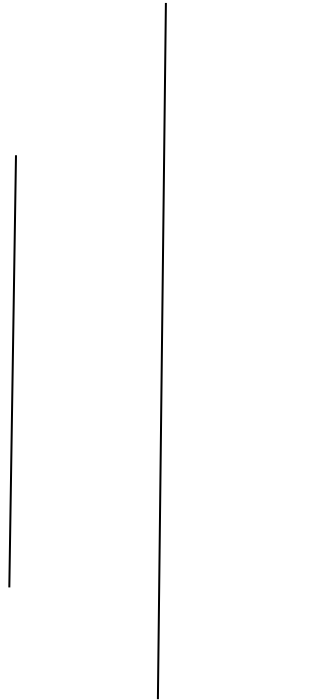


Data Mining Report on Predicting Vehicle Sales Based on Fuel Price



Cape Breton University

Course Number: 2024 Winter Data Mining (MGSC-5126-23)

Submitted by:

Md Sahid Parvez (Student Id:)

Submitted To:

Professor Samiul Islam

1. Introduction:

The relationship between gasoline prices and new motor vehicle sales is critical to understand in current energy situation. Canada committed to cut the emissions of greenhouse gas under the Paris Agreement of 2015, effective laws focusing on vehicle emissions are vital, particularly in transportation industry as it accounts for a large portion of Canada emissions profile. That's why It is important to understand the correlation between fuel prices and motor vehicle sales for reaching the emission reduction targets.

The variations in the fuel prices have in the past affected the buyer habits and choice making for new vehicles. Variations in fuel prices cause the consumers to constantly adjust their behaviour towards buying the new vehicles. Therefore, it is crucial to examine the current rules and developing the policies to decrease the emission from motor vehicles.

This report evaluates the relationship between the changing fuel prices and the customer behaviours when buying a new vehicle. This report explores the motor vehicle (Passenger cars and Trucks) sales data from 2000 to 2022 in different provinces of Canada to figure out how changes in fuel price influence consumers buying decisions for new motor vehicle.

Several studies have underscored a robust correlation between gasoline prices and market shares, particularly during periods of high or escalating fuel costs. For instance, from 2003 to 2007, the surge in gasoline prices elucidated approximately half of the transition from large sport utility vehicles (SUVs) to smaller crossovers. However, recent years have witnessed elevated and relatively volatile gas prices. My findings suggest that new vehicle sales exhibit greater responsiveness to ascending prices as opposed to descending prices. Nonetheless, the precise reasons behind the reduced impact of falling gas prices remain an unresolved query warranting further investigation.

The data used for this report has been sourced from a Canadian statistics website. Specifically, I extracted data pertaining to new motor vehicle sales and monthly average retail prices for gasoline from the website. Both datasets contained numerous unnecessary columns for the project's purposes and missing values. Therefore, I refined the data by editing it and then merged both datasets to align with the requirements of the project. For refining of the datasets, I used Microsoft Power Bi tool.

The first dataset contains information about the year, month, geography, vehicle type, and number of sales from the years 2001 to 2022. The table below displays the first five records of new motor vehicle sales data.

Table 1: New Motor Vehicle Sales

	Year	Month	GEO	Vehicle type	Number Of Sales	Fuel Price(In Cents)
0	2001	January	Alberta	Passenger cars	4263	72.6
1	2001	January	Alberta	Trucks	8102	72.6
2	2001	January	British Columbia	Passenger cars	5005	78.9
3	2001	January	British Columbia	Trucks	5595	78.9
4	2001	January	Manitoba	Passenger cars	1106	72.8

The second dataset contains information about the year, month, geography, and fuel price (in cents) from the years 2001 to 2022. The table below displays the first five records of new Monthly average fuel price data.

Table 2: Monthly average retail prices for gasoline

	Year	Month	GEO	Fuel Price(In Cents)
0	2001	April	Alberta	71.0
1	2001	April	Alberta	71.0
2	2001	April	British Columbia	80.8
3	2001	April	British Columbia	80.8
4	2001	April	Manitoba	73.6

The final dataset, after cleaning and merging, is shown in the table below, containing all the required information.

Table 3: Merged table

	Year	Month	GEO	Vehicle type	Number Of Sales	Fuel Price(In Cents)
0	2001	January	Alberta	Passenger cars	4263	71.0
1	2001	January	Alberta	Trucks	8102	71.0
2	2001	January	British Columbia	Passenger cars	5005	80.8
3	2001	January	British Columbia	Trucks	5595	80.8
4	2001	January	Manitoba	Passenger cars	1106	73.6

Table 4: Summary Statistics

	Year	Number Of Sales	Fuel Price(In Cents)
count	5280.00000	5280.000000	5280.000000
mean	2011.50000	7222.840152	119.362955
std	6.34489	9807.783236	27.955777
min	2001.00000	44.000000	56.400000
25%	2006.00000	1203.000000	100.400000
50%	2011.50000	2436.000000	118.000000
75%	2017.00000	10066.250000	136.325000
max	2022.00000	63433.000000	249.100000

The average number of motor vehicle sales is 7,222, and the average fuel price is \$119.36 for the entire period in the data. The maximum number of vehicle sales is 63,433 in a year, and the fuel price reached a maximum of 249.1 cents from 2001 to 2022.

