

# Analysis on the Basis Spread Between Crude Oil Spot and Futures Price

Yichao, Dai

College of Business and Public Management  
Wenzhou-Kean University  
Wenzhou, China  
yichaod@kean.edu

**Abstract**— West Texas intermediate Crushing, Oklahoma Crude Oil as one of the main benchmarks in oil pricing, becomes a popular higher commodity used by producers and refiners. Traders in the futures market might suffer some losses or gains caused by basis risks. This project aims to use the multi-linear regression analysis to predict the basis spread between the West Texas intermediate Crushing, Oklahoma Crude Oil spot price, and its 4-month futures contract price. Hedgers can partly eliminate the basis risks or even gain some arbitrage profits by knowing the basis in advance. In this project, researchers try to pick the most relative explanatory variables to build a multi-linear regression model to predict the basis of WTI OK crude oil. Oil production level, oil consumption level, speculative index, economic conditions, and several financial market factors had been taken into consideration. The research also excludes the impact of the Covid-19 outbreak on a WTI basis spread to reduce the impact on the accuracy of the prediction model. The predictive model built in this project is powerful to predict the movement direction or even the basis spread when there is a significant event in the world market such as a financial crisis, however, when there is no significant event over the world market, the predicted basis cannot be explained by the linear method. The final model has a relatively good adjusting R square equal to 25.6%, which can explain 25.6% of the WTI crude oil basis spread. **Keywords**- basis spread; crude oil; model selection; OLS regression; LASSO regression

## I. INTRODUCTION

West Taxes Intermediate crude oil has a relatively low density since it only has 0.24% sulfur content, therefore, WTI crude oil is also referred to as light crude oil. It's a specific grade of crude oil from Texas that spot and futures prices serve as one of the main benchmarks in oil pricing. The West Taxed Intermediate is the commodity of New York Mercantile Exchange's oil futures contract and is the main oil benchmark in North American. And the main delivery and price settlement point for West Taxed Intermediate crude oil in Cushing, Oklahoma. Also, WTI crude oil futures prices are also included in the Bloomberg Commodity Index and S&O GSCI commodity index. Therefore, WTI crude oil futures contracts gradually develop as hedging tools used by producers and refiners. Basis usually occurs when the hedgers are uncertain of the exact date when the asset will be bought or sold, and the hedgers may be required to close their current position before maturity days. And basis risk refers to financial risk

that will happen while using hedging strategies which brings the potential for gains or losses. When there exist large amounts of shares or contracts in transactions, basis risk may bring significant losses or gains. In this project, the asset to be hedged and the asset under the futures contract is the same. Consequently, the basis should be zero at the maturity date. However, the basis may be positive or negative prior to maturity. Hence, investors invested in WTI crude oil could be assisted by exploring the basis spread of WTI crude oil spot and futures to eliminate some basis risks and minimize losses or even gain some arbitrage profit. However, previous researches have shown that basis prediction is a complex process, involving different factors in different markets [5]. Some of the literature shows that the non-linear method is an optimal choice to predict the movement of basis. The multi-regression analysis helps determine correlations between one dependent variable and more independent variables and to have some predictions among these variables. (Unver & Gamgam, 1999, as cited in Uyanik & Güler, 2013). The purpose of this research is to predict the basis spread between WTI Crushing, Oklahoma Crude spot price, and its 4-months futures contract price by using different multi-linear regression methods and chose the most fitted one.

## II. LITERATURE REVIEW

Previous researches showed many great ideas on how to predict the futures price of crude oil and the factors influencing the spot market. However, only a few researches are exploring the basis spread between the crude oil spot price and the futures price. As a result, in this literature review, the researchers aim to explore the factor that can both affect the spot market and futures market of the WTI crude oil.

### A. Oil Production Capacity & Oil Consumption Level

Crude oil is one kind of exhaustible resource in the world [16]. The Middle East and North Africa, as the main supplier, play an important role in the spot market. In 2009, Kaufmann and Ullman [14] try to explore in their article "Interpreting causal relations among spot and futures prices" that how the innovations in oil prices first enter the market, after comparing the spot price and futures price in a different region, they find that the innovation first appears in the Middle East region and then spread to other spot and futures price, which suggest the situation in Middle East region plays an important role in the fluctuation of the WTI spot and futures market.

### *B. Situation in Middle East and North Africa (MENA)*

Miao and his team [16] discussed the influential factors of the crude oil price forecasting and indicated the oil production capacity in the Middle East area had a strong relationship with the WTI crude oil price. Miao [16] also found that the WTI spot price and futures price fluctuate because of some political factors, such as the total amount of terrorist attacks in the Middle East and North Africa. Considering that oil is actually an important energy resource, it would be easily targeted by terrorism, and it seems that the spot price and the futures price of crude oil would be vulnerable to be influenced by terrorist attacks. However, Holwerda and Scholtens [11] found that the spot price and the futures price of crude oil did not have a significant reflection to the terrorist attacks, since the market had already included the terrorist attacks in the risk premium previously. It means that the market had already made the reaction to the terrorist attacks, and it makes the spot price and futures price of crude oil did not have too many changes. Of course, the spot price and the futures price of crude oil in the different markets would also have different reactions to the terrorist attacks. Kollias et al. [15] tested the co-movement between oil returns and four market indexes under the influence of war and terrorism. It shows that the spot price and futures price of crude oil in S&P500, FTSE100 would have more capacity in absorbing the risk of the terrorist attacks, and the spot price and the futures price of crude oil CAC40, DAX did not.

