

Table of Content

CHAPTER 1	1
1.1 BACKGROUND	1
1.2 SCOPE	3
1.3 OBJECTIVE	3
1.3.1 Aim	3
1.3.2 Vision and Mission	4
1.4 STRUCTURE	4
1.4.1 Chapter 1 : Introduction	4
1.4.2 Chapter 2 : Theoretical Foundation	5
1.4.3 Chapter 3 : System Design	5
1.4.4 Chapter 4 : Solution Design	5
CHAPTER 2	6
2.1 SOFTWARE DEVELOPMENT LIFE CYCLE	6
2.2 VIRTUAL DATA ROOM	10
2.3 OIL AND GAS IN A RESERVOIR	11
2.3.1 Oil Formation	11
2.3.2 Gas Formation	12
2.3.3 Pressure and Temperature in Oil and Gas Formation	14
2.4 MACHINE LEARNING	15
2.5 TREE-BASED ALGORITHMS	16
2.6 DATA ANALYTICS PIPELINE	18
2.7 DATA IMPUTATION	18
2.7.1 Mechanism of Missingness	18
2.7.2 Central Value Imputation	19
2.7.3 Forward Filling	20
2.8 OUTLIERS	21
2.9 CORRELATIONS	22
2.9.1 Pearson Correlation	24
2.9.2 Spearman Correlation	24
2.10 FEATURE SELECTION	25
2.11 MODEL EVALUATION	26
2.11.1 Root Mean Square Error	26
2.11.2 Coefficient of Determinant	27
2.12 REST API	27
CHAPTER 3	29
3.1 PROBLEM STATEMENT	29
3.2 RELATED WORKS	29
3.3 PROPOSED SOLUTION	30
3.3.1 Model Selection	30
CHAPTER 4	32
4.1 DATA COLLECTION	32
4.2 DATA PRE-PROCESSING	32
4.2.1 Empty Data Analysis	32
4.2.2 Data Imputation	33
4.2.3 Correlation in Dataset	33
4.3 FEATURE SELECTION AND CONVERSION	34
4.4 FEATURE STATISTICS	36
REFERENCES	39
APPENDICES	43

List of Figures

FIGURE 2.1: PHASE DIAGRAM OF OIL AND GAS [22]-----	14
FIGURE 2.2 : DATA ANALYTICS PIPELINE-----	17
FIGURE 2.3 : POSITIVE CORRELATION [42]-----	23
FIGURE 2.4 : NEGATIVE CORRELATION [42] -----	23
FIGURE 2.5 :NO CORRELATION [42]-----	23
FIGURE 3.1 : FLOWCHART FOR PROPOSED SOLUTION -----	31
FIGURE A.1 : FEATURE CORRELATION FOR VOLVE AND KYLE MASTER DATASETS WITH FORWARD FILLING-----	43
FIGURE A.2 : FEATURE CORRELATION FOR VOLVE AND KYLE MASTER DATASETS WITH MEAN IMPUTATION -----	44
FIGURE A.3: KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT FOR ON_STREAM_HRS-----	45
FIGURE A.4 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR AVG_DOWNHOLE_PRESSURE-----	46
FIGURE A.5 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR AVG_DOWNHOLE_TEMPERATURE -----	47
FIGURE A.6 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR BORE_OIL_VOL-----	48
FIGURE A.7 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR BORE_GAS_VOL-----	49
FIGURE A.8 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR AVG_WHP_P-----	50
FIGURE A.9 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR AVG_WHT_P-----	51
FIGURE A.10 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT FOR HOURS ONLINE -----	52
FIGURE A.11 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR Av. DHP (BAR)-----	53
FIGURE A.12 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR Av. DHT (DEG C)-----	54
FIGURE A.13 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR OIL (M3)-----	55
FIGURE A.14 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR GAS (M3)-----	56
FIGURE A.15 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR Av. WHT (DEG C)-----	57
FIGURE A.16 : KERNEL DENSITY ESTIMATION PLOT, HISTOGRAM, AND BOXPLOT WITH AND WITHOUT OUTLIERS FOR Av. WHP (BAR)-----	58

