# Crude Oil Price Volatility and the Impact on the Construction Industry
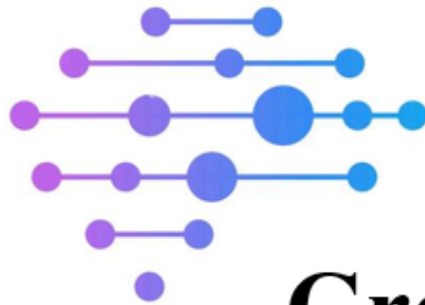
Abdulrehman Khan (60695343)

Justin Lam (77333383)

# 1.0 Introduction

The last few years, the world has experienced a chaotic series of events such as the COVID-19 pandemic, global conflicts, recession, and unprecedented government monetary policy decisions. These events have heavily impacted the economies and industries of many countries. According to a report by Deloitte, with COVID-19 alone, the construction industry suffered from supply chain disruptions, raw material shortages, and increased energy prices leading to loss of profits in the range of $61 billion worldwide. The pandemic exposed an economic flaw of the industry: vulnerability to market volatility of construction resource prices.

The spike in global inflation and the Russia-Ukraine crisis further confirmed these findings. During the global economic slowdown, governments in North America such as the USA, flooded the global market with printed money leading to a CPI jump of 8.5% from 2021 to 2022. In the wake of Russia's invasion of Ukraine, access to Russian gas was shut off that caused a 38% increase in oil prices since January 2022 due to Russia controlling 30% of oil and 35% of natural gas exports. Overall, StatsCanada reported a surge in the building construction price index of 21.7% for residential buildings and 11.2% for non-residential post-pandemic.

Oil is a construction resource that heavily influences the performance of projects and the underlying revenues and expenses provided. This occurs for reason such as:

1. **Supply Chain**: delivery of construction materials produced around the world are heavily dependent on import/export costs of transportation vehicles such as trucks, rail, aviation, and cargo ships. Delivery costs and material availability are heavily influenced by oil prices
2. **Construction Materials and Manufacturing**: manufacturing is often done in regions using fossil-fuel based energy sources. In addition, many construction materials are produced from petroleum, fossil fuel-based polymers: lubricants, insulation, PVC, composite building materials
3. **Construction Equipment**: heavy equipment, machinery, and power source generators is currently powered by gasoline and diesel

The goal of our project is to use predictive analysis to help governments and companies minimize the risk of oil price fluctuations on their construction projects. The project will accomplish this through a two part process:

1. **Oil Price Fluctuation Predictor**: a framework to predict oil price fluctuations will be provided through a list of top indicators to track and an associated model to predict fluctuations. Inputs will be the top indicators: OPEC, OECD, commodities, and financial market instruments. Outputs will be direction and magnitude of price change month-on-month.
2. **Oil Price - Construction Industry correlation**: correlate fluctuations of oil prices with aspects of the construction industry. 5 Features considered in this report are mentioned in the table below. Outputs used are PPI (Producer Price Index) and CPI(Consumer price index) change from month to month.

By providing a framework to predict oil price fluctuations and correlating its fluctuations to the performance of the construction industry, the report will provide clients with the foresight and actionable information to minimize oil price risk on their underlying project costs. The data sources for these models are outlined in the table below:

| Part 1: Oil Price Fluctuation | Source |
|---|---|
| Oil and Commodity Prices | St Louis Fed: https://fred.stlouisfed.org/ |
| Financial Asset Markets | St Louis Fed: https://fred.stlouisfed.org/ |
| Crude Production, Inventory | US EIA: https://www.eia.gov/ |
| Part 2: Oil-Construction Correlation | Source |
| CPI Engineering Services | https://fred.stlouisfed.org/series/GOV54133TAXABL157QNSA#0 |
| PPI Hourly Earnings | https://fred.stlouisfed.org/series/CWSR0000SEHE |
| PPI Power Related Products (construction machinery) | https://fred.stlouisfed.org/series/PPIENG |
| PPI Ready Mix Concrete | https://fred.stlouisfed.org/series/PCU32733273#0 |
| PPI Oil Transportation | https://fred.stlouisfed.org/series/PCU486110486110 |

