# TIME SERIES ANALYSIS OF WTI CRUDE OIL PRICES

Pauline Mwaura

April 15, 2025

**Abstract**

This project employs an extensive time series analysis approach to forecast WTI Crude Oil prices, focusing on understanding the dynamics influencing oil prices and developing accurate predictive models for informed decision-making. The methodology encompasses various statistical tests and modeling techniques.

The stationarity of the WTI Crude Oil price time series is assessed using both the Augmented Dickey-Fuller (ADF) test and the Kwiatkowski-Phillips-Schmidt-Shin (KPSS) test. This initial step is crucial in determining the stability of the statistical behavior over time.

Autoregressive Integrated Moving Average (ARIMA) models are utilized for forecasting, with model parameters determined through autocorrelation function (ACF) and partial autocorrelation function (PACF) assessments. The selected ARIMA model, ARIMA(2, 1, 2), is indicative of second-order differencing, incorporating autoregressive and moving average terms.

To account for time-varying volatility, Generalized Autoregressive Conditional Heteroskedasticity (GARCH) and Autoregressive Conditional Heteroskedasticity (ARCH) models are considered. These models capture the clustering of volatility observed in financial returns.

Extreme Value Analysis (EVA) is applied to identify and model extreme events in the WTI Crude Oil price data. The Block Maxima Method is employed, dividing the time series into non-overlapping blocks to assess the distribution of block maxima.

In conclusion, this comprehensive approach contributes significant insights into the volatility, trends, and extreme events associated with WTI Crude Oil prices. The project's outcomes are valuable for stakeholders in the energy and financial sectors, providing a holistic understanding for forecasting and risk assessment.

Keywords: WTI Crude Oil, ARIMA model, GARCH, ARCH, Extreme Value Analysis, Stationarity Testing, ADF, KPSS, ACF, PACF, Time Series Analysis, Forecasting, Financial Markets.

# Contents

# Chapter 1

# Introduction

This report presents a comprehensive time series analysis of West Texas Intermediate (WTI) crude oil prices from 2019 to the present day, a period marked by significant global economic changes and geopolitical events impacting the energy sector. The objective of this analysis is to understand the dynamics of crude oil prices through various statistical models, offering insights that are crucial for economic forecasting, policy making, and investment decision-making. The relevance of this study is underscored by previous research in the field, such as the work by Hamilton (2009) and Kilian (2009), who demonstrated the profound impact of oil price fluctuations on the global economy, underscoring the need for accurate predictive models in this domain.

The first phase of the study, Exploratory Data Analysis (EDA), lays the groundwork by identifying patterns, trends, and anomalies in the WTI crude oil prices. This phase is crucial for determining the appropriate modeling approach, as it reveals underlying structures and behaviors in the time series data. Following EDA, the report delves into the fitting of stationary time series models including Autoregressive (AR), Moving Average (MA), and Autoregressive Moving Average (ARMA) models. These models are pivotal in understanding the short-term dependencies and regular patterns in the data, as evidenced in studies like Box and Jenkins (1976), which revolutionized time series analysis with these methodologies.

The analysis then progresses to more complex non-stationary models, namely Generalized Autoregressive Conditional Heteroskedasticity (GARCH), Autoregressive Conditional Heteroskedasticity (ARCH), and Seasonal Autoregressive Integrated Moving-Average (SARIMA) models. These models are particularly adept at capturing the volatility and seasonal variations in the crude oil market, features that are often observed in financial time series data as highlighted in the seminal works of Engle (1982) and BOLLERSLEV (1986) on ARCH and GARCH models. The report culminates with an Extreme Value Analysis (EVA), an approach critical for assessing the risk of rare but high-impact events in oil prices, as demonstrated by the 2008 financial crisis and the 2020 pandemic-induced market shocks. This comprehensive methodological approach ensures a robust analysis of WTI crude oil prices, offering valuable perspectives for a range of stakeholders in the energy sector.

The primary data source for this analysis is Yahoo Finance, utilizing the quantmod package in R. this is what it looks like:

The dataset includes daily historical prices of WTI Crude Oil futures, going back five years. The ticker symbol used for retrieval is "CL=F," representing the continuous contract for WTI Crude Oil.

The structure for this data is achieved with the following R code:

```
# Specify the ticker symbol for WTI Crude Oil
oil_symbol <- "CL=F"
# View the structure of the loaded data
str(get(oil_symbol))
```