Overall, this report contributes to the existing literature by offering detailed insights into the interplay between fuel prices and consumer behaviour in the Canadian automotive market. The comprehensive analysis spanning model years 2000 through 2023 captures various economic conditions, including periods of fluctuating gasoline prices and significant shifts in market dynamics, thus providing valuable implications for policymakers, industry stakeholders, and researchers interested in promoting sustainable transportation practices.

2. Literature review:

Answering the question "How do gasoline prices affect gasoline usage?" is complex and multifaceted, as there are numerous factors and adjustments that come into play across various timeframes. Drivers, car buyers, and automobile manufacturers have several options to adapt to changes in gasoline prices, each of which influences gasoline usage differently.

In the short term, drivers can quickly adjust their driving habits in response to fluctuations in gasoline prices. Studies by Donna (2010) and Goldberg (1998) examine how changes in gasoline prices impact public transportation usage and vehicle miles travelled, respectively, shedding light on immediate driving responses.

Conversely, in the long run, automobile manufacturers can modify the fuel efficiency of vehicles by altering characteristics such as weight, power, and combustion technology, or by transitioning to hybrid or electric models. Gramlich (2009) explores these manufacturer responses by analysing changes in the MPG (miles per gallon) of car models over time in relation to gasoline prices.

This paper is akin to a group of studies that address an intermediate horizon question: How do gasoline prices influence the prices or sales of car models? These studies investigate how consumer choices among available car models are affected by gasoline prices, examining outcomes such as prices, sales, and market shares. Some focus on the impact of gasoline prices on car quantities, contributing to the understanding of fleet fuel economy, while others explore the effects on car prices, which relates to consumer behaviour and decision-making.

Two notable studies examining the relationship between gasoline prices and car quantities are Klier and Linn (forthcoming) and Li, Timmins, and von Haefen (2009). While both papers investigate similar questions, they utilize different datasets. Klier and Linn analyse the impact of national average gasoline prices on the sales of new cars at a detailed model level, finding that higher gasoline prices lead to reduced sales of low-MPG cars compared to high-MPG cars. In contrast, Li, Timmins, and von Haefen combine data on new car sales with vehicle registrations to assess the effect of gasoline prices on both inflow and outflow from the vehicle fleet. They observe differential effects based on car fuel economies, with increased gasoline prices associated with higher sales of fuel-efficient new cars and improved survival probabilities for fuel-efficient used cars, while sales of fuel-inefficient cars decline.

Several studies investigate the relationship between car prices and gasoline prices to understand if car buyers exhibit myopia regarding future usage costs. Kahn (1986) and Kilian and Sims (2006) find evidence of some degree of myopia, with used car prices adjusting partially to changes in gasoline prices. Allcott and Wozny (2010) and Sallee, West, and Fan (2009) similarly find that car buyers undervalue fuel costs, although to varying degrees. However, Goldberg (1998) concludes that car buyers are not myopic, as they adjust demand for cars in response to both purchase prices and fuel costs.

Our paper diverges from previous literature in two key aspects. Firstly, our dataset enables us to examine the effects of gasoline prices on both car prices and quantities, across both new markets, within a single data source. Secondly, we employ more flexible functional forms in our analysis compared to many existing studies.

3. Methodology:

Exploratory Data Analysis (EDA) and Predictive Analytics has been used for the methodology part.

3.1 EDA

Exploratory Data Analysis (EDA) is a process of describing the data by means of statistical and visualization techniques to bring important aspects of that data into focus for further analysis. This involves inspecting the dataset from many angles, describing & summarizing it without making any assumptions about its contents.

“Exploratory data analysis is an attitude, a state of flexibility, a willingness to look for those things that we believe are not there, as well as those we believe to be there” – John W. Tukey

I plotted the time series of average monthly sales for each vehicle type over the years to observe trends and seasonality.