List of Tables

TABLE 1.1 : SCOPE OF ACTIVITIES -----	3
TABLE 2.1 : SAMPLE DATASET-----	20
TABLE 2.2 : SAMPLE DATASET AFTER FORWARD FILLING -----	21
TABLE 4.1 : OBSERVATIONS FOR MISSING DATA-----	33
TABLE 4.2 : FEATURES WITH HIGH PEARSON CORRELATION IN VOLVE AND KYLE MASTER -----	34
TABLE 4.3 : COLUMN HEADINGS IN VOLVE AND KYLE DATASET -----	36
TABLE 4.4 : FEATURE STATISTICS FOR VOLVE DATASET-----	36
TABLE 4.5 : FEATURE STATISTICS FOR KYLE MASTER DATASET-----	37
TABLE 4.6 : FEATURE STATISTICS FOR VOLVE + KYLE MASTER DATASET -----	38

Chapter 1

INTRODUCTION

This chapter introduces the project the author worked on alongside the author's team. It also includes the background of the project as well as the aims, vision, and mission of carrying out this project. It will also describe the structure and provide insights into the remaining chapters.

1.1 Background

There is no doubt that oil and gas are vital elements to the growth of the economy. There have been traces of oil trade ever since 1875 BC [1]. In this modern, technologically-advanced society, the demand for oil and gas has only continued to grow stronger. It is used for many modern inventions enjoyed by a vast majority of people, such as vehicles, fuels, medical equipment, agriculture, and many more [2]. Additionally, the oil and gas industry has also provided jobs to thousands of individuals [3].

There are many oil and gas reserves in different corners of the world. In Indonesia, in particular, the Energy Ministry has recorded that in January 2021, there is a total reserve of 2.44 billion barrels of oil and 43.6 trillion cubic feet of gas [4]. However, due to the rapidly increasing population and a growing economy, the demand for oil and gas in Indonesia is rising [5]. Furthermore, 50% of Indonesia's energy is derived from oil [5]. This reliance on oil results in Indonesia importing nearly 350,000 barrels per day (BPD) and 50,000 barrels of fuel per day from other countries [5].

Oil and gas have many uses and have a substantial impact on the economy of a country. Therefore, oil and gas industries often make use of software applications in order to help them manage it, such as a Virtual Data Room (VDR). A VDR is an online repository that can store data securely and can be accessed by multiple users simultaneously [6]. These kinds of applications can help the oil and gas industries discover which areas could have more oil and gas. It can also help clients visualize the oil and gas data. Lynx and INTViewer are examples of software applications capable of data visualization. These applications are similar, yet they also have their differences. Lynx offers petroleum data services, geophysical and Geographical Information System (GIS) services [7]. It offers 2D and 3D seismic viewers and costs at least £250 per user per year [7]. On the other hand, INTViewer is a platform that allows users to check seismic data, geospatial integrity and also process datasets [8]. It can cost up to \$60,000 per year [9]. These types of applications can benefit oil and gas industries greatly; however, they tend to be expensive. Therefore, the goal of this project is to develop a VDR website application with similar features intended for the oil and gas industries in Indonesia at a lower cost.

In order to enhance the VDR website application, data science could be incorporated. Data science is the method of obtaining meaningful insights from a large set of data [10]. The data will be analysed and processed so that high-level data analysis can be performed [10]. The data analysis will then reveal patterns in the data, thus enabling users to draw conclusions regarding the data [10]. This is useful for the VDR as users will be able to understand the oil and gas data, thus gaining meaningful insights from it.

1.2 Scope

In this project, the author's responsibility was to create a predictive model capable of predicting oil and gas production. The author had to collect and scrape valuable data in order to make a dataset. This dataset would then be processed and cleaned to train the machine learning model. After training the model, the author will connect the model to the website created by the other members of the author's team. Table 1.1 shows the scope of activities for the author and the author's team.

Student	Role
Kotrakona Harinatha Sreeya Reddy	Collecting and Processing Data Using the data collected to develop predictive models Visualizing data through diagrams, such as charts, as well as performing data analytics
Elizabeth Chan	Design the frontend of the proposed VDR application that uses GIS Testing (e.g., unit test & integration test)
Vicky Vanessa	Designing UI/UX of the frontend of the website application Visualizing the data of oil, gas, and water Testing

Table 1.1 : Scope of Activities

1.3 Objective

This section will describe the objectives of the author, as well as the vision and mission of the author and the author's team.

1.3.1 Aim

The main aim of this VDR website application is to help the oil and gas industry discover more profitable areas of resources by visualizing oil and gas volume as well

as visualizing reserve resources. This VDR website application aims to obtain latent information from oil and gas production data which will be used to build a predictive model. This predictive model will help the oil and gas industry by showing areas that are more likely to contain more oil and gas.