Table 1.1: Data Sources

# 2.0 Background

The oil markets are a globally traded commodity that is sensitive to changes in economic conditions and political events. Oil is a unique commodity in that it heavily affects the economic prosperity of exporting and important nations. With the diversity of factors influencing prices, predicting price shocks is a major challenge [1]. A paper named *Predicting Oil Price Movements: A Dynamic Artificial Neural Network Approach* by Ali Abbasi Godarzi [2] proposes that the state of the art model for the prediction of oil price movements is a dynamic artificial neural network approach. The study develops an ANN approach known as Nonlinear Auto Regressive model with eXogenous input (NARX) and uses the Mackinnon-White-Davidson (MWD) test to analyze and compare different models of oil price prediction. The study hypothesizes that the application of the NARX model enhances dynamic performance of the model and improves the ability of ANN methodology to predict oil price, particularly the occurrence of price shocks. ANN is a good technique for dynamic models for modeling and forecasting behavior of nonlinear economic variables.

The methodology consists of first creating a time series model to identify factors affecting oil price. The ANN model then verifies the data and the NARX model is used to factor time delays in the analysis. This iterative approach is used to improve the R-squared score. The time series model analyzes previously observed values in indicators, creates a trend, and predicts future values. The modeling of the relationships between variables and verifying the credibility and relationship between them. This will be done by identifying most relevant variables, selecting optimal models, and improving model performance [3]. Using concepts learned in CIVIL 498A, our model will try to optimize this first Time Series step of this state of the art model. Concepts used include data processing and cleaning, explorative analysis, feature selection, correlation analysis, regression analysis, classification algorithms, and cross validation.

The application of ANN and NARX in the models, the ANN is a multi-layer perceptron neural network that takes inputs from external sources, combines them with a nonlinear operation, and produces final results. Weights for neurons are determined using error back-propagation until the desired degree of

accuracy is achieved. The NARX model is a dynamic version of the ANN model that considers the factors of time. Due to limitations of knowledge scope in CIVIL 498A, our team was unable to execute these applications.

The variables used in Godarzi's model are supply and demand variables from categories such as economic growth, energy demand, oil prices, and OECD countries. The time series model used a combination of linear and log models with gold prices, energy production, and energy use as the top 3 indicators. The best model was the linear linear model. Applying the NARX and ANN applications, an R-squared score of 97% was achieved with an MAE of 4.96%, this yields a stellar prediction accuracy.

# 3.0 Methodology

**Part 1: Oil Price Fluctuation Predictor**

The first part of the methodology is the prediction of oil prices changes for both direction of movement and magnitude of change. In order to accomplish this, the model must be further split into 2 separate parts: prediction of direction and prediction of magnitude change. The former requires a classification based model whereas the latter uses a regression model. The overall workflow of the models highlighting the supervised learning method is shown in the chart below.
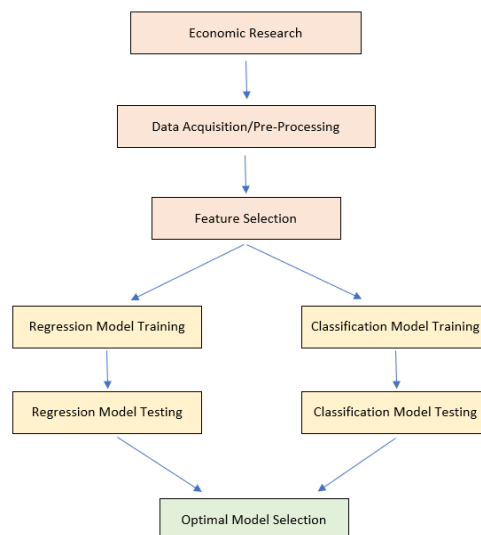


Figure 3.1: General Workflow

The process began with rigorous economic research to understand the underlying economic principles that affect oil prices such as supply and demand. Further economic research was conducted to narrow down the list of applicable data sets: geo-political stability, non-OPEC fuel production, OPEC fuel production, natural gas production, crude inventories, crude production capacity, interconnected commodities, and financial assets.