Results:

```
> str(get(oil_symbol))
An xts object on 2019-01-02 / 2024-01-24 containing:
    Data:    double [1276, 6]
     Columns: CL=F.Open, CL=F.High, CL=F.Low, CL=F.Close, CL=F.Volume ... with 1 more column
    Index:   Date [1276] (TZ: "UTC")
    xts Attributes:
```
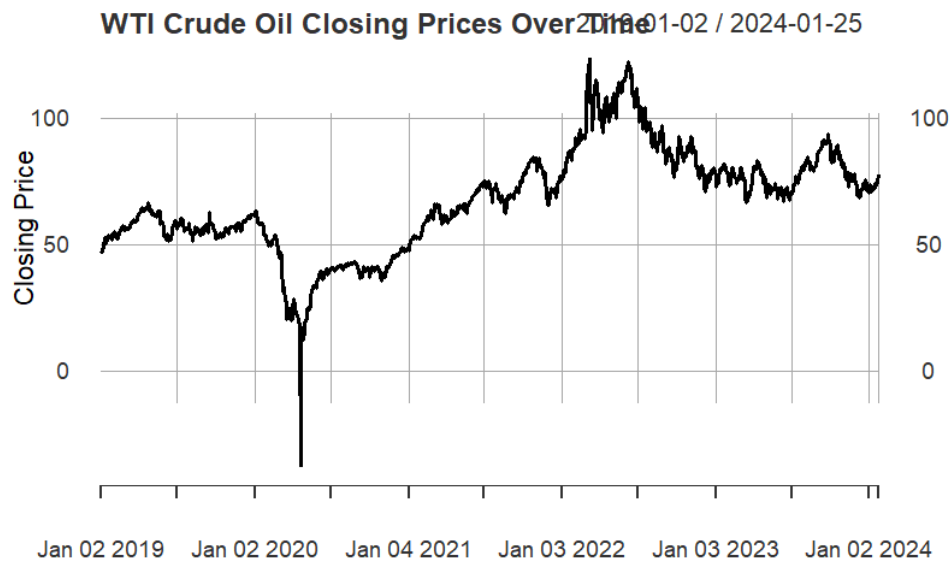
Figure 1.1: WTI Crude closing prices in the past 5 years

```
$ src    : chr "yahoo"
$ updated: POSIXct[1:1], format: "2024-01-25 19:25:05"
>
```

A breakdown:

- The data contains 1276 rows and 6 columns.

- Each column represents different attributes such as Open, High, Low, Close, Volume, and an additional unnamed column. The column names are "CL=F.Open," "CL=F.High," "CL=F.Low," "CL=F.Close," "CL=F.Volume," and an unnamed column. These correspond to the different aspects of the financial instrument's price and volume.

- The index represents the date range from January 2, 2019, to January 19, 2024, in UTC timezone

- The attributes include information about the source of the data ("yahoo") and the last update time.

# Chapter 2

# Literature Review

Yusof et al. [10] endeavored to construct an ARIMA model utilizing monthly data on Malaysian crude oil production with the aim of forecasting output for the upcoming three months. They successfully formulated an ARIMA model with parameters (1,0,0). However, this investigation revealed a notable gap in the existing literature regarding the modeling of crude oil production. Most prior research has primarily focused on crude oil exports and the volatility of crude oil prices.

Fluctuations in oil prices are influenced by a complex interplay of various factors, encompassing economic considerations, as well as dynamics associated with the supply and demand mechanism and their intricate interactions. The global oil supply's magnitude is shaped by factors such as productive capacity and the allocation of production quotas among worldwide producers, both within and outside OPEC, as highlighted by Kosakowski (2020).

Prices also experience influence from geopolitical shifts on both the global and regional scale within oil-producing nations and key oil-consuming countries. Any occurrence of political instability, such as the outbreak of war, rebellion, riots, revolution, or coup, poses a threat to production, transportation, distribution routes, or consumption hubs. This volatility impacts production volumes, the overall global supply, demand levels, and consequently, short-term oil prices. An illustrative instance is the situation in July 2008 when the barrel price surged to 128 dollars due to instability and consumer apprehension regarding potential war in Afghanistan and Iraq, as noted by Ganti (2020).

Various elements might be characterized as behavioral, contributing to the fluctuations in oil prices by aligning with the actions of customers and financial investors in their choices to transact oil and gas contracts. These determinants are anchored in trust, expectations, speculation, and the pursuit of profit. Additionally, these decisions are influenced by geopolitical conditions and assumptions about the anticipated direction of price changes, whether ascending or descending (Quan, 2014).

# Chapter 3

# DATA ANALYSIS

We will start with the Data Collection, then we perform Exploratory Data Analysis, Stationarity Testing, fit both stationary and non-stationary models and perform extreme value analysis. We then select a suitable model to use, then we do model Validation.

## 3.1 Data Collection

The primary dataset for this study is the historical daily closing prices of WTI Crude Oil. Data is sourced from Yahoo, which is a reputable financial database, spanning from January 1, 2019, to the current date. The choice of this period allows for a comprehensive analysis of recent trends and fluctuations in WTI Crude Oil prices. the plot: 3.8

## 3.2 Exploratory Data Analysis (EDA)

Before delving into model development, an exploratory data analysis (EDA) is conducted to gain insights into the characteristics of the time series.