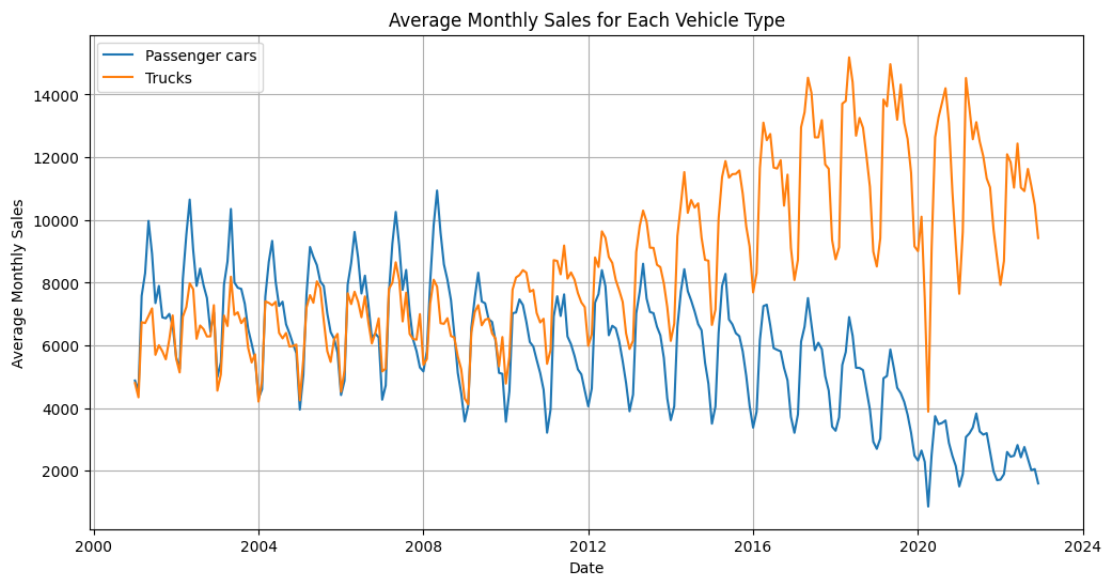


Figure 1: Average monthly sales for each vehicle types.

Interpretation: In the initial years, passenger cars enjoyed greater popularity, but over time, trucks have surpassed them in popularity. Both vehicle sales experienced a significant drop in the years 2020-2021, but they recovered afterward.

Plotting the time series of average monthly fuel prices over the years to observe trends and fluctuations.

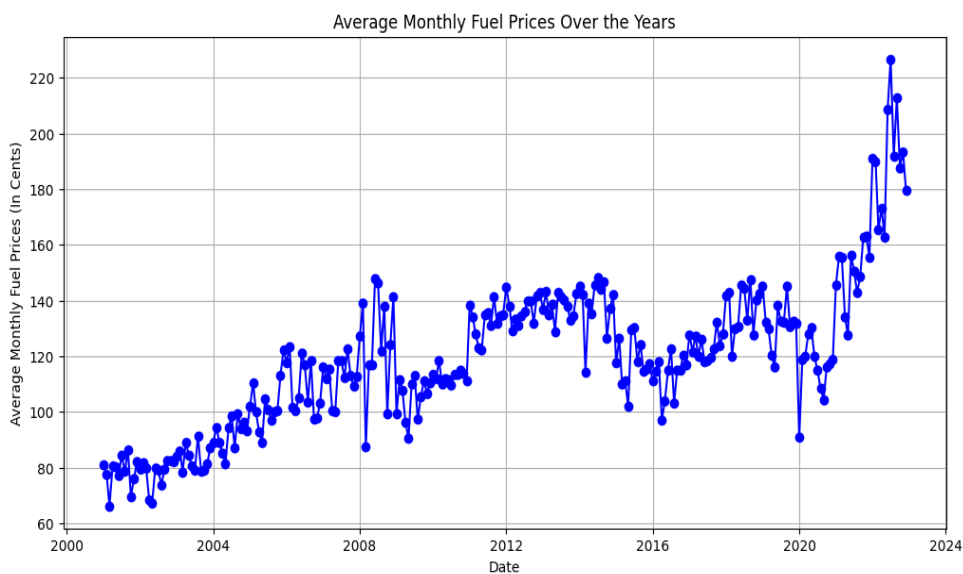


Figure 2: Average monthly fuel prices over the years.

Interpretation:

The average price of oil increased over the years, reaching its recent low in 2020, and then surged to its highest point in 2022. The highest average price increase occurred from 2020 to the end of 2022.

Now, exploring the correlation between monthly oil prices and average monthly sales for each vehicle type using scatter plots.