1.3.2 Vision and Mission

The vision of the author's team is to increase the use of local services in the oil and gas industries to gain more profits and indirectly increase the national income. Another vision of the author and the author's team is to make this VDR website known internationally. The mission of the author and the team is to develop a high-quality yet affordable web-based application. This application will only consist of essential features and exclude unnecessary features, which would lower the cost for the buyer. In addition to this, this application will also help engineers comprehend complex data and gain better insights into how the data can be used.

1.4 Structure

This thesis consists of four chapters which will be briefly described in this section

1.4.1 Chapter 1 : Introduction

Chapter 1 introduces the author's topic, the scope, objectives, aims, vision, and mission of this project.

1.4.2 Chapter 2 : Theoretical Foundation

Chapter 2 describes the fundamental theories behind the predictive models designed by the author. It defines specific terms and provides further insights into the problem.

1.4.3 Chapter 3 : System Design

Chapter 3 will detail the problem even further and describe the works related to the author's project while also briefly describing the model the author intends to make.

1.4.4 Chapter 4 : Solution Design

Chapter 4 focuses on the design of the solution devised by the author; it includes data pre-processing as well as how the models will be manipulated.

Chapter 2

THEORETICAL FOUNDATION

This chapter will delve into the theories and techniques the author used while developing this project. It discusses the process of the Software Development Life Cycle (SDLC), which the author's team will follow while developing the project. Afterwards, it delves into the specifics of a VDR, which the author's team intends to develop. It will then probe into how oil and gas are produced in the reservoir. Additionally, it will discuss how the author intends to build the model to predict oil and gas production using machine learning. Afterwards, this chapter will discuss missing data, outliers, and feature correlation in the dataset. It will also examine how to evaluate the performance of the model and how the model can be connected to a website.

2.1 Software Development Life Cycle

SDLC is the process that is made up of steps that a particular software can follow in order to develop in a proper manner [11]. This would make it more likely for the project to be accomplished on time whilst ensuring the quality of the product is suitable for the user [12]. The activities for a specific SDLC can be labelled as [13] :

- 1) understanding the case,
- 2) deciding solution scheme,
- 3) coding based on the solution decided,
- 4) testing.

However, these activities are quite broad; therefore, they can be broken down even further to illustrate the SDLC process [13] better. The phases of SDLC are

requirements analysis, design, development, testing, and deployment and maintenance.

Requirements analysis is the first phase of SDLC. In this phase, the business requirements of the project are gathered. The project managers and stakeholders will discuss to define the requirements of the software. These requirements could include answering questions such as “who will use the software” or “how will the system be used” [13]. After the discussion, a Software Requirement Specification (SRS) document will be created, which will contain the results of the discussion [13].

The main objective of the design phase is to turn the requirements specified in the first phase into an architecture [13]. In this phase, the hardware and system requirements are specified so that the architecture of the software can be defined [13]. Additionally, this phase is where testers are required to define what needs to be tested and how it should be tested [13].

In the development phase, the results of the design phase are converted into a system that meets the user requirements. A common name for this phase is the coding phase. All the developers and engineers play an active role in this phase, and they are required to follow the required guidelines defined beforehand [14]. It is the most extensive yet most crucial phase in the entire SDLC process. Additionally, the process of the development phase will be recorded in a document entitled Source Code Document (SCD) [14].

The next phase is the testing phase, where the software developed in the previous phase will be tested. There is usually a specific team whose purpose is solely for testing the software; their job is to conduct a series of tests on the software [14]. The testing team will document any errors they encounter and send this report to the development team so that the developers can attempt to remove the errors [13]. The testing phase is one of the most essential phases as it decides whether the software is eligible to be released to the users [14].

In the deployment and maintenance phase, the software has passed the testing phase and is bug-free; therefore, it is now deployed and useable by the client [14]. Additionally, in this phase, there are possibilities that the software needs to be updated due to technological advancements. Therefore, the developers need to maintain the software to ensure that its performance will not decline [11].

Over the years, the SDLC model has been adapted into different kinds of models. These models include the *Waterfall Model*, *V-shaped Model*, *Incremental Model*, *Agile Methodology*, and many more [13]. The agile methodology, in particular, is known for constant iterations for software testing and development [11]. In this methodology, it is typical for the development phase and the testing phase to occur concurrently [13]. The Agile methodology contains twelve core principles, which are [15] :

- 1) customer satisfaction,
- 2) adaptive to changing requirements,
- 3) regular software delivery, the faster the better,
- 4) productive collaboration between developers and stakeholders,

- 5) support developers by supplying an ideal work environment and believe that they will accomplish the project,
- 6) direct face to face communication for team discussion,
- 7) assess progress by checking on working software,
- 8) encourage maintainable development,
- 9) constant focus on technical quality and design,
- 10) simplicity is vital,
- 11) working units that can organize themselves will provide the ideal output (design, software architecture, requirements),
- 12) occasional reflection so that the team can improve.