Each data set was modified and cleaned for data processing purposes. The data sets were then compared individually with oil prices to narrow down specific features of the data set to use. The methodology used was a correlation matrix and regression analysis using linear, ridge, and lasso regression. Correlations were computed using the Pearson Correlation . Regression analysis used an 80/20 train-test split with a random state of 42. The filter method was then used for feature selection based on ranking of correlation, RMSE, and R-squared scores. The top features are outlined in the table below. Specific features used for

the model were: OPEC/Non-OPEC crude production, Saudi crude production, OECD crude inventories, oil market futures spreads, wheat/copper/corn commodities, US Dollar strength basket, 1mo/1yr/5yr/10yr US government bonds.

| | Correlation | R2 | | | RMSE | | |
|---|---|---|---|---|---|---|---|
| | | Linear | Ridge | Lasso | Linear | Ridge | Lasso |
| Non-OPEC Production | 0.2 | 0.0351 | 0.0349 | 0.0319 | 23.044 | 23.0421 | 23.009 |
| Natural Gas Production | 0.67 | 0.25743 | 0.25743 | 0.2574 | 25.482 | 25.482 | 25.4825 |
| OPEC Production | 0.4-0.64 | 0.57252 | 0.4378 | 1.0688 | 5.697 | 21.4297 | 6.537 |
| Inventory vs Futures Spread | 0.64 | 0.4135 | 0.4135 | 0.4132 | 6.445 | 6.445 | 6.447 |
| Commodities | 0.16-0.3 | 0.01453 | 0.01458 | 0.02812 | 23.1814 | 23.18079 | 23.0225 |
| Financial Instruments | 0.18-0.34 | 0.1489 | 0.1231 | 0.2243 | 25.7985 | 25.8615 | 30.567 |

Table 3.1: Features vs Correlation/R2/RMSE Score

The top features were compiled along with the price of WTI crude oil into a single data set for classification and regression model analysis. WTI was chosen as it is the oil produced by the USA and used most in North America. All prices and features are calculated as month-on-month changes using data ranging from 1990 - 2021. Further details can be found in section 4.0 Metrics.

The workflow of the classification and regression models are highlighted in the diagram below.
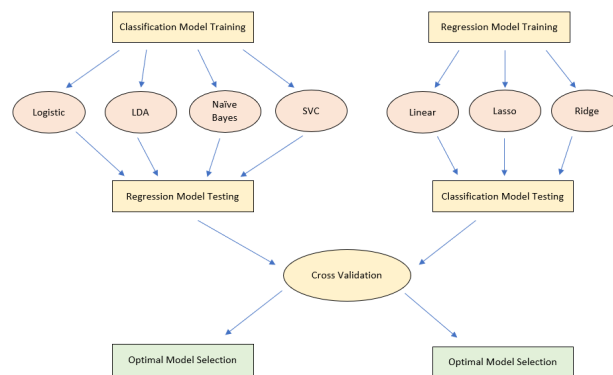


Figure 3.2: Modeling Workflow

The classification models predict price change direction and use an 80/20 train-test split with a random state of 42. Logistic, linear discriminant analysis, naive bayes, and support vector classification models were tested using ten fold cross validation. The SVC model was modeled under 15 combinations of: regularization parameter C [0.5, 1, 2], kernels [linear, poly, rbf], and poly kernel degrees [2,3,5]. Precision, accuracy, recall, and log loss scores were computed and the models were ranked based on F1 test score. The model with the maximum F1 score was chosen and a confusion matrix was plotted to test the models ability to predict a price increase or decrease based on feature information.