### 3.2.1 Visualization

**1. Histogram**
Histograms can reveal the shape of the distribution, whether it's symmetric, skewed, or has multiple peaks.
If your data is not normally distributed, you might explore transformations to achieve a more normal distribution.
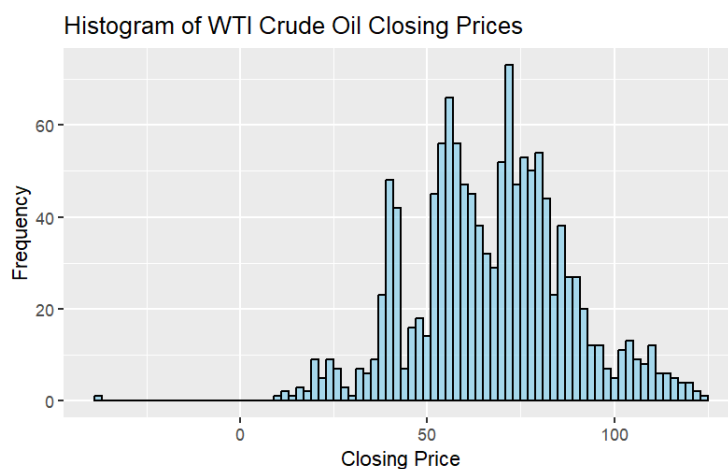Our graph is assymetric and the data is multimodal



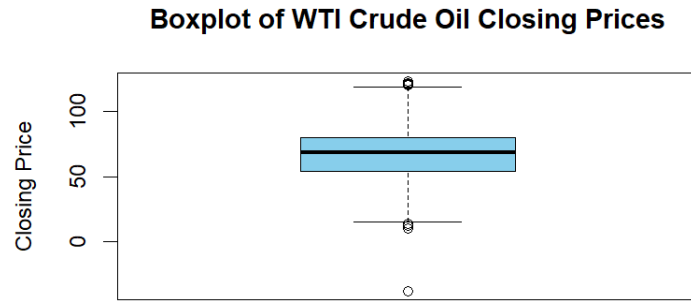Figure 3.1: histogram for WTI closing price

## 2. Boxplots



**Boxplot of WTI Crude Oil Closing Prices**

Figure 3.2: boxplot for WTI clossing price

The boxplot is vertically oriented
It is notable that there is one outlier present in the dataset.
The median (Q2) is positioned at around 60 indicating the central tendency of the dataset

**Summary statistics**

Table 3.1: Summary Statistics

| Index | CL=F.Close |
|---|---|
| Min. | -37.63 |
| 1st Qu. | 54.16 |
| Median | 68.36 |
| Mean | 67.31 |
| 3rd Qu. | 80.35 |
| Max. | 123.70 |
| Std Dev | 20.4637 |

## 3.3 Stationarity Testing

Stationarity is a crucial assumption for time series modeling. We employ the Augmented Dickey-Fuller (ADF) and Kwiatkowski-Phillips-Schmidt-Shin (KPSS) tests to assess the stationarity of the time series. If the data is found to be non-stationary, differencing techniques are applied to achieve stationarity. The differenced series plot:

### 3.3.1 ADF Test

Results of the **ADF test**:

```
    Augmented Dickey-Fuller Test

data:  closing_prices
Dickey-Fuller = -1.765, Lag order = 10, p-value =
0.6778
alternative hypothesis: stationary
```

**Interpretation**: Null Hypothesis (H0): The null hypothesis of the ADF test is that the time series has a unit root and is non-stationary. Alternative Hypothesis (Ha): The alternative hypothesis is that the

time series is stationary (lacks a unit root). The test statistic is -1.765 in this case.The more negative the test statistic, the stronger the evidence against the null hypothesis. (non-stationary time series) Lag Order: The lag order is 10. It indicates the number of lags used in the test. Lag order is chosen to minimize information criteria, and it plays a role in controlling for autocorrelation in the data.

p-value: The p-value associated with the test statistic is 0.6778. The p-value is crucial for determining the significance of the test. In this case, the p-value is relatively high (greater than the common significance level of 0.05).it's larger than a common significance level like 0.05. This suggests that you do not have enough evidence to reject the null hypothesis of a unit root, indicating non-stationarity.

### 3.3.2  kpss test

(Kwiatkowski-Phillips-Schmidt-Shin) test . Unlike the Augmented Dickey-Fuller (ADF) test, which tests for the presence of a unit root (non-stationarity), the KPSS test focuses on detecting trends in the data. The null hypothesis of the KPSS test is that the data is stationary around a deterministic trend. The test involves estimating the trend component in the data and examining whether the residuals (detrended series) are stationary. The KPSS test provides valuable information about the stationarity properties of a time series:

If the p-value is less than a chosen significance level (e.g., 0.05), you may reject the null hypothesis, suggesting that the data is not stationary around a deterministic trend.

