Logistic	KNN	SVM
Metrics	1166	Regression	121111	S V IVI
Accuracy				
(in %)	100%	100%	78.74%	76.63%

The decision tree and logistic regression has the highest accuracy (100%). KNN and SVM have significantly lower accuracies of 78.74% and 76.63% respectively.

Decision tree provides the highest accuracy of all the models according to the data and will be the appropriate model for the customer classification. Decision tree is able to handle both numerical and categorical which does benefit in this data as the data contains a combination of variables which are categorical and continuous in nature.

Managerial insights according to the appropriate model (Decision Tree)

Managerial insights according to the Decision Tree model for car prices:

Market differentiation: The model could identify distinct customer segments based on preferences for body style (sedan vs. SUV) and manufacturer (Ford vs. Honda vs. Chevrolet). This can inform targeted marketing campaigns for each cluster.

Pricing strategy: The decision tree can help establish price ranges based on the clusters. For instance, Ford Altima sedans (cluster 0) might have a different pricing strategy compared to Honda Camry SUVs (cluster 1).

Inventory management: The model can help predict demand for specific car types (clusters). This can guide decisions on stocking the right car models and trims to meet customer preferences.

Sales strategy: Insights from the model can help tailor sales approaches to different customer segments. For example, targeting features like fuel efficiency for Honda CR-V buyers (if present in the data) might be more effective than for Chevrolet F-150 buyers (cluster 2) who might prioritize power.

5.2. Relevant or Important Variables or Features

The relevant or important variables that are used in the decision tree supervised learning algorithm are: -

State	
Color	
Interior	
Transmission	
Selling Price	
Make	
Model	