The regression model predicts the $ of price change and it follows the same train-test split and cross validation parameters as the classification model. Linear, ridge, and lasso models were used with a max iteration of 1000. The model with the largest R squared score was chosen to test the model's ability to predict the $ value of price increase/decrease.

**Part 2: Construction Features Predictive Analysis**

This part of the model involved performing predictive analysis on input features to determine their correlation with fluctuating oil prices. Although there is an exhaustive list of different factors that could affect the construction industry, going too broad in this predictive analysis could've skewed our data results. The best way to start implementing predictive analytics solutions for your next construction project is by first honing in on your area of focus. Moreover, getting the most out of predictive analytics requires you to centralize and standardize your data. The higher quality your data input is, the higher quality, and thus better able to predict, your data output is. Standardizing the data provided an accurate summary of the amount of data where the goal is to see patterns. Moreover, The data that we use to compute correlations often contain missing values. This can either be because we did not collect this data or don't know the responses. Various strategies exist for dealing with missing values when computing correlation matrixes. The practice adopted was to fill in the num values with mean values for that column. The null values for some columns were very few and wouldn't significantly affect future and therefore was implemented

Following this, our team intended to build a classification model which could provide the strongest relation with mentioned features. For this a correlation matrix, using pandas class, was used to determine this. This model could be used by construction industries to:
- Analyze oil price fluctuations to forecast future trends.

- Include other relevant features that the industry might relevant

Generally the correlation matrix is used to determine features with the higher correlation and remove features which could be discounted and used for future predictive analysis. However, based on the limited features included, all of these were used. To select the best model for prediction, the models were used with the intent to select one with the highest regression ($R^2$) and R.M.S.E score to provide a greater accuracy for prediction. The results from each model are outlined below:

| Regression Model | R2 | R.M.S.E |
|---|---|---|
| Linear Regression | 0.980 | 5.210 |
| Ridge Regression | 0.981 | 5.212 |
| Support Vector Regression | 0.890 | 5.201 |

Table 3.2: Model vs. R2/RMSE Score

Based on these results, **Ridge Regression** was chosen for the predictive analysis. All features were used in this analysis.

# 4.0 Metrics

**Part 1: Oil Price Fluctuation Predictor**

The features mentioned in section 3.0 were used to train the model. The definitions for these metrics are provided in the table below.

| Feature | Units |
|---|---|

| Crude Production (OPEC, non-OPEC) | Million barrels/day |
|---|---|
| Crude Inventory (OECD) | Million barrels |
| Futures Spread | 12-month future - current month future |
| Commodities | $/unit |
| Bond Yields | % yield |
| US Dollar Strength (DXY) | Base = 100 (value relative to foreign currency bask |

Table 4.1: Feature Descriptions

As mentioned in section 3.0, the performance metric of our models were the F1 and R-squared scores for the classification and regression models respectively. F1 score was chosen as an all-encompassing metric as it observes both the precision and recall of a classifier and combines it into a single metric to measure accuracy. The prediction of an increase or decrease in oil production is a purely binary measure and therefore F1 score is a suitable metric. R-squared is an appropriate metric because it measures the proportion of variance between true and predicted oil price changes. The best classification model produced an F1 score of 0.694 using Logistic Regression and the best regression model produced an R2 score of 0.124 using ridge regression.

**Part 2: Construction Features Predictive Analysis**

The features mentioned in Table 1.1 were used to train the model. The features used for this included trends in Price Indexes. Engineering Services was the Consumer Price Index whilst Hourly Earnings, Power Related Products (construction machinery), Ready Mix Concrete, Oil Transportation cost were PPIs respectively.

The performance metrics used were the R-Squared scores and R.M.S.E scores for the regression model. R-square is a comparison of residual sum of squares (SSres) with total sum of squares(SStot). This was used to assess the goodness of fit for the model when measuring variance between true and predicted values.