If the p-value is greater than the significance level, you may fail to reject the null hypothesis, indicating that the data is stationary around a deterministic trend. Resultsof the **KPSS Test**

```
   KPSS Test for Level Stationarity

data:  closing_prices
KPSS Level = 8.7445, Truncation lag parameter = 7,
p-value = 0.01
```

**Interpretation** KPSS Level: The test statistic for the level component is 8.7445.

Truncation Lag Parameter: The truncation lag parameter is specified as 7.

p-value: The p-value associated with the test is 0.01. With a p-value of 0.01, if the significance level is set at 0.05, you would reject the null hypothesis. This suggests that the data is not stationary around a deterministic level, indicating the presence of a trend in the series.

Hence our Time Series Data is non-stationary and we need to difference it to achive stationarity before chosing a model.

## 3.4   Differencing and ACF/PACF

we difference our time series Then we check for missing values and omit them. Next we Compute ACF
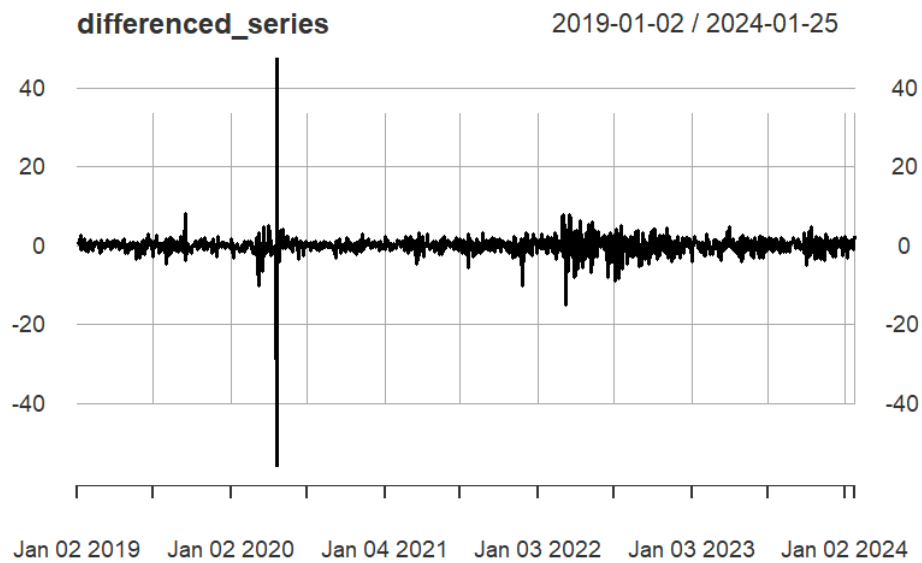


Figure 3.3: Differenced series

and PACF for the differenced series without missing values

**ACF** • A significant spike at lag 0 is expected and represents the autocorrelation of the series with

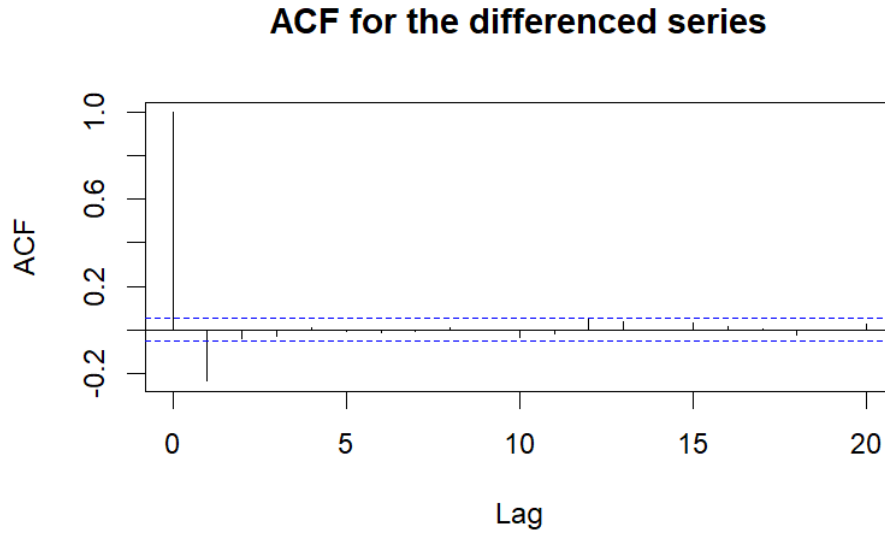## ACF for the differenced series



Figure 3.4: ACF for the Differenced series

itself at the same time point. • A significant spike at lag 1 suggests a strong positive autocorrelation at the first lag, which may indicate the presence of an autoregressive (AR) component.