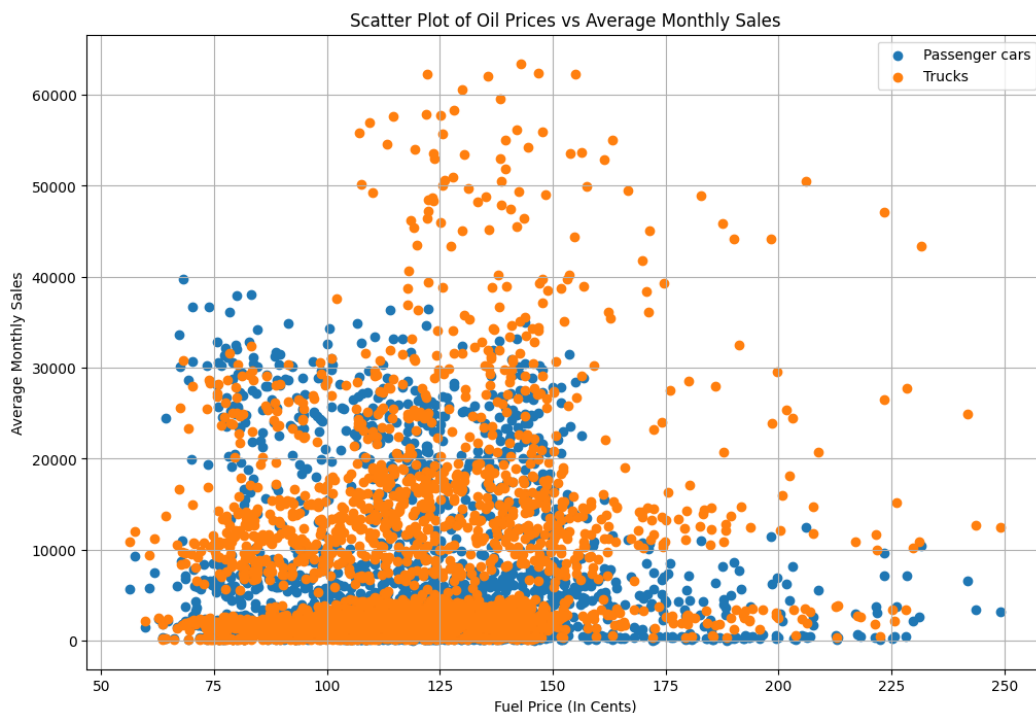


Figure 3: Scatter plot Oil Prices and Average Monthly Sales

Interpretation:

Before the fuel price reached 150 cents per litre, the sales of both motor vehicles were high. However, as the fuel price increased beyond 150 cents per litre, the average sales decreased. Additionally, a noticeable trend was observed: the average sales of trucks were higher than those of passenger cars as the fuel price increased beyond 150 cents per litre.

We explored the impact of oil prices by comparing sales trends before and after significant changes in oil prices using bar chart.

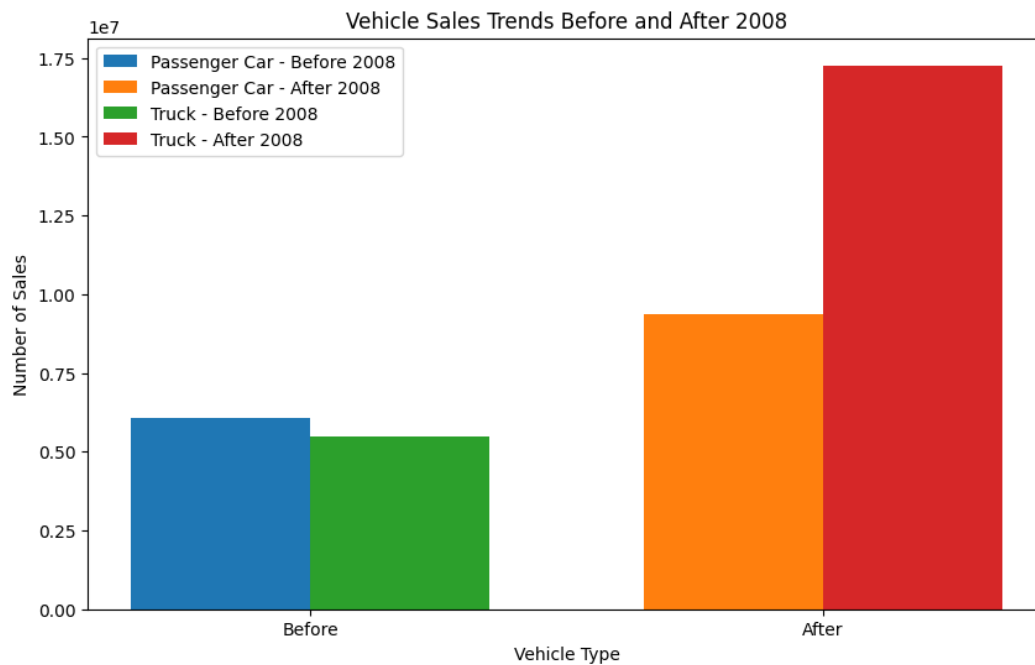


Figure 4: Vehicle Sales Trends Before and After 2008

Calculate correlation coefficients between oil prices and car sales for each vehicle type to quantify the strength and direction of the relationship.

Oil Prices and Vehicle Sales	Correlation coefficient
Passenger Cars	-0.07333
Trucks	0.173702