# 5.0 Results

**Part 1: Oil Price Fluctuation Predictor**

The 3 main types of crude oils traded globally are West Texas Intermediate, Brent, and Dubai crude. The correlations were plotted between the 3 oil price movements, yielding a correlation of 99%. Given the strong correlation, West Texas Intermediate was chosen as the benchmark oil price for analysis.

Non-OPEC oil production accounts for 40% of the world's oil exports. Non-OPEC producers are often public companies and price takers while OPEC suppliers are nationalized oil organizations with a strong interest in managing the supply of oil. When non-OPEC production is reduced, a greater reliance on OPEC is needed causing OPEC to have greater influence on oil prices which generally push up prices. Overall the 20% correlation and 9% R2 score is low but a considerable indicator that affirms the general trend and can be factored into the model.

Natural gas is an up and coming alternative fuel that is cleaner, safer, denser. New discoveries of gas deposits may put a damper on crude oil price increases. The 68% correlation between natural gas production increases and price decreases and the 26% R2 score affirms the general trend and can be factored into the model. Companies should monitor this correlation; as natural gas technology continues to develop and improve its production/use capabilities in the future, this correlation may increase.

OPEC nations provide 60% of the world's oil exports. Through central coordination, the big 5 OPEC countries (Saudi Arabia, Iraq, Iran, UAE, Kuwait) have a significant effect on global crude prices. The correlation matrix between oil prices and OPEC production was created.



Figure 5.1: OPEC Correlation Matrix

Many interesting relationships are shown in the graph. Most notably, Kuwait, Saudi Arabia, UAE, and Iraq have the highest impact on WTI prices in descending order. The political stability and oil policy decisions of these nations should be heavily monitored to forecast prices. The matrix also shows that coordination and compliance between OPEC countries and the big 5 policy makers is around 45%. This indicates that OPEC announcements of supply changes does not necessarily translate to equal movements by OPEC countries. Using ridge regression, the R2 score for price predictions was 0.458 which shows a general linear trend as shown in the figure below.
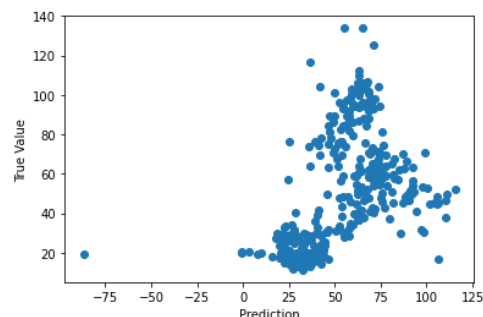


Figure 5.2: Price Prediction, OPEC Data

Global inventories have a strong effect on oil prices as it directly impacts the supply and demand curves. If inventories rise, the market expects a strong future demand and the price of futures contracts for oil will

increase leading to increased spreads between futures oil contracts and existing oil prices. The change in spread for a 12 month - 1 month spread had a 64% correlation with inventory changes. Monitoring inventory levels and inventory projections for OECD and OPEC countries affect the future/current spreads the most with a linear regression R2 score of 41.4%.

Due to many industries being independent on oil to function, there are correlations between commodity prices and prices of crude oil. Analyzing the data, wheat/corn/copper have the highest correlations of 0.16/0.17/0.3. A study done by the US Energy Information Administration [4] showed that commodity correlations are relatively inconsistent as shown in the graphic below. The best model output was a Lasso regression with an low R2 score of 2.8%

As a result of weak commodity correlation, general market sentiment was a better alternative to monitor. Crude oil is a stable investment asset that is used to diversify risk in case of inflation or global crisis. As government bond rates rise, this indicates an economic slowdown and therefore prices of oil increase. US government bonds yield changes were plotted against changes in crude prices. Longer-term bonds such as the 5 and 10 year bonds had higher correlations with oil changes than the short term 0.34 vs 0.18. This is most likely attributed to longer-term bonds being an indicator of recession outlook. The best Lasso regression model yielded a 22.4% R2 score.