**PACF** • A significant spike at lag 1 in the PACF suggests a direct relationship between observations
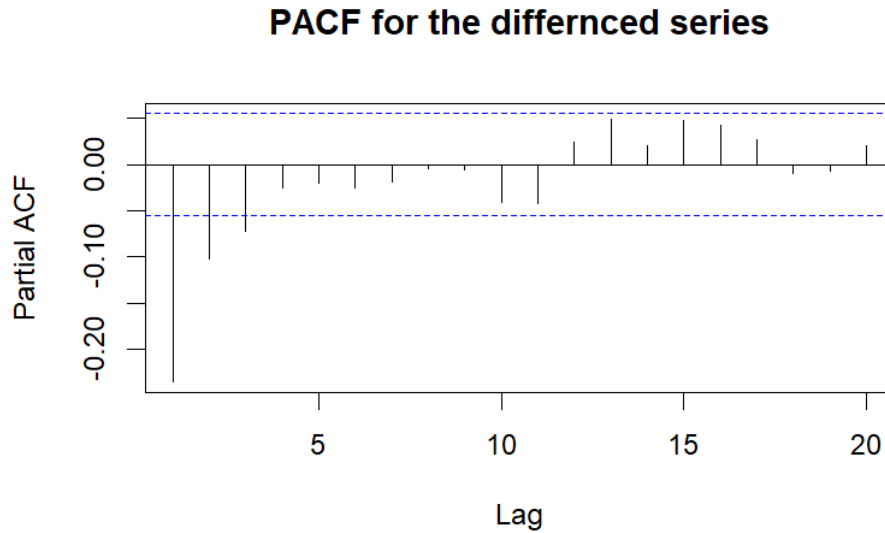
## PACF for the differnced series



Figure 3.5: PACF for the Differenced series

one time unit apart. This is common when differencing a series once. • Significant spikes at lag 2 and lag 3 in the PACF might suggest the potential presence of autoregressive terms beyond the first lag.

## 3.5   Fitting Stationary models

**Fitting AR** Autoregressive (AR) models are used in time series analysis when the current value of a variable is believed to be linearly dependent on its past values. In an AR model, the current observation is

modeled as a linear combination of its past values, and the term "autoregressive" refers to the dependence on the past values of the variable itself.

based on the PACF 3.5, we considered an AR of order 3, as it seems that the correlation is present up to lag 3. We use PACF to evaluate the AR model because the AR model is built on the (real) correlation between time spots.

```
arima(x = differenced_series_no_na, order = c(3, 0, 0))
```

```
Coefficients:
          ar1      ar2      ar3  intercept
      -0.2654  -0.1207  -0.0717     0.0240
s.e.   0.0279   0.0287   0.0279     0.0517
```

```
sigma^2 estimated as 7.245:  log likelihood = -3076.42,  aic = 6162.84
```

Interpretation: all three coefficients are negative, suggesting a negative correlation with the series' recent past values. intercept: 0.0240 - represents the mean of the differenced series. $\sigma^2$ estimated as 7.245 - estimate of the variance of the residuals. It provides a measure of the variability of the unexplained part of the time series after accounting for the AR components. Log likelihood = -3076.42 - how well the model explains the observed data. AIC (Akaike Information Criterion) = 6162.84 Overall, this ARIMA(3, 0, 0) model suggests that the current value of your time series is influenced by its values at lags 1, 2, and 3, with negative correlations.

one day prediction using the ar model Predicted value for the next day: -0.1443777 meaning that the model is predicting a decrease or negative change from the current value to the next day.

**Fitting MA**

```
  Coefficients:
          ma1  intercept
      -0.2842     0.0243
s.e.   0.0291     0.0540
```

```
sigma^2 estimated as 7.273:  log likelihood = -3078.94,  aic = 6163.88
```

The ma1 coefficient represents the impact of the lagged error term on the current observation. In this case:

For each one-unit increase in the lagged error term, we expect the current observation to decrease by approximately 0.2842 units.

one day prediction using ma Predicted value for the next day using MA: -0.1991262 predicting a decrease by 0.1991262

**Fitting ARMA** we use an ARMA of order(3,0, 1), based on AR(3), ma(1). the results:

```
arima(x = differenced_series_no_na, order = c(3, 0, 1))
```

```
Coefficients:
          ar1     ar2     ar3      ma1  intercept
      0.3762  0.0468  0.0019  -0.6453     0.0235
s.e.  0.1931  0.0603  0.0430   0.1912     0.0464
```

```
sigma^2 estimated as 7.232:  log likelihood = -3075.29,  aic = 6162.59
```

ar1 (0.3762): The coefficient for the first lag of the differenced series is positive.indicating that the current value of the series is influenced by its value at the previous time step (lag 1). ar2 (0.0468): for the second lag, also positive, suggesting a smaller influence from the value two time steps ago (lag 2). ar3 (0.0019): for the third lag, very close to zero, indicating a minimal influence from the value three time steps ago (lag 3).

**ma1** (-0.6453): The coefficient for the first lag of the residual series (the white noise term) is negative. It signifies that the model is incorporating information from the past error terms to predict the current value.

10

Intercept (0.0235): represents the constant term in the model, accounting for the mean of the differenced series. Sigma ($\sigma^2$ estimated as 7.232): an estimate of the variance of the residuals or the white noise term.

Log likelihood (-3075.29): A measure of how well the model explains the observed data. AIC (6162.59)

one day prediction using arma Predicted value for the next day using ARMA: -0.2945799 . expecting the current value to drop by 0.2945799 the next day.