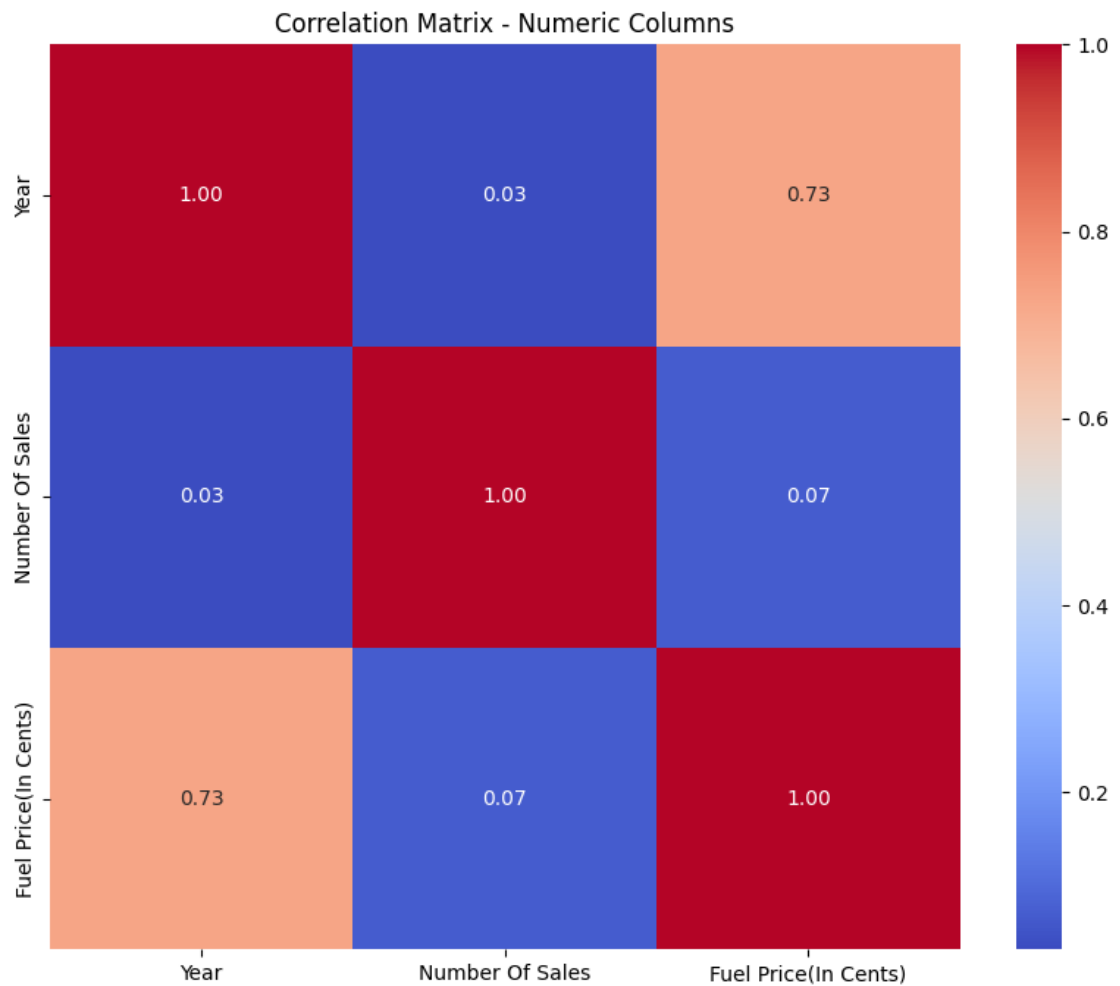
For Passenger cars:

The correlation coefficient between oil prices and car sales for passenger cars is approximately -0.073. This value is close to zero, indicating a very weak negative linear relationship between oil prices and car sales for passenger cars. In other words, there is almost no discernible pattern in the relationship between oil prices and passenger car sales.

For Trucks:

The correlation coefficient between oil prices and car sales for trucks is approximately 0.174. This value is also relatively close to zero, but slightly positive. It suggests a very weak positive linear relationship between oil prices and truck sales.

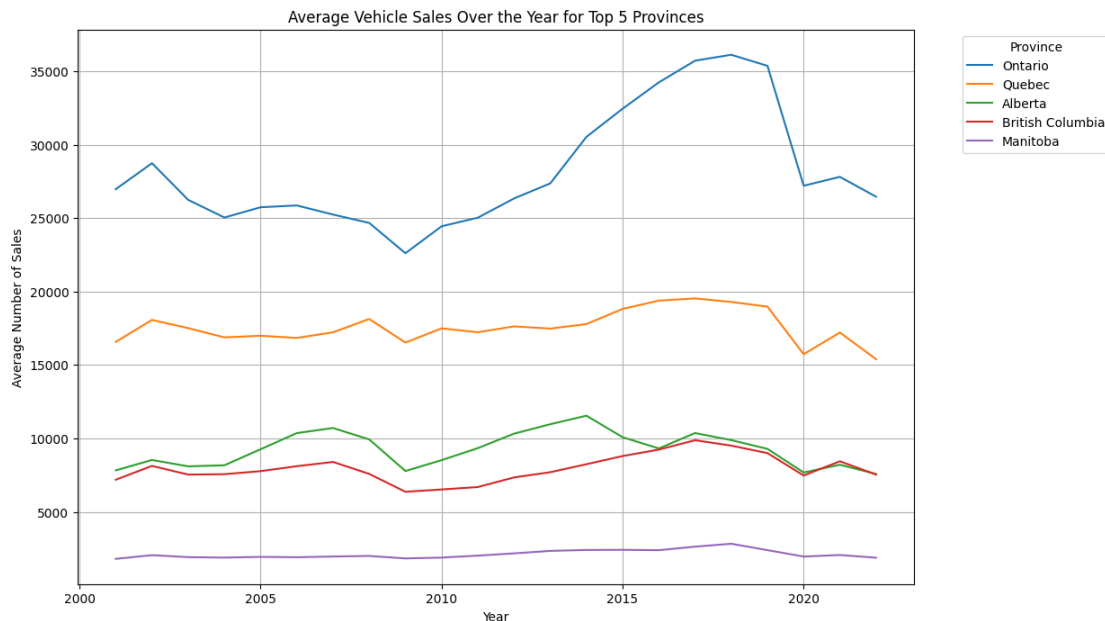
Created heatmap visually represent the correlations between oil prices and car sales for different vehicle types.