### *C. Speculative Index: NTMI-CNCN Index*

Kaufmann and Ullman [14] also stated that the relationship between the spot market and the futures market was relatively weak, however, when the market initiated a long-term increase in oil price, this trend would be exacerbated heavily by the speculators. As a result, the number of speculators in the WTI oil futures market may have some contribution to the variance of the basis of the WTI crude oil. Other researchers [5] also argue that the relationship between co-movement of the spot & futures prices and the net non-commercial futures position, which could be measured as a speculative index, they found that the speculative index played an uncertain role in the spot and futures market. Although the previous research before 2008 has shown the price collapse would be accompanied by a significant drop in the speculative index, the oil price collapse in 2008 did not indicate a large drop in the speculative index [5]. Moreover, Bu [7] stated in his research that “reveal that the position changes held by speculative traders will cause crude oil price movement”, especially when the financial crisis occurred, speculative traders would inject the futures price large vitality. Moreover, Stoll and Whaley [20] found in their research that excessive speculation will cause the same direction movement in the spot and futures price of crude oil.

### *D. US Economics Condition: US GDP Growth Rate, US CPI, Dollar Index*

US GDP Growth Rate & Dollar Index Economics condition plays a significant role to determine the spot price of crude oil. Miao [16] stated in his research that the world economic growth is closely related to crude oil demand. Higher global economics would raise the spot price of the oil, however, a lower global economic growth rate would lead to the fall of the spot price. He also stated several factors such as

the steel production and ISM manufacturing index which might have some effect on the spot price. Study [22] also show that a weaker dollar exchange rate may result in a higher spot price. The study also showed a similar relationship in the oil futures price, Algieri [2] found in his research that a weaker US/Euro exchange rate could bring a decline in the futures commodity return. However, the relationship is not as strong as the price in the spot market. Furthermore, Amendola et al. [4] explored that the expansionary monetary policy would have a negative impact on the crude oil future price and the fluctuation of the industrial production would make the crude oil future price have the opposite correlation. Then, Frondel et al. [8] explained in their research that the production of crude oil in the US would also create an opposite influence on the spot price of crude oil. Furthermore, Basistha and Kurov [6] also found that the changes in the crude oil price would create a significant influence on the federal funds target rate, but, with the analysis of the VAR model, it seems that the crude oil price might not have contemporaneous feedback with the federal funds rate shocks. Studies have shown that the financial factors could produce effects both on the spot price and futures price of the crude oil. Algieri [2] stated in his research that Standard and Poor Index 500 positively affected the commodity futures market. Miao [16] also found a significantly positive relationship between the oil stock price and the oil spot price. Moreover, according to the article by Jones and Kaul [13] and Sadorsky [19], they found that the stock market and oil prices tended to move in the same direction.

## III. METHODOLOGY

The project would conduct time-series quantitative research, which aimed to explore the influencing factors on the basis risk of WTI Crude Oil. The overall approach was to use the unit root test and cointegration test to build a rational and valid ordinary least square (OLS) multi-regression model, and then the project will use the features-deduction method to avoid model overfitting. The project would also deal with the multicollinearity problem by checking the variance inflation factor (VIF). Moreover, The project would make some assumptions about the initial model, such as normality of the error, to fulfill the condition of use of the model and method. Last but not least, the project would use the adjusted R square and root mean square error (RMSE) to evaluate the overall model. All the data are quantitative data, which is collected through the existing data sources, Bloomberg and the U.S. Energy Information Administration [12], which is the official energy statistics from the U.S. Government. These two data sources are reliable data sources for many researchers working on econometrics and finance. Because of Coronavirus-19, the crude oil spot market and the futures market are affected heavily in some time period [1]. To eliminate the effect and corresponding effect of Coronavirus-19, the project will only select the data between January 2000 and December 2018, a 19-year range data set. All the explanatory variables would also be selected in this specific time period. The project will typically use the STATA and R to do the data cleaning, data analysis, data visualization, features selection, and model building

### A. Multi-regression Linear Model Selection

In this project, the multi-regression model is built based on the ordinary least square method (OLS), which is a typical method to use several explanatory variables to predict a response variable. By using the ordinary least square method (OLS), we want to minimize the sum of square error (SSE) between the observed response variable and the value of the response variable predicted by the multi-regression model [21].

In this project, we will use features deduction method to update the OLS model. We would use the backward Akaike information criterion (AIC) method and least absolute shrinkage and selection operator (LASSO) method to do the dimension reduction to avoid the overfitting or multicollinearity problems. AIC method start with a full OLS model and remove one variable at a time to get a new AIC, which is calculated as follow:

$$AIC = 2k - \ln(L)$$

Where:

k: number of estimated parameters in the model.

L: maximum value of the likelihood function for the model.

The method aims to explore whether the model can achieve a lower AIC, model with a lower AIC would be better. The final model will be selected if deleting any one of the remaining variables cannot achieve a lower AIC.