Another similar indicator is the strength of the US dollar relative to other countries. As oil is traded mainly in US dollars, dollar depreciation puts upward pressure on prices for several reasons. First, the cost of oil outside the US decreases leading to increased demand. Second, effective oil revenue decreases leading to price hikes to maintain real revenue. Finally, the commodity of oil is an asset to hedge against the risk of dollar depreciation. The US dollar basket and price of oil had a 23% correlation and an R2 value of 0.224 for Lasso regression which confirms the trend.

After evaluating the usefulness of these features and the methods with which to interpret them, the features were loaded into our regression and classification models. The Logistic Regression model had the best performance for the classification model after ten-fold cross-validation. The results of other models in comparison are shown in the table and confusion matrix for logistic regression below.

| Model | Precision | Accuracy | Recall | F1 | Log Loss |
|---|---|---|---|---|---|
| Logistic | 0.6356 | 0.6057 | 0.7722 | 0.6943 | -0.653 |
| LDA | 0.6325 | 0.5986 | 0.7833 | 0.6936 | -0.6997 |
| Naïve Bayes | 0.6429 | 0.5429 | 0.4708 | 0.5362 | -1.4429 |
| SVC | 0.6203 | 0.5845 | 0.7824 | 0.6839 | - |

Table 5.1: Classification Model Metrics

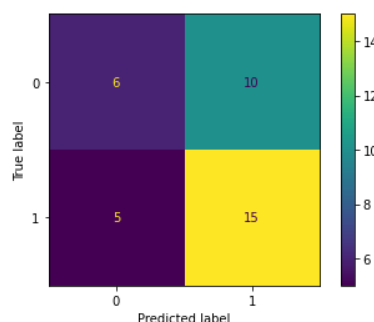The confusion matrix for our Logistic regression model is shown in the figure below (1 = increase)

Figure 5.3: Confusion Matrix

The Ridge regression model had the best performance for the regression model after ten-fold cross validation. The results of other models in comparison are shown in the table below.

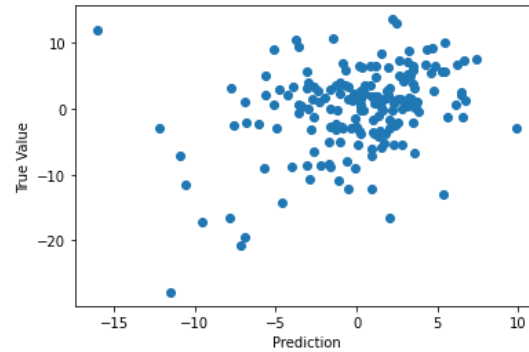| Model | R2 | RMSE |
|-------|--------|----------|
| Linear | 0.0513 | -32.288 |
| Ridge | 0.1241 | -29.662 |
| Lasso | 0.1222 | -30.5817 |

Table 5.2: Regression Model Metrics



Figure 5.4: Oil Price Change, Predicted vs True

Overall, the price change classification model had a 61% accuracy and 69% F1 score while the price change prediction model had a low R2 score of 12.4%. From the scatter plot, it appears that the model more accurately predicts price increases as opposed to price decreases. The correlation matrix for feature prediction can be seen below. PPI for Power related products (construction machinery) obtained the highest score.

**Part 2: Construction Features Predictive Analysis**

Results of the correlation matrix discussed in section 3.0 are shown below with PPI for Power related products (construction machinery) scoring the highest results.
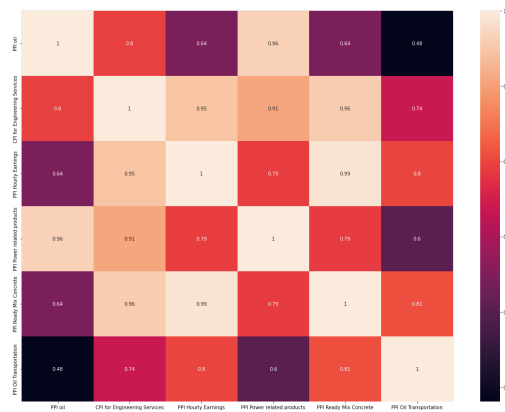