Interpretation:

Changes in the year or fuel price tend to significantly affect the vehicle sales as year and fuel price are highly correlated.

Created line graph to represent average vehicle sales over the year for top 5 provinces.



Interpretation:

Ontario is the province having the greatest number of vehicles sales with respect to other provinces.

3.2 Predictive Analysis:

Predictive analytics is the process of using data to forecast future outcomes. The process uses data analysis, machine learning, artificial intelligence, and statistical models to find patterns that might predict future behaviour. Organizations can use historic and current data to forecast trends and behaviours seconds, days, or years into the future with a great deal of precision.

Data scientists use predictive models to identify correlations between different elements in selected datasets. Once data collection is complete, a statistical model is formulated, trained, and modified to generate predictions.

```
[71] X = final_df[['Year', 'GEO', 'Vehicle type', 'Fuel Price(In Cents)', 'Month']]
      y = final_df['Number Of Sales']

      X = pd.get_dummies(X)

      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

      # Function to evaluate RMSE
      def evaluate_rmse(y_true, y_pred):
          rmse = np.sqrt(mean_squared_error(y_true, y_pred))
          return rmse

      # Function to evaluate MAPE
      def evaluate_mape(y_true, y_pred):
          mape = mean_absolute_percentage_error(y_true, y_pred)
          return mape

      future_data = [[2025,172.6,1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0,0,0,0,0,1]]
```

Figure 5: Code Snippet for train and test, dummies creation and future data

I used three different forecasting methods on the dataset to predict the Vehicle sales for future years. They are:

1. Linear regression
2. Gradient Boosting
3. Random Forest Regressor

Linear regression:

Linear regression analysis is used to predict the value of a variable based on the value of another variable. The variable you want to predict is called the dependent variable. The variable you are using to predict the other variable's value is called the independent variable.

This form of analysis estimates the coefficients of the linear equation, involving one or more independent variables that best predict the value of the dependent variable. Linear regression fits a straight line or surface that minimizes the discrepancies between predicted and actual output values. There are simple linear regression calculators that use a “least squares” method to discover the best-fit line for a set of paired data. You then estimate the value of X (dependent variable) from Y (independent variable).

Reason to use:

Providing simplicity and interpretability, Linear Regression is effective for understanding linear relationships and serving as a baseline model; although it may not capture complex patterns as effectively as ensemble methods, it offers fast training and prediction times.

```

# Linear Regression
start_time = time.time()
lr_model = LinearRegression()
lr_model.fit(X_train, y_train)
lr_train_time = time.time() - start_time

# Evaluate models
lr_predictions = lr_model.predict(X_test)
lr_rmse = evaluate_rmse(y_test, lr_predictions)
lr_mape = evaluate_mape(y_test, lr_predictions)

print("Linear Regression RMSE:", lr_rmse)
print("Linear Regression MAPE:", lr_mape)
print("Linear Regression Training Time:", lr_train_time)

start_time = time.time()
future_predictions_lr = lr_model.predict(future_data)
lr_prediction_time = time.time() - start_time

print("Linear Regression Prediction Time:", lr_prediction_time)
print("Future Linear Regression Prediction:", future_predictions_lr)

Linear Regression RMSE: 4210.590662507497
Linear Regression MAPE: 1.5107060332491677
Linear Regression Training Time: 0.01924610137939453
Linear Regression Prediction Time: 0.002560853958129883
Future Linear Regression Prediction: [4959.81363424]

```

Figure 6: Code Snippet for Linear Regression

Gradient Boosting:

Gradient Boosting is a powerful boosting algorithm that combines several weak learners into strong learners, in which each new model is trained to minimize the loss function such as mean squared error or cross-entropy of the previous model using gradient descent. This algorithm calculates the gradient of loss function based on the prediction of the current ensemble and then trains a new model to minimize the gradient at each iteration. The outcomes of the new models are delivered to the ensemble. This method is repeated till the required conditions are achieved.

This algorithm operates on tabular form of data having a collection of independent columns (X) and a target variable (y). The Target variable in this analysis is 'Number of Sales'. The purpose of this algorithm is to collect statistics enough from training data correctly predict new data points.

Reason to use:

This algorithm improves model predictions through iteration, making it best for figuring out the complicated patterns inside the data. However, because of its ensemble nature it may need greater computational power and tuning.

```

# Gradient Boosting
start_time = time.time()
gb_model = GradientBoostingRegressor(n_estimators=100, random_state=42)
gb_model.fit(X_train, y_train)
gb_train_time = time.time() - start_time

gb_predictions = gb_model.predict(X_test)
gb_rmse = evaluate_rmse(y_test, gb_predictions)
gb_mape = evaluate_mape(y_test, gb_predictions)

print("Gradient Boosting RMSE:", gb_rmse)
print("Gradient Boosting MAPE:", gb_mape)
print("Gradient Boosting Training Time:", gb_train_time)

start_time = time.time()
future_predictions_gb = gb_model.predict(future_data)
gb_prediction_time = time.time() - start_time

print("Gradient Boosting Prediction Time:", gb_prediction_time)
print("Future Gradient Boosting Prediction:", future_predictions_gb)

```

```

Gradient Boosting RMSE: 2022.519205288318
Gradient Boosting MAPE: 0.5132015692707662
Gradient Boosting Training Time: 0.38396549224853516
Gradient Boosting Prediction Time: 0.0006480216979980469
Future Gradient Boosting Prediction: [2396.47108651]

```

Figure 7: Code for Gradient Boosting

Random Forest Regressor:

Random Forest Algorithm is famous and is used widely as its nature is user friendly and flexibility. It may be used successfully for each Regression and Classification tasks. This algorithm can method the complex datasets and it can save you overfitting. It is powerful algorithm used for numerous predictive tasks.