Another dimension reduction method is LASSO. LASSO method will do the adjustment to the corresponding coefficients of each variables by using following formula:

$$Lasso = \min \left[ \sum_{i=1}^n \left( y_i - \sum_j x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right]$$

When lambda is equal to 0, no parameters are eliminated. The estimation of the coefficient is exactly equal to the previous one with OLS linear regression. However, As lambda increases, more and more coefficients are influenced and trend to be zero, when lambda is tend to infinite, all coefficients will be tend to 0 and be eliminated. LASSO method selection process will choose the best lambda with its corresponding model as the predictive model.

### B. Unit Root Test

In a time-series project, it is important to avoid spurious regression, which is caused by a similar local trend. It is possible to build a significant model even if the relationship none exist. In time-series data, it is a common occurrence when the data is not stationary [10]. As a result, the project should do the unit root test to make sure all the variable should be stationary, including the response variable. The common method is (Augmented) Dickey-Fuller tests [9]. The null hypothesis of the ADF test is the presence of the unit root (non-stationary), and the alternative hypothesis state that the variable is stationary. According to the empirical significant level in past researches, the significance level for this project was set as 5%, as a result, when the p-value is higher than the significant level, we cannot reject the null hypothesis and state the variable is non-stationary. If the variable is non-stationary, we may need to transform the variable to a stationary variable to adjust the model.

### C. Assumption of the Model

Although ordinary least squares (OLS) are used in a variety of economic and financial analysis, to be able to reasonably explain the parameters and outputs of the model, the least-squares method needs to satisfy three main assumptions: normal error, homoscedastic, and without multicollinearity. Assumptions are sometimes very difficult to satisfy, which needed certain tests that would allow the project to fulfill those assumptions. Firstly, the observed data should have a normal distribution error. Non-normality residuals can be misunderstood by the confidence interval model for its predictive ability and create skewed distributions by increasing the appearance of outliers [18]. In technical terms, “the Assumption of Normality claims that the sampling distribution of the mean is normal” [17]. Nearly all the inferential statistics such as t-test and ANOVA reply on the assumption of normality Secondly, the linear relationship should exist between the response variable and explanatory variables. With other explanatory variables holding constant, there should be a linear function to decide the relationship between the response variable and each explanatory variable. Moreover, The effects of each variable on the predictive value of the dependent variable should be additive [18], otherwise, the estimation will be misleading. Thirdly, multicollinearity refers to the highly correlated relationship among the explanatory variables. Multicollinearity destroys the statistical significance of an independent variable [3], which means we may construct a variable as a linear function of other variables. The common evaluation of Multicollinearity is to use the variance inflation factor (VIF) [18], if VIF is equal to 1, we can state that no explanatory variable are correlated. However, when VIF is more than 10, It indicates a high correlation and is cause for concern. The ideal condition should be VIF less than 3, however, it does not cause concern when VIF is less than 10.

### D. Evaluation of The Models

Based on the output of the model, we would like to evaluate the overall model performance by checking the Adjusted R square. The basic R square formula is shown as following formula:

$$R^2 = \frac{\text{Expected Variance}}{\text{Total Variance}}$$

The formula suggest that  $R^2$  can indicate to what extent can the model explain the variance of the response variable. For example, a R-square value 0.8 indicates that 80% of the variance of dependent variables can be explained by the selected explanatory variables. However, the project would use the adjusting R square rather than simple R square in that the project use multi-regression model with multi variables; consequently, adjusting R square is the modified simple R square. The adjusting R square formula is showing as:

$$R_{adj}^2 = 1 - \left[ \frac{1 - R^2}{n - k - 1} (n - 1) \right]$$

In Finance field, the model with an adjusted R square above 0.7 would generally be seen as a high predictive model. However, for times series data, because all the variables should be converted to the stationary

variables, as a result, an adjusting R square higher than 0.25 can be considered as a quite good model.  
Another way to know the model performance is calculating the root mean square error, which is a general way in machine learning to evaluate the model accuracy.

#### IV. RESULTS

In this section, the research will show the result of the unit root test, model assumption checking, model selection processs, model performance and the comparison among different models.

##### A. Unit root test and cointegration test

Times-series data need the unit root test to avoid spurious regression. These variables need to check whether there exists a unit root. The following table (table 1), shows the ADF test statistics and the significance level. The significant level is overall 5% in this project. The null hypothesis on the stationary test is that there exists a unit root,

and the alternative hypothesis is the variable is a stationary variable. The null hypothesis on the drift test is that there exists a unit root with drift, the alternative hypothesis is that one of these two conditions in the null hypothesis does not meet. If the test statistics is larger than the 5% critical value (absolute value), then we can reject the null hypothesis and accept the alternative. As result, we have found that the response variables and four other explanatory variables including NTM1-CNCN Index, OCED Oil Consumption, US GD, and CPI US are stationary variables, which is also named as I(0) variable. Other explanatory variables are considered as the non-stationary variables. We would like to convert the non-stationary variables to their first difference and check whether the first differences are stationary. The table 2 shows the test statistics for their first difference. We found all the first different of the stationary variables are stationary. As a result, we have converted the whole model dataset to a stationary dataset. Because the response variable is an I(0) variables, as a result, there is no need to do the cointegration test.