### 3.5.1   What model to choose?

```
> AIC(ar)
[1] 6162.844
> AIC(ma1)
[1] 6163.882
> AIC(arma)
[1] 6162.587
>
```

Lower AIC values are preferred. so the best model to fit this data with stationary models is the ARMA model

## 3.6   Fitting Non-Stationary models

### 3.6.1   ARIMA

We can choose the arima order a number of two ways: 1. using the PACF and ACF of the differenced series Based on the observations in figure 3.4 and figure 3.5, • p (AR order): The order of the AR component might be considered as 1, given the significant spike at lag 1 in both ACF and PACF. • q (MA order): The order of the MA component could be initially considered as 1, given the significant spike at lag 1 in the ACF. hence, our preliminary ARIMA order might be something like ARIMA(1,1,1). or: we could use the auto.arima() function to automatically select the ARIMA model. it selects an ARIMA of order(1,1,1)

to validate the significance of the identified order, we then perform a box-ljung test. The Box-Ljung test is a statistical test used to assess whether there are significant autocorrelations in a time series at various lags. It helps in evaluating whether the residuals of a time series model exhibit serial correlation, which is an important aspect to consider when fitting models such as ARIMA (AutoRegressive Integrated Moving Average). So we extract the residuals and perform the test on them. results:

```
    Box-Ljung test

data:  model_residuals
X-squared = 34.492, df = 20, p-value = 0.02298
```

The test statistic is calculated as X-squared, and a low p-value (in this case, 0.02215) suggests evidence against the null hypothesis of no autocorrelation.

An acf of the residuals: A significant spike at lag 2 in the ACF plot indicates a correlation between
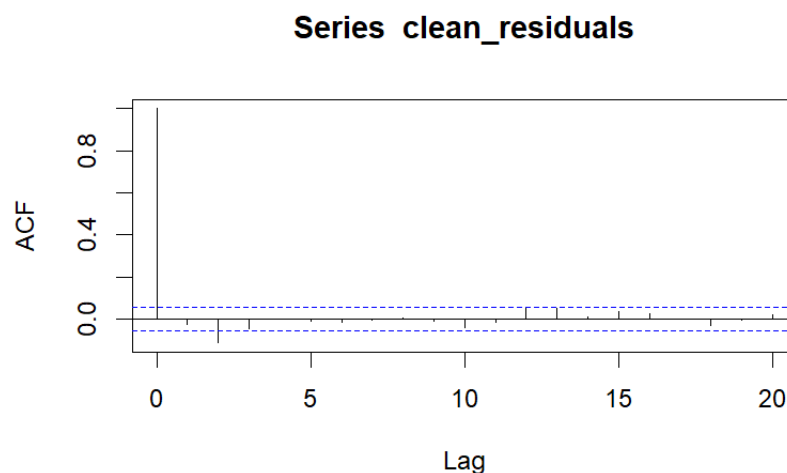
**Series  clean_residuals**



Figure 3.6: ACF for the Residuals

the residuals at the current time point and the residuals at the second lag.

given these findings,we might consider Refining our Model Order, as there is correlation in our residuals.

Choosing an ARIMA of order 2,1,2 no significant spikes in our ACF and our box-ljung test yields the following:

```
    Box-Ljung test

data:  model_residuals
X-squared = 15.934, df = 20, p-value = 0.7207
```

indicating, a higher p-value, ¿ 0.05, indicating no autocorrelation between the residuals hence our ARIMA model seems to be a good fit. the AIC for this ARIMA order is 6164.881

### 3.6.2 ARCH

ARCH (Autoregressive Conditional Heteroskedasticity) is a statistical model used to analyze and forecast volatility in time series data. ARCH models assume that the variance of the error term in a time series is not constant but rather varies over time. It incorporates lagged values of the squared residuals to capture this time-varying volatility. We started by fitting an ARCH model of order 1 (ARCH(1)).

- mu (Constant Term): 67.23358

- omega: 2.60378

- beta1 (ARCH Coefficient): 0.99412

The ARCH(1) model provides a good fit to the data based on diagnostic tests and goodness-of-fit measures. The significant p-values in the Weighted Ljung-Box tests indicate the absence of serial correlation in the standardized residuals:

- Lag[1]: Statistic = 1256, p-value = 0

- Lag[2*(p+q)+(p+q)-1][2]: Statistic = 1878, p-value = 0

- Lag[4*(p+q)+(p+q)-1][5]: Statistic = 3722, p-value = 0

low p-values (all zeros in this case) indicate rejection of the null hypothesis, suggesting evidence against serial correlation in the residuals, suggesting that the model adequately captures volatility patterns in the squared residuals:

- ARCH Lag[2]: Statistic = 897.1, p-value = 0

- ARCH Lag[4]: Statistic = 1959.2, p-value = 0

- ARCH Lag[6]: Statistic = 2838.5, p-value = 0

The low p-values (all zeros in this case) indicate rejection of the null hypothesis, suggesting evidence against the absence of autocorrelation in the squared residuals. These results indicate that the model provides a statistically valid representation of autocorrelation and volatility dynamics in the data.