Whether the dataset carries continuous variables or categorical variables, this algorithm is powerful enough to tackle each. Therefore, it was extensively used for Regression and Classification. It performs better for classification and regression tasks. In this tutorial, we will understand the working of random forest and implement random forest on a classification task.

Reason to use:

Offering robustness to overfitting and reliable predictions, Random Forest is well-suited for large datasets with diverse feature types; its interpretability may be limited compared to simpler models, and tuning hyperparameters could be necessary for optimal performance.

```

# Random Forest
start_time = time.time()
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)
rf_train_time = time.time() - start_time

rf_predictions = rf_model.predict(X_test)
rf_rmse = evaluate_rmse(y_test, rf_predictions)
rf_mape = evaluate_mape(y_test, rf_predictions)

print("Random Forest RMSE:", rf_rmse)
print("Random Forest MAPE:", rf_mape)
print("Random Forest Training Time:", rf_train_time)

start_time = time.time()
future_predictions_rf = rf_model.predict(future_data)
rf_prediction_time = time.time() - start_time

print("Random Forest Prediction Time:", rf_prediction_time)
print("Future Random Forest Prediction:", future_predictions_rf)

Random Forest RMSE: 1298.602442110225
Random Forest MAPE: 0.11149943440215529
Random Forest Training Time: 3.59944748878479
Random Forest Prediction Time: 0.007151365280151367
Future Random Forest Prediction: [1415.92]

```

Figure 8: Code Snippet for Random Forest

4. Result:

Model	RMSE	MAPE	Training Time (s)	Prediction Time (s)
Linear Regression	4210.59	1.51%	0.025	0.00079
Gradient Boosting	2022.52	0.51%	0.691	0.00094
Random Forest	1298.60	0.11%	2.303	0.015

Evaluation of models based on above table:

1. Linear Regression:

- RMSE: Highest among the three models, indicating larger prediction errors on average.
- MAPE: Relatively high, suggesting a higher percentage of error compared to the other models.
- Training Time: Fastest training time.
- Prediction Time: Fastest prediction time.

- Future Prediction: The output indicates that for the year 2025, with a fuel price of 172.6 cents per litre, and focusing on passenger cars in Alberta during the months of January, March, and December, the total number of vehicle sales is 4960.

2. Gradient Boosting:

- RMSE: Lower than Linear Regression but higher than Random Forest.
- MAPE: Moderate, indicating moderate percentage of error.
- Training Time: Moderate training time.
- Prediction Time: Fastest prediction time.
- Future Prediction: The output indicates that for the year 2025, with a fuel price of 172.6 cents per litre, and focusing on passenger cars in Alberta during the months of January, March, and December, the total number of vehicle sales is 2396.

3. Random Forest:

- RMSE: Lowest among the three models, indicating smallest prediction errors on average.
- MAPE: Lowest, suggesting the smallest percentage of error.
- Training Time: Highest training time.
- Prediction Time: Moderate prediction time.
- Future Prediction: The output indicates that for the year 2025, with a fuel price of 172.6 cents per litre, and focusing on passenger cars in Alberta during the months of January, March, and December, the total number of vehicle sales is 1416.

After evaluation I found:

- Random Forest outperforms the other models with the lowest RMSE and MAPE, indicating the highest accuracy.
- However, Gradient Boosting offers a good balance between accuracy and computational efficiency, with relatively lower training time compared to Random Forest.
- Linear Regression, while the simplest and fastest model, exhibits the highest prediction errors and may not capture complex patterns as effectively as ensemble methods.

5. Discussion:

If accuracy is the top priority, especially for future predictions, I'd go with Random Forest. It consistently shows the lowest prediction errors, making it a solid choice for reliable forecasts. However, if time is of the essence and I need quick results, Linear Regression might be more suitable due to its fast training and prediction times, even though it's not as accurate as Random Forest.

The observed regional variations in fuel pricing and vehicle types prompt a deeper exploration into the underlying factors shaping these dynamics. This could involve an investigation of the condition of economy of each province, rules and regulations, consumer purchasing patterns and investment in infrastructure of each province. To formulate actions and policies to meet the specific needs and requirements of province, we need to derive insights from these factors.

In addition, the temporal pattern in fuel prices and motor vehicle buying choices highlight issues about the short term and long-term effects of external factors like global energy market, geopolitical instability, and technical breakthroughs.

To formulate strategic planning and risk management efforts, we need to understand these trends over time that can lead to reduce the potential threats and increase opportunities for stakeholders in energy and vehicle sectors and they may be able to make confident decisions.