The agile methodology also consists of a framework entitled SCRUM. SCRUM is an agile development methodology that is based on an iterative as well as an incremental process [13]. One of the main features of SCRUM is that it focuses more on feedback, revisions, and frequent customer engagement rather than documenting procedures and predicting a plan of action for accomplishing the project [13]. In SCRUM, there are three prominent roles which are Product Owner (PO), Scrum Master (SM), and Scrum Team (ST). There is a lack of guidelines or descriptions for how the project should be accomplished in SCRUM; most of the decision-making is left to the team doing the project as the team knows best [13]. There are three constants in SCRUM, which are Product Backlog, Sprint Backlog, and Sprint Goal [13]. Product Backlog is the list of things that need to be done by the PO, Sprint Backlog is the list of things selected by the ST that needs to be done in the current sprint cycle, whereas Sprint Goal is the endgame of the current sprint [13]. SCRUM Methodology is beneficial for

complicated projects, and this methodology greatly helps the project progress efficiently [13].

2.2 Virtual Data Room

A VDR will be developed through a website in which the concept of SDLC will be used during the development process. A VDR is based on the concept of a data room. A data room is a valuable tool for the oil and gas industry as the companies can use it whenever they desire to dilute equity in assets [16]. The company places the data in the data room where it can be assessed [16]. If the company wishes to sell the data, buyers can visit the data room and inspect the data [16]. There are different types of data rooms, namely, Physical Data Room (PDR), VDR, and a PDR – VDR combination [16].

PDR is a physical secure room where the data is placed by the seller [16]. However, PDR has mostly been replaced by VDR as it has several drawbacks. A few of the drawbacks of a PDR are that it is expensive, burdensome, and time-consuming compared to VDR [16]. A VDR is a website where documents are uploaded that users can access and assess at ease [16]. A VDR is as secure as PDR and is also always available for as long as the client desires [16]. Furthermore, as a VDR can be accessed online, there is no need for the client or their representative to travel to access the data [16]. The information in the VDR can also be updated immediately to show new information and is immediately accessible [16]. Based on [17], the features of a VDR could include

- 1) Accessible anywhere and anytime, disregarding operating systems,
- 2) Downloading and generating documentations such as reports, and

- 3) Petrotechnical solutions such as reservoir analysis or exploration and production tools.

Figure 2.1 shows some of the features available in a VDR.