Fig 5.5: Correlation matrix: Oil prices vs. feature prediction

An important note is the high correlation of all features included because of which all of them are used for future forecasting. 2 different PPI values for oil were assumed 18 and 27 respectively (index 100=1980). The results from the predictions are shown in the figure below:

| | |
|---|---|
| Engineering Services (CPI) | [132, 141] |
| Hourly Earnings (PPI) | [150, 159] |
| Power related products (PPI) | [87, 102] |
| Ready Mix Concrete (PPI) | [153, 164] |
| Oil Transportation (PPI) | [140, 150] |

# 6.0 Limitations

**Part 1: Oil Price Fluctuation Predictor**

Overall, our model exceeded the limitations of the classification model while there is a must to be desired for the price change regression model. However, it is clear that the dynamic nature of oil price fluctuations is too complex for the simple concepts learned in an introduction to machine learning course. One limitation of the study was in economic research and feature selection. The dynamic nature of the features raises two main questions: does correlation amount to causation? If yes, is the change of price or the change in a feature the independent variable? Or does it depend on other situational factors? As we observed through cross-validation of models, high correlation scores do not amount to high R-squared scores in modeling. Another challenge with dynamic economic modeling is the need for data from a wide variety of sources that may not always provide the same accuracy or consistency due to barriers such as: technological inaccessibility, censorship, poor reporting practices, etc. The limited data points in our models severely impacted the training ability of our models.

Another major limitation of the model is the difficulty in fitting data for time delays between the output of information and given market reaction. For example, the feature of oil future spreads is information that is constantly updated by the minute as it is traded on international stock exchanges. In contrast, non-OPEC production of oil is difficult to track consistently due to the lack of communication and central coordination between countries. Market prices may adjust for oil future spreads immediately whereas the market reaction to non-OPEC production may be triggered several months later; yet due to the limitations of our modeling, both these feature changes will be monitored on the same timeline. As Ali Godarzi proposes in his state-of-the-art model proposal, a neural network structure may need to be implemented to adjust for time variations which he proposes can be done through a NARX ANN model.

**Part 2: Construction Features Predictive Analysis**

During the data selection process, it was preemptive that the data would be incomplete for data analysis. Essentially, based on the project scope, a much larger variety of data sets needs to be included to replicate an industry-standard model. The effect of this could be a potential limitation on data fitting. Data analysts can add and remove construction features based on construction requirements and include larger data sets for a better prediction.

# 7.0 References

[1] Tang, L., & Hammoudeh, S. (2002). An empirical exploration of the world oil price under the target zone model. Energy Economics, 24(6), 577-596

[2] Godarzi, A. A., Amiri, R. M., Talaei, A., & Jamasb, T. (2014). Predicting oil price movements: A dynamic Artificial Neural Network approach. Energy Policy, 68, 371-382.

[3] Gujarati, D. N., Porter, D. C., & Gunasekar, S. (2012). Basic econometrics. Tata mcgraw-hill education

[4] U.S. Energy Information Administration - EIA - independent statistics and analysis. EIA. (n.d.). Retrieved April 25, 2022

[5] Machine Learning vs Predictive Modelling: Top 8 Vital Differences." *EDUCBA*, 30 Apr. 2021, https://www.educba.com/machine-learning-vs-predictive-modelling/.

[6] Kearns, Shannon. "The Limitations of Predictive Analytics Tools and Why Execs Should Care." River Logic,https://www.riverlogic.com/blog/the-limitations-of-predictive-analytics-tools-and-why-execs-should-care#:~:text=While%20this%20is%20a%20useful,certain%20threshold%20that%20requires%20action.