**Table 1: Unit root test before the first-difference transformation**

Variable	ADF Test on Stationary	ADF Test On Drift	Constant 5% Critical Value	Constant and Trend 5% Critical Value
Basis	-4.5437***	6.9471***	-3.43	4.75
Dow Jones Oil Index	-1.8973	1.6993	-3.43	4.75
NTM1-CNCN Index	-3.4675***	4.0867	-3.43	4.75
Open Interest-Crude Oil	-3.0976	4.0676	-3.43	4.75
CPI US	-4.4786***	6.7045***	-3.43	4.75
US Federal Rate	-1.263	1.8419	-3.43	4.75
OCED Oil Consumption	-4.7136***	7.4353***	-3.43	4.75
S&P 500 Index	-1.8281	2.5567	-3.43	4.75
US Dollar Index	-1.3244	1.0478	-3.43	4.75
US GDP Growth Rate	-3.8136***	4.8894***	-3.43	4.75
OPEC Total Surplus	-2.5598	2.2842	-3.43	4.75
Terrorism Attack	-1.7006	1.1564	-3.43	4.75

**Table 2: Unit root test after the first-difference transformation**

Variable	ADF Test on Stationary	ADF Test On Drift	Constant 5% Critical Value	Constant and Trend 5% Critical Value
Dow Jones Oil Index	-8.8393***	26.0969***	-3.43	4.75
Open Interest-Crude Oil	-10.0222***	33.4821***	-3.43	4.75
US Federal Rate	-13.1882***	57.9829***	-3.43	4.75
S&P 500 Index	-10.1739***	34.5773***	-3.43	4.75
US Dollar Index	-9.9518***	33.0143***	-3.43	4.75
OPEC Total Surplus	-8.9708***	26.8305***	-3.43	4.75
Terrorism Attack	-15.5329***	80.4342***	-3.43	4.75

##### B. Checking Model Assumption

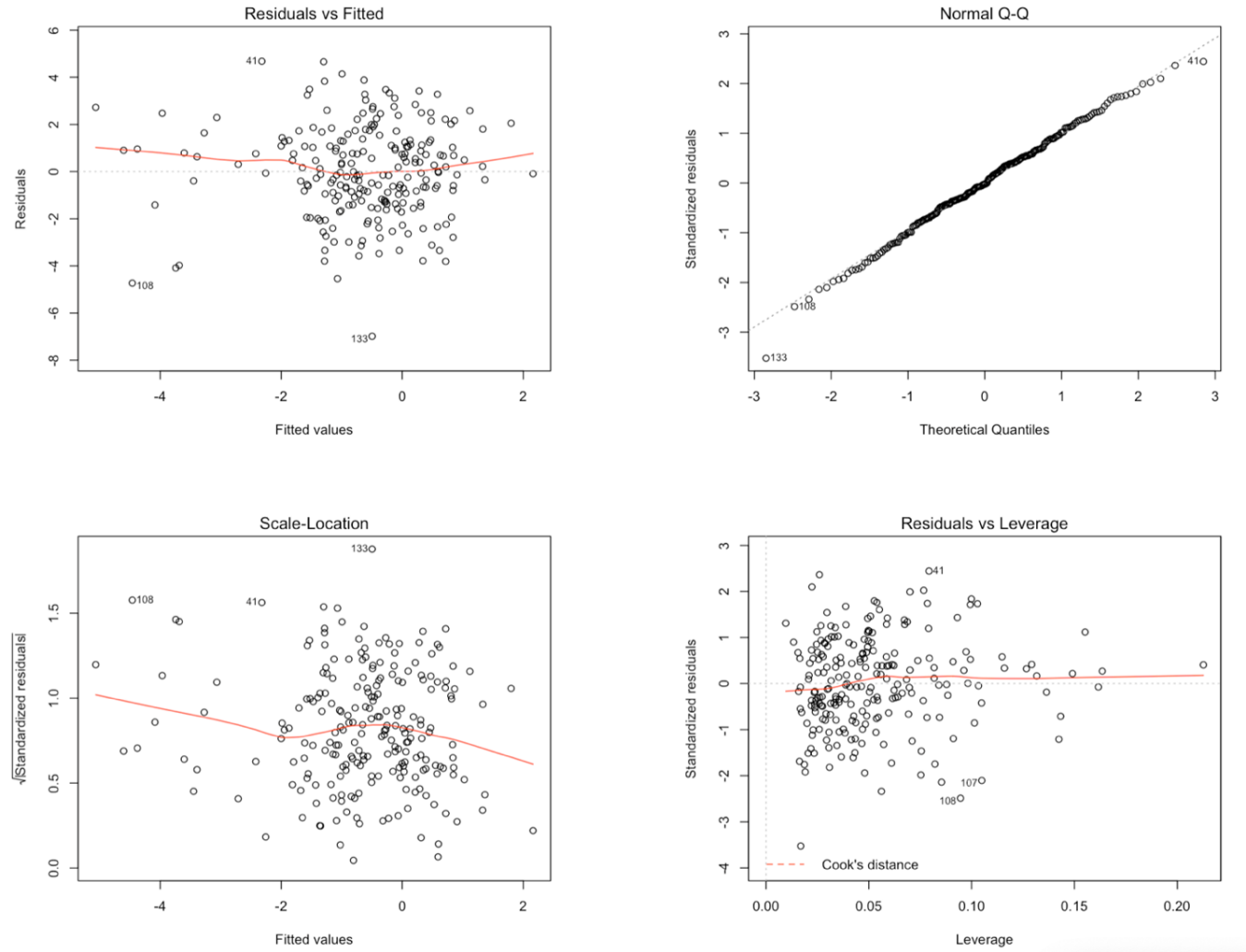
In this section, we want to check the assumption of the multi-linear regression model. It will mainly focus on whether there exist multicollinearity problems, whether there is a normal residual, and