6. Conclusion and References:

The relationship between fuel prices and motor vehicle sales gives ideas that are important to achieve the targets of reducing the emissions from transportation sector of Canada. This report highlights how fuel prices influences the consumers behaviour towards buying new motor vehicles. As a result, conclusions can derive about the structure of Canada motor vehicle market.

This report shows substantial differences in the fuel prices and motor vehicle sales regionally with Alberta being the focus point. The variations in the fuel prices and motor vehicle sales over the time reveals the demand for adaptation to economic situation and seasonal trends.

The importance of customizing the rules and regulations to solve the specific needs of various motor vehicles are highlighted by the consideration of vehicle type characteristics.

References:

Baumeister, C., & Kilian, L. (Forthcoming). Understanding the Decline in the Price of Oil since June 2014. *Journal of the Association of Environmental and Resource Economists*.

Busse, M., Knittel, C., & Zettelmeyer, F. (2013). Are Consumers Myopic? Evidence from New and Used Car Purchases. *American Economic Review*, 103(1), 220–256.

Council of Economic Advisers. (2015). *Annual Report of the Council of Economic Advisers*. Washington, DC: Government Printing Office.

Gillingham, K. (2014). Identifying the Elasticity of Driving: Evidence from a Gasoline Price Shock in California. *Regional Science and Urban Economics*, 47, 13–24.

Goldberg, P. K. (1998). The Effects of the Corporate Average Fuel Efficiency Standards in the U.S. *Journal of Industrial Economics*, 46(1), 1–33.

Ito, K., & Sallee, J. (2014). The Economics of Attribute-Based Regulation: Theory and Evidence from Fuel Economy Standards. NBER Working Paper 20500. Cambridge, MA: National Bureau of Economic Research.

Jacobsen, M. (2013a). Evaluating U.S. Fuel Economy Standards in a Model with Producer and Household Heterogeneity. *American Economic Journal: Economic Policy*, 5, 148–187.

Jacobsen, M. (2013b). Fuel Economy and Safety: The Influences of Vehicle Class and Driver Safety. *American Economic Journal: Applied Economics*, 5(1), 1–26.

Jacobsen, M., & van Benthem, A. A. (2015). Vehicle Scrappage and Gasoline Policy. *American Economic Review*, 115, 1312–1338.

Klier, T., & Linn, J. (2010). The Price of Gasoline and New Vehicle Fuel Economy: Evidence from Monthly Sales Data. *American Economic Journal: Economic Policy*, 2, 134–153.

Klier, T., & Linn, J. (2013). Fuel Prices and New Vehicle Fuel Economy—Comparing the United States and Western Europe. *Journal of Environmental Economics and Management*, 66, 280–300.

Klier, T., & Linn, J. (2016). The Effect of Vehicle Fuel Economy Standards on Technology Adoption. *Journal of Public Economics*, 133, 41–63.

Leard, B., & McConnell, V. (2015). New Markets for Pollution and Energy Efficiency: Credit Trading under Automobile Greenhouse Gas and Fuel Economy Standards. Discussion paper 15-16. Washington, DC: Resources for the Future.

Li, S., Timmins, C., & von Haefen, R. H. (2009). How Do Gasoline Prices Affect Fleet Fuel Economy? *American Economic Journal: Economic Policy*, 1, 113–137.

Linn, J. (Forthcoming). The Rebound Effect for Passenger Vehicles. *Energy Journal*.

Sallee, J. (2014). Rational Inattention and Energy Efficiency. *Journal of Law and Economics*, 57(3), 781–820.

Bresnahan, T. F. (1987). Competition and Collusion in the American Automobile Industry: The 1955 Price War. *The Journal of Industrial Economics*, 35(4), 457–482.

Breusch, T. S., & Pagan, A. R. (1980). The Lagrange Multiplier Test and Its Applications to Model Specification in Econometrics. *The Review of Economic Studies*, 47(1), 239–253.

Copeland, A., & Hall, G. (2005). The Response of Prices, Sales, and Output to Temporary Changes in Demand. Cowles Foundation Discussion Paper No. 1543, Yale University, December.

Corrado, C., Dunn, W., & Otoo, M. (2006). Incentives and Prices for Motor Vehicles: What Has Been Happening in Recent Years? 2006-09 Finance and Economics Discussion Series, Divisions of Research & Statistics and Monetary Affairs, Federal Reserve Board.

General Motors Corporation. (2006, September 19). Press release.

General Motors Corporation. (2006, March 27). U.S. Sales and Market Update Conference Call - Media Briefing. Mark LaNeve, GMNA Vice President of Vehicle Sales, Service, and Marketing and Paul Ballew, Executive Director, Global Market and Industry Analysis.

Griliches, Z. (1961). Hedonic Price Indexes for Automobiles: An Econometric Analysis of Quality Change. In *The Price Statistics of the Federal Government*, National Bureau of Economic Research.

Plehn-Dujowich, J. M. (2006). The Role of Entry Deterrence in Explaining Why Prices Fall During Times of High Demand. Working Paper, Department of Economics, SUNY.