### 3.6.3 GARCH

We fit GARCH(1,1) model

- estimated mean is 2.20908, with a standard error of 0.166384.

- omega (Intercept of Conditional Variance): The estimated intercept is 10.76773, with a standard error of 0.972090. The t-value is 11.077, and the p-value is 0, indicating a significant intercept.

- alpha1 (Coefficient for Lagged Squared Residual): The estimated coefficient is 0.50287, with a standard error of 0.042905. The t-value is 11.720, and the p-value is 0, indicating a significant coefficient.

- beta1 (Coefficient for Lagged Conditional Variance): The estimated coefficient is 0.49613, with a standard error of 0.023483. The t-value is 21.127, and the p-value is 0, indicating a significant coefficient.

Weighted Ljung-Box Test on Standardized Residuals: All p-values are non-significant, indicating no evidence of serial correlation. Weighted Ljung-Box Test on Standardized Squared Residuals: Tests for serial correlation in the squared residuals. Lag[1] and Lag[2*(p+q)+(p+q)-1][5] have non-significant p-values. In summary, the GARCH(1,1) model appears to be well-specified based on the significance of the parameters, lack of serial correlation in residuals, and non-significant results in various diagnostic tests. The model captures volatility clustering and adequately describes the conditional variance dynamics in the data.

the QQ plot is often applied to the standardized residuals to check whether the standardized residuals approximately follow a normal distribution.
the points are close to the line, suggesting that the residuals are approximately normally distributed. indicating that the GARCH model seems to adequately capture the conditional heteroskedasticity and volatility patterns in the data. This is validated by a lack of significant autocorrelation in the ACF in the standardized residuals
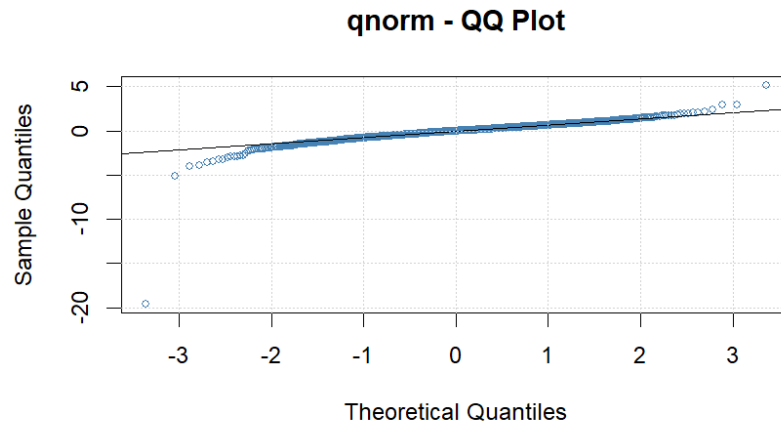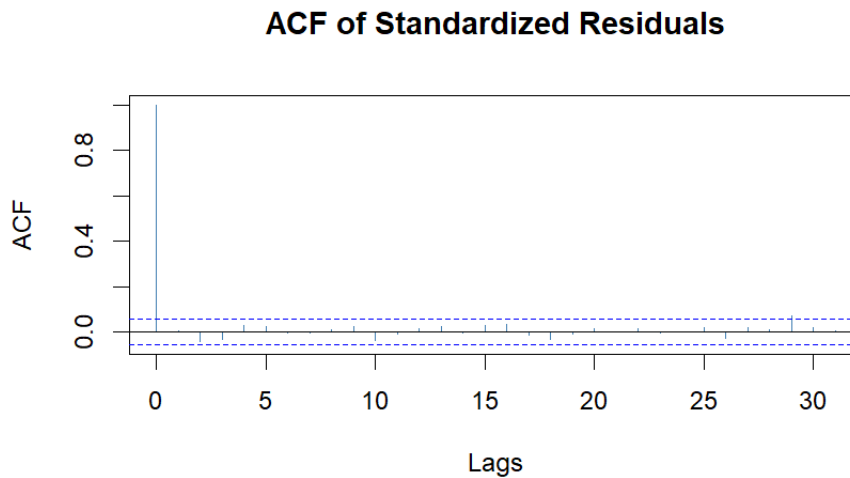


Figure 3.7: QQ plot



Figure 3.8: ACF of the standardized residuals

# Chapter 4

# Extreme Value analysis

Extreme Value Analysis (EVA) is a statistical technique employed in time series analysis to investigate and model the extreme values or outliers within a dataset. This method is particularly useful in understanding rare events or extreme conditions, which might have significant implications in various fields such as finance, environmental science, and engineering. In the context of time series analysis, extreme values refer to unusually high or low observations that deviate significantly from the typical behavior of the series.