whether there is homoscedasticity. The following matrix (figure 3) is the correlation matrix of the explanatory variables. We can see that most of the correlation between explanatory variables, however, Dow Jones Oil Index and S&P 500 have a relatively high correlation. All

the VIF of explanatory variables is between 1 to 10, which indicates there is not multicollinearity problem in the OLS model. We can see the following pictures (figure 1) to see the performance and diagnostics of the model. The residual and fitted plot shows that the residual is around 0, also, the Quantile-quantile plot shows the residuals almost follow a normal distribution, which meets the

assumption of the model. Also, by looking at the Scale-location plot, we found a nearly horizontal line in the plot, which suggests the model does meet the assumption of equal variance (homoscedasticity). The Residual vs Leverage plot does not show any extreme value at the upper right corner, which suggests that there is not any spot that can be influential to the regression line.

Figure 1: Diagnostic Plot of OLS



### C. Model Selection

In this section, we will first build up the full model by using all the explanatory variables by using the ordinary least square (OLS) method. Then the research will do the dimension reduction by using the backward Akaike information criterion (AIC) method and least absolute shrinkage and selection operator (LASSO) method. The section will also show the performance of each model by using the adjusted R square and the root means square error (RMSE). The AIC model will delete one variable at one time, to achieve a better model with a lower AIC. Backward AIC final model through

the stepwise-backward method eliminates multicollinearity problems. The selection process is showing as the following table (Table 3):

The final model achieved by the stepwise backward method is as following table (table 3) with a lowest AIC = 963. The estimated column show the estimated coefficients of the model. All the explanatory variables are linear significant. Moreover, the whole multi-regression model is in a significant level (F-statistics), and the adjusted R square is about 0.23, which means the model can explained 23% of variance of the WTI basis spread.

**Table 3: AIC selection Process**

	Adj R-Square	AIC	RMSE
DOW JONES OIL INDEX DIFF1	0.2251	970.0052	1.9929
SP500 INDEX DIFF1	0.2281	968.1909	1.9892
OCED OIL CONSUMPTION	0.2301	966.6417	1.9866
Terrorism_attack DIFF1	0.2313	965.3361	1.9850
FEDERAL RATE DIFF1	0.2326	963.9847	1.9834

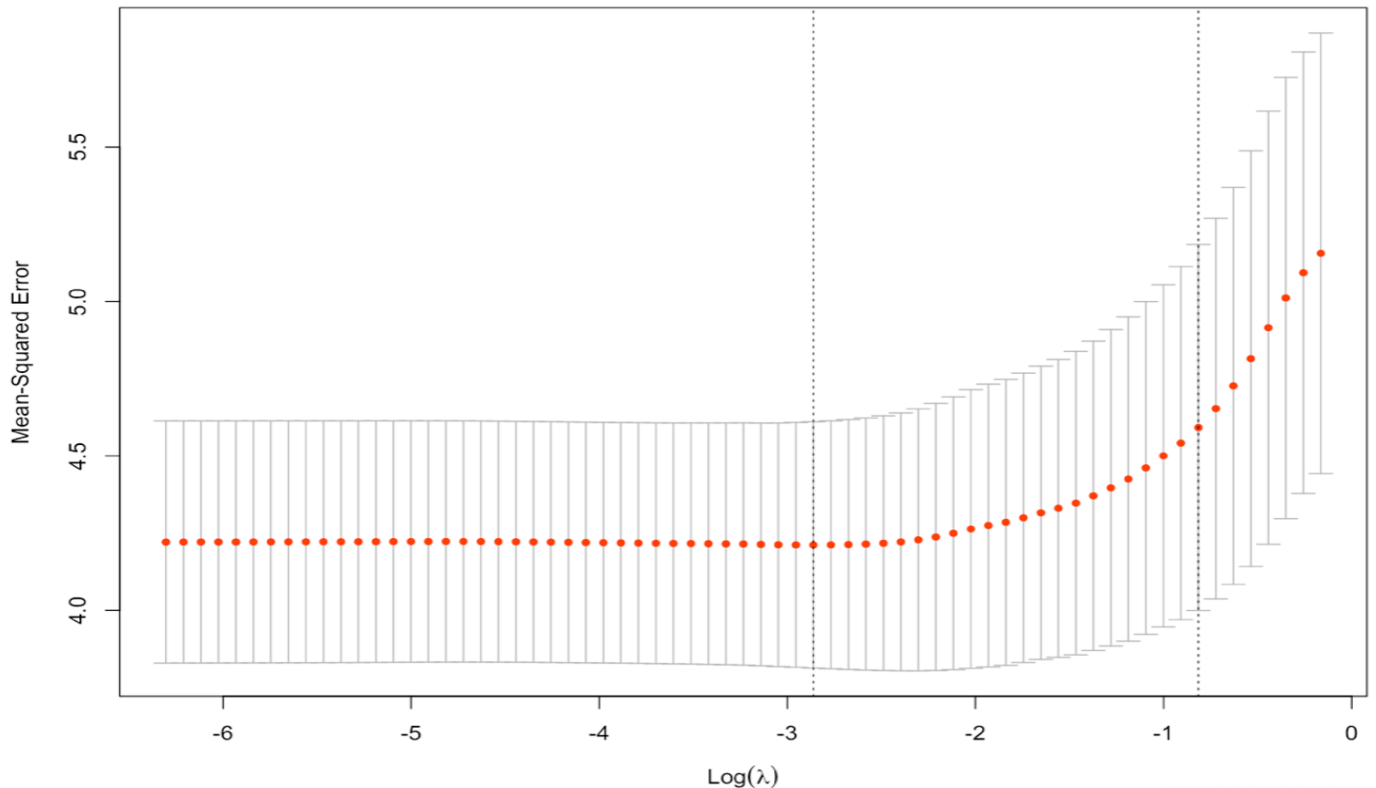
**Table 4: AIC Linear Regression Model**