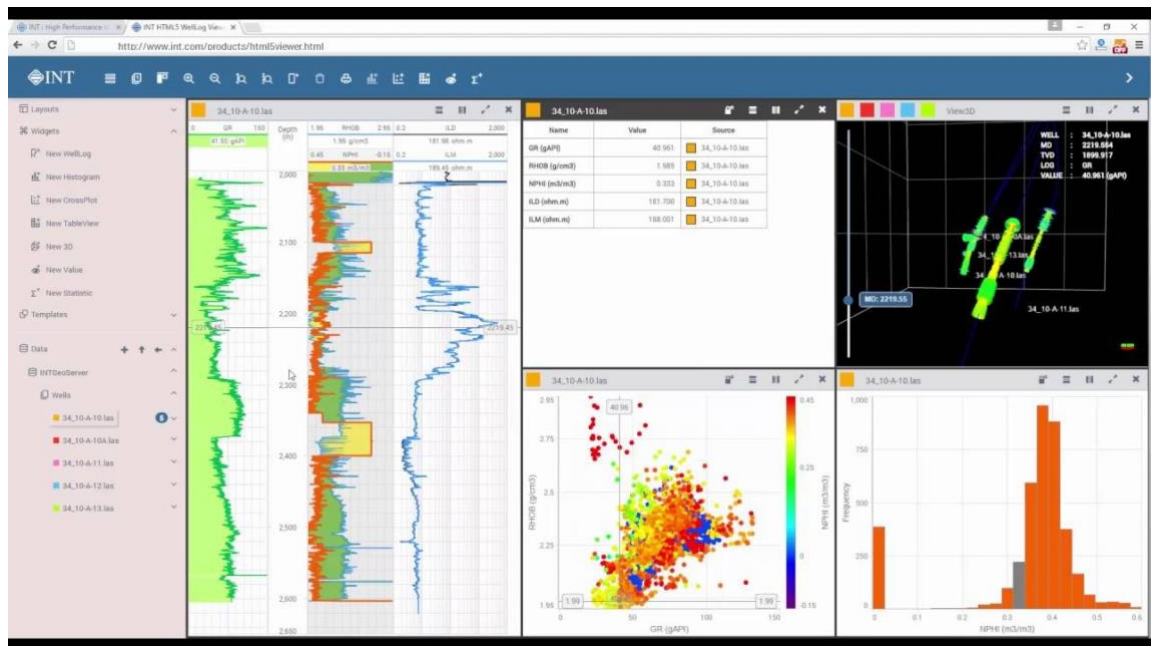


Figure 2.1: Dashboard in a Virtual Data Room

To summarize, a VDR is faster, cheaper, and more efficient compared to PDR [16].

2.3 Oil and Gas in a Reservoir

In order to build an oil and gas predictive model, it is vital to understand how oil and gas are formed in a reservoir and the factors that affect its formation.

2.3.1 Oil Formation

A formula that can be taken into account for oil formation is the oil formation volume factor (B_o). It is the ratio of oil volume and dissolved gas at a specific temperature and pressure that is needed to make one barrel of oil [18]. B_o is either greater than or equal to unity [19].

The equation for the oil formation volume factor is :

$$B_o = \frac{(V_o)pT}{(V_o)_{sv}}. \quad (2.1)$$

In Equation 2.1, B_o is the oil volume factor, V_o is the volume of oil, $(V_o)_{sc}$ is the volume of oil measured under standard conditions, p is the pressure at the reservoir, whereas T is the temperature at the reservoir [18]. From Equation 2.1, it can be inferred that temperature and pressure are essential factors in the formation of oil. Once the oil reaches the surface, it loses the dissolved gas, which leads to changes in the reservoir oil obtained. First of all, the mass of the oil will reduce as it loses the dissolved gas, then the oil will also contract as temperature decreases on the surface [18]. Afterwards, the oil will again expand as the pressure increases [18]. Often the effect of the temperature and pressure changes when the oil reaches the surface is minimal and will cancel out each other [18].

2.3.2 Gas Formation

A formula that can be taken into account for gas formation is the gas formation volume factor (B_g). It is the ratio of the volume of gas at a specific temperature and pressure that is needed to manufacture one standard volume of gas [20]. This equation for gas formation volume factor can be expressed as :

$$B_g = \frac{V_{p,T}}{V_{sc}}. \quad (2.2)$$

In Equation 2.2, B_g is the gas formation volume, $V_{p,T}$ is the volume of gas at the reservoir pressure and temperature and V_{sc} is the volume of gas at standard conditions.

In real life, gases follow the real gas law, which can be expressed mathematically as :

$$pV = znRT, \quad (2.3)$$

where p is the pressure, V is the volume, n is the number of moles of gas, R is the universal gas constant, T is the temperature, and z is the gas compressibility factor [21]. Variable z can be expressed as :

$$z = \frac{V_a}{V_i}, \quad (2.4)$$

where V_a is the actual volume of n -moles of gas at a certain temperature and pressure, and V_i is the ideal volume of n -moles of gas at the same temperature and pressure [21]. Therefore, the equation for real gas law should be applied to Equation 2.2. Equation 2.3 is applied onto Equation 2.2 by substituting for the volume (V), which will result in Equation 2.5.

$$B_g = \frac{zTP_{sc}}{T_{sc}P}. \quad (2.5)$$

In Equation 2.5, B_g is the gas formation volume, P is the pressure, T is the temperature, P_{sc} is 1 atm, T_s is 60°F, and z is the gas compressibility factor at standard conditions (1.0) [21]. With the assumption that the standard conditions are represented by $P_{sc} = 14.7 \text{ psia}$ and $T_{sc} = 520$, Equation 2.5 can be reduced to :

$$B_g = 0.0283 \frac{zT}{P}. \quad (2.6)$$

2.3.3 Pressure and Temperature in Oil and Gas Formation

Figure 2.2 shows the phase diagram of oil and gas in a reservoir. As stated previously in Section 2.3.1, when oil is drilled, it also contains dissolved gas. Therefore, in a reservoir, there exist 2 phases, namely liquid and gas. Based on the current pressure and temperature, the phase diagram shows that there is a region where the mixture will be either liquid or gas only and a region where both liquid and gas are at equilibria. The black line, known as the Bubble Point Line, denotes where both phases begin to appear [22]. Before the bubble point, the only phase that exists is liquid. However, at a constant temperature, as pressure decreases, the total volume of gas increases, whereas the volume of oil decreases [22]. This property is supported by Le Chatelier's Principle, which states that an increase in volume or decrease in pressure would increase the formation of the gaseous product.