Fitting a GEV distribution to block maxima
results:

```
fevd(x = block_maxima)
```

[1] "Estimation Method used: MLE"


 Negative Log-Likelihood Value:  4.62668


 Estimated parameters:
   location        scale        shape
 2.18952670   0.26465111 -0.03273685

 Standard Error Estimates:
  location        scale       shape
0.06737862 0.04912427 0.17623757

 Estimated parameter covariance matrix.
             location         scale          shape
location  0.004539879  0.001224986 -0.004615555
scale     0.001224986  0.002413194 -0.003280437
shape    -0.004615555 -0.003280437  0.031059682

 AIC = 15.25336

 BIC = 18.24056

**Estimated Parameters:** Location: The location parameter (2.1895) represents the location of the distribution.It indicates the threshold above which extreme events are considered. Scale: the scale parameter (0.2647) characterizes the spread or variability of the extreme values. It essentially measures the standard deviation of the distribution. Shape: The shape parameter (-0.0327) determines the shape of the distribution. For a negative shape parameter in the GEV distribution, it indicates a slight tendency towards a lighter tail.

**Standard Error Estimates:** These values (0.0674 for location, 0.0491 for scale, 0.1762 for shape) represent the standard errors associated with the estimated parameters. Lower standard errors generally indicate more precise parameter estimates.

**Estimated Parameter Covariance Matrix:** provides insights into the relationships between the estimated parameters. For example, a negative covariance between location and shape indicates that as one parameter increases, the other tends to decrease.

To predict extreme values using the fitted Generalized Extreme Value (GEV) distribution:

```
    # Generate quantiles for a range of probabilities
probabilities <- c(0.95, 0.99, 0.999)
predicted_quantiles <- qgev(probabilities, loc = 2.18952670, scale = 0.26465111, shape = -0.03273685
```

These predicted quantiles obtained represent the estimated quantiles at different probabilities from the Generalized Extreme Value (GEV) distribution.
Quantile at Probability 0.95: 2.938585 - this means that, based on the fitted GEV distribution, there is a 95% probability that an extreme value will be less than or equal to approximately 2.9386.

Quantile at Probability 0.99: 3.319725
For a higher probability of 99%, the extreme value is estimated to be less than or equal to approximately 3.3197.
Quantile at Probability 0.999: 3.825598
At an even higher probability of 99.9%, the extreme value is estimated to be less than or equal to approximately 3.8256.

# Conclusion

In this thorough analysis of WTI crude oil prices time series data, we explored various models tailored to different characteristics. We considered models for stable patterns, such as Autoregressive (AR), Moving Average (MA), and Autoregressive Moving Average (ARMA), models accommodating changing patterns, like Autoregressive Integrated Moving Average (ARIMA), and those addressing volatility, including Autoregressive Conditional Heteroskedasticity (ARCH) and Generalized Autoregressive Conditional Heteroskedasticity (GARCH). Essential steps included differencing the data to achieve stationarity, examining patterns using the Autocorrelation Function (ACF) and Partial Autocorrelation Function (PACF), and selecting models based on the Akaike Information Criterion (AIC). Additionally, we employed Extreme Value Analysis (EVA) with the Generalized Extreme Value (GEV) distribution to predict extreme events in the context of crude oil prices. Among all the models assessed, the ARMA model particularly stood out for effectively capturing patterns in the WTI crude oil prices time series. The selection of ARMA was driven by its alignment with the observed characteristics of the data, providing a balanced mix of simplicity and explanatory power.

# Bibliography

**REFERENCES**

1.Ganti, A. (2020). How OPEC (and Non-OPEC) Production Affects Oil Prices. Retrieved from https://www.investopedia.com/articles/investing/012216/how-opec-and-nonopec-production-affects-oil-prices.asp

2.Kosakowski, P. (2020). What Determines Oil Prices? Retrieved from https://www.investopedia.com/: https://www.investopedia.com/articles/economics/08/determining-oil-prices.asp

3.N. M., Yusof, R. S. A., Rashid, and, Z. Mohamed. Malaysia crude oil production estimation: an application of ARIMA model. In 2010 International Conference on Science and Social Research (CSSR 2010) (2010, December) (pp. 1255- 1259). IEEE.

4.Quan, L. (2014). Daqing crude oil price forecast based on the ARIMA model. BioTechnology.

5.Hamilton, J. D. (2009). Understanding Crude Oil Prices. The Energy Journal, 30(2), 179-206. https://doi.org/10.5547/ISSN0195-6574-EJ-Vol30-No2-9

6. Box, G.E.P. and Jenkins, G.M. (1976) Time Series Analysis: Forecasting and Control. 2nd Edition, Holden-Day, S. Francisco.

7. Engle, R. F. (1982). Autoregressive Conditional Heteroscedasticity with Estimates of the Variance of United Kingdom Inflation. Econometrica, 50(4), 987–1007. https://doi.org/10.2307/1912773

8. [Bollerslev, 1986] Bollerslev, T. (1986). Generalized autoregressive conditional heteroskedasticity. Jour- nal of Econometrics, 31(3):307–327