	Estimated	Std. Error	Pr(> t )	Signif.
NYM1CNCN INDEX	1.358e-06	7.771e-07	0.081997	.
OPEC TOTAL SURPLUS DIFF1	-5.044e-04	2.365e-04	0.034016	*
DOLLAR INDEX DIFF1	-1.273e-01	6.335e-02	0.045692	*
GDP GROWTH RATE	3.374e-01	9.352e-02	0.000382	***
CPI US	4.838e-01	1.194e-01	7.02e-05	***
OPEN INTERESTS OIL DIFF1	-5.976e-06	2.134e-06	0.005550	**

Regression Summary, Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1, Multiple R-squared: 0.2529, Adjusted R-squared: 0.2326, F-statistic: 12.41, p-value: 4.915e-12

The figure (figure 2) shows the model selection process of the LASSO regression, each lambda value will have the corresponding whole path of coefficients. The cross validation process will pick up the best model with lowest MSE. The model shown the relationship between the value of logarithm of lambda and the mean squared

error (MSE) of the model. As a result, we can directly observe the best model is when the value of logarithm of lambda is near -3. The table (table 5) shows the coefficient of variables in LASSO model. We can see some variable have been eliminated.

**Figure 2: LASSO Selection Process**

**Table 5: LASSO Model Coefficients**

Variables	Estimated
OCED OIL CONSUMPTION	2.274374e-02
NYM1CNCN INDEX	1.071267e-06
OPEC TOTAL SURPLUS DIFF1	-4.130028e-04
DOLLAR INDEX DIFF1	-8.878033e-02
SP500 INDEX DIFF1	.
GDP GROWTH RATE	3.212513e-01
CPI US	4.252777e-01
FEDERAL RATE DIFF1	-1.162570e-01
DOW JONES OIL INDEX DIFF1	.
OPEN INTERESTS OIL DIFF1	-4.952636e-06
Terrorism_attack DIFF1	7.457815e-04

## V. DISCUSSION

By comparing the adjusted R square and the root mean square error (RMSE), we can see the performance of each model. The model with a higher adjusting R square can explain more variability of the response variable. By comparing three different models, the LASSO model has the highest adjusting R square equal to 0.256, which indicates it can explain almost 26% of the variability of the basis spread between the WTI crude oil spot price and WTI 4-month futures price. Moreover, the model with a lower root mean square error means the model is accurate. By comparing the RMSE, the LASSO model also has the lowest RMSE equal to 1.9529. As a result, the LASSO model has the lowest adjusting R square and lowest RMSE, we should choose the LASSO model among these three to predict the basis spread in the future.

Following figure (figure 3) show the time series of predicted basis spread by using LASSO model and the real basis spread. We can see the predicted basis is less fluctuated than the real basis spread. This is fair enough because the LASSO model only explains about 26% of the variability of real basis. Because the original data set are time-series data, by removing the random effect, all the variables used in the predicted models are stationary. As a result, it is hard to have an adjusted R square which is higher than 0.7, actually, according the empirical rule of other previous time-series research, a predicted model with an adjusted R square higher than 25% can be considered as a quite good model.

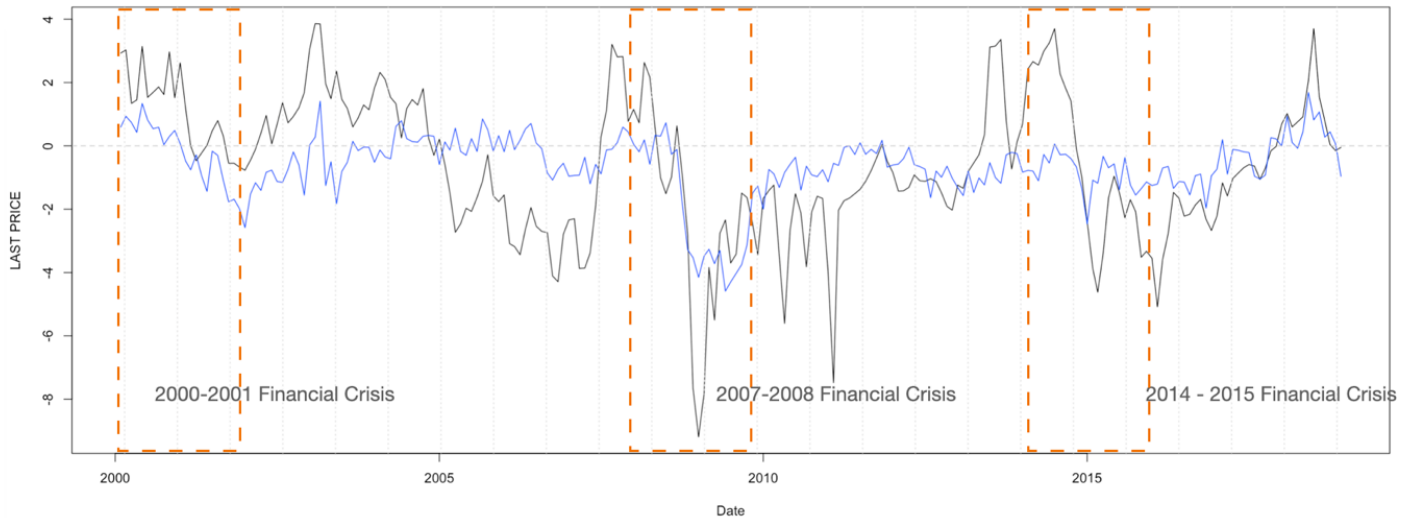
Moreover, by looking close to the figure, we can find the LASSO regression model does not do well in predicting the relatively large spread. Starting at the beginning, the real basis spread and predicted basis is moving in the roughly same direction at the same time. However, at around 2003 - 2004, the predicted basis start to move behind the real basis spread and fluctuate around 0 after 2005. At this period, the predicted model becomes inaccurate, even sometimes the predicted basis spread move in the opposite direction of the real basis. Around the 2008 financial crisis, the movement of both is highly correlated. Later it becomes inaccurate and unregulated again. Then around the 2014 financial crisis, the movement of both basis is similar again. There is three main financial crisis in the 21 century. The first time is the 2001–2002 Argentine Economic Crisis, the second time is the 2007–2009 Global Financial Crisis, and the last time is the 2014 Russian Financial Crisis. The figure shows the highly similar co-movement of the real basis and predicted basis.

One possible reason is that most of the significant explanatory variables can be easily affected by the financial crisis and put these influences in the same direction of the basis. Moreover, basis spread has a large uncertainty, it depends on both the WTI crude oil spot market and the WTI crude oil futures market, some of the variability cannot be explained by the linear model. As a result, the predicted basis becomes more accurate and reasonable if there is a large event that happened in the market, such as the financial crisis. Whereas, when there are not influential events happening in the market, the predicted basis become more like random work.

**Table 5: Model Comparison**

	Adjusting R Square	MSE	RMSE
OLS FULL Model	0.222	3.989727	1.9974
AIC BACKWARD Model	0.233	3.933876	1.9834
LASSO Model	0.256	3.813779	1.9529

**Figure 3: Best Prediction Model with Major Financial Crisis**



## VI. CONCLUSION

In this research, we considered the 11 different variables from a different aspect, including demand factors, supply factors, economic condition, financial market performance, and the dollar strength, to predict the basis spread between WTI OK Crushing crude oil spot price and the 4-month WTI OK Crushing crude oil futures price. We start with the full multi-linear regression model to predict the basis spread and then do the feature reduction to reduce the effect of overfitting. The final result shows that the LASSO model has the best model performance, which can explain 26% of the basis variability with the lowest root mean square error. The adjusting R square is a little bit higher than the empirical value of 25% for times-series and stationary explanatory variables. As a result, the LASSO regression model can be considered as a quite good model. By comparing the predicted basis and real basis, we have seen that the predictive model is powerful to predict the movement direction or even the basis spread when there is a significant event in the world market. The time period during the three main financial crises in 21 century shows strong evidence that the predicted model is more accurate when there are influential events. On the other side, WTI crude oil basis spread is more predictable when there is an influential event—However, if there are not influential factors in the market, the predicted basis spread becomes inaccurate and more like random work. Due to COVID-19, the financial market has been suffered for some time. As a result, further research can conduct more research about how basis of crude oil is driven by the COVID-19. In this research, we used the data ranged from January 2000 to December 2018 excluding the driven power by COVID-19. However, by applying the LASSO predictive model to the data during COVID-19, we have found the huge error and smallest adjusting R square. It may indicate the further research should separate their data and observe the COVID-19' effect respectively. Moreover, further research also can be conducted to use the non-linear way to predict the basis of WTI crude oil.

## ACKNOWLEDGMENT

Thanks for financial support from the SpF funding of Wenzhou-Kean University (WKU201920026) to Yichao Dai.

## REFERENCES

- [1] Ajifowoke, M. G. (2020). Weekly economic index: MPC raises cash reserve requirement, coronavirus affects oil prices.
- [2] Algieri, B. (2014). The influence of biofuels, economic and financial factors on daily returns of commodity futures prices. *Energy Policy*, 69, 227–247. <https://doi.org/10.1016/j.enpol.2014.02.020>
- [3] Allen, M. P. (Ed.). (1997). The problem of multicollinearity. In *Understanding Regression Analysis* (pp. 176–180). Springer US. [https://doi.org/10.1007/978-0-585-25657-3\\_37](https://doi.org/10.1007/978-0-585-25657-3_37)
- [4] Amendola, A., Candila, V., Candila, V., Scognamillo, A., & Scognamillo, A. (2017). On the influence of US monetary policy on crude oil price volatility. *Empirical Economics*, 52(1), 155–178. doi:10.1007/s00181-016-1069-5
- [5] An, H., Gao, X., Fang, W., Ding, Y., & Zhong, W. (2014). Research on patterns in the fluctuation of the co-movement between crude oil futures and spot prices: A complex network approach. *Applied Energy*, 136, 1067–1075. <https://doi.org/10.1016/j.apenergy.2014.07.081>
- [6] Basistha, A., & Kurov, A. (2015). The impact of monetary policy surprises on energy prices. *The Journal of Futures Markets*, 35(1), 87–103. doi:10.1002/fut.21639
- [7] Bu, H. (2011). Price Dynamics and Speculators in Crude Oil Futures Market. *Systems Engineering Procedia*, 2, 114–121. <https://doi.org/10.1016/j.sepro.2011.10.014>
- [8] Frondel, M., & Horvath, M. (2019). The U.S. fracking boom: Impact on oil prices. *The Energy Journal* (Cambridge, Mass.), 40(4), 191. doi:10.5547/01956574.40.4.mfro
- [9] Hanck, C., & Czudaj, R. (2015). Nonstationary-volatility robust panel unit root tests and the great moderation. *ASTA Advances in Statistical Analysis*, 99(2), 161–187. Retrieved from:
- [10] Hill, R., Griffiths, W., Lim, G. (2017). *Principles of Econometrics*, 5Th Edition. Chapter 12: Regression With Time-Series Data: Non-stationary Data. LCCN 2017056927 (eBook).
- [11] Holwerda, D., & Scholtens, B. (2016). The financial impact of terrorist attacks on the value of the oil and gas industry: An international review. (pp. 69–80). Cham: Springer International Publishing. doi:10.1007/978-3-319-32268-1\_5



- [12] Homepage—U.S. Energy Information Administration (EIA). (2020). Retrieved November 22, 2020, from <https://www.eia.gov/index.php>
- [13] Jones, C. M., & Kaul, G. (1996). Oil and the Stock Markets. *The Journal of Finance*, 51(2), 463–491. <https://doi.org/10.1111/j.1540-6261.1996.tb02691.x>
- [14] Kaufmann, R. K., & Ullman, B. (2009). Oil prices, speculation, and fundamentals: Interpreting causal relations among spot and futures prices. *Energy Economics*, 31(4), 550–558. <https://doi.org/10.1016/j.eneco.2009.01.013>
- [15] Kollias, C., Kyrtou, C., & Papadamou, S. (2013). The effects of terrorism and war on the oil price-stock index relationship. *Energy Economics*, 40, 743–743. <https://doi.org/10.1016/j.eneco.2013.09.006>
- [16] Miao, H., Ramchander, S., Wang, T., & Yang, D. (2017). Influential factors in crude oil price forecasting. *Energy Economics*, 68, 77–88. <https://doi.org/10.1016/j.eneco.2017.09.010>
- [17] Mordkoff, J. T. (n.d.). The Assumption(s) of Normality. 6.
- [18] Pham, C. S. (2018). Multiple regression model for cotton price returns: Analysis of the impact of weather, oil price return, and China's economy. <https://aaltodoc.aalto.fi:443/handle/123456789/33972>
- [19] Sadorsky, P. (1999). Oil price shocks and stock market activity. *Energy Economics*, 21(5), 449–469. [https://doi.org/10.1016/S0140-9883\(99\)00020-1](https://doi.org/10.1016/S0140-9883(99)00020-1)
- [20] Stoll, H. R., & Whaley, R. E. (2015). Commodity Index Investing and Commodity Futures Prices (SSRN Scholarly Paper ID 2693084). Social Science Research Network. <https://papers.ssrn.com/abstract=2693084>
- [21] Tufféry, S. (2011). *Data Mining and Statistics for Decision Making*. John Wiley & Sons, Incorporated. Retrieved from: <http://ebookcentral.proquest.com/lib/kean/detail.action?docID=792450>
- [22] Wang, Y., & Wu, C. (2012). Energy prices and exchange rates of the U.S. dollar: Further evidence from linear and nonlinear causality analysis. *Economic Modelling*, 29(6), 2289–2297. <https://doi.org/10.1016/j.econmod.2012.07.005>