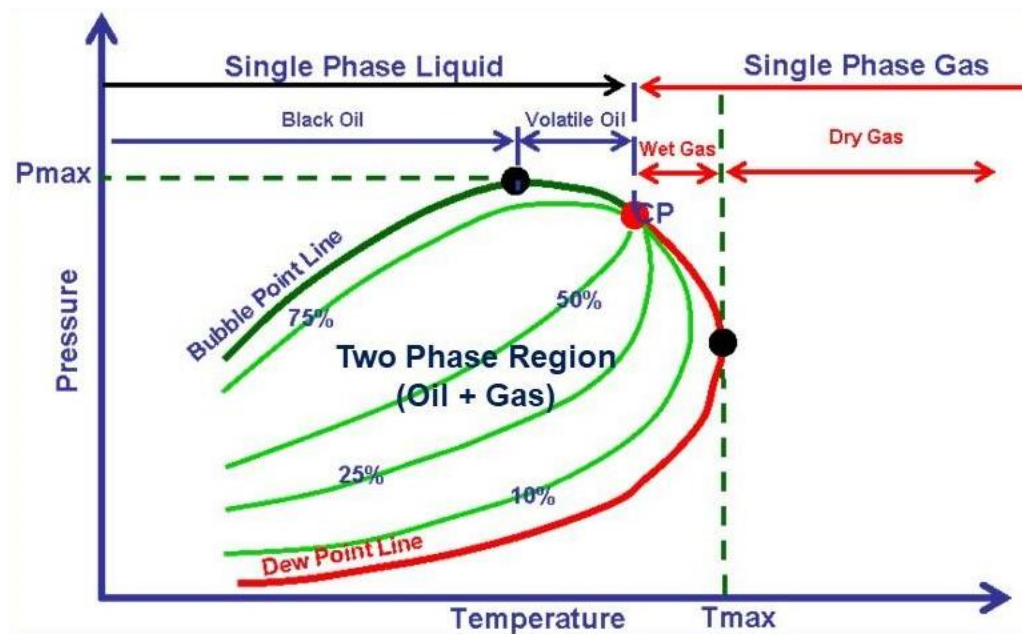


Figure 2.2: Phase Diagram of Oil and Gas [22]

As the pressure continues to decrease, more heavier molecules become gaseous, increasing the density and viscosity of the gas [22]. Subsequently, there will be a point where only a small portion of liquid remains; this is called the Dew Point [22]. If the pressure drops below this point, the only phase that exists is gas [22].

2.4 Machine Learning

Upon briefly explaining the oil and gas formation process, this section will delve into machine learning. Machine learning is the central machinery in building a prediction model of the oil and gas production data. It is defined as *the capability of a system to be able to learn from data and algorithms to automate the process of solving certain tasks* [23]. It is a branch of Artificial Intelligence (AI) that centers on using data and algorithms to echo the way humans act and learn [24]. Machine learning helps uncover insights, make classifications and predictions from the data given in order to aid users [24]. Machine learning depends on a dataset, which is a collection of data that will be regarded as one unit by the machine [25]. This dataset will act as the “training data” for the machine to learn. It is preferable to have large amounts of data as this means the machines would learn more efficiently and be able to solve problems with better accuracy. However, the quantity of the dataset is not the only significant factor in machine learning; the quality of the dataset is also a notable factor. A machine would perform significantly better with a high-quality dataset in contrast to a poor-quality dataset. In the context of the author’s problem, a high-quality dataset would mean a dataset that provides the oil production value, the gas production value, and data on sensors such as pressure and temperature.

Machine learns in different ways, namely, supervised learning, unsupervised learning, and reinforcement learning. Supervised learning is a part of machine learning and artificial intelligence; it is learning by means of mapping between a set of input variables and output variables [26]. The input variables are fed into the machine learning model, and after the training phase, it will apply what it learned to unknown data [27]. This type of machine learning is one of the most common methods and is usually used for classification and regression problems [27]. There are several types of supervised machine learning models, namely Naive Bayes, linear regression, support vector machine (SVM), KNN, and others [28]. On the other hand, unsupervised learning aims to derive meaningful information from unlabelled data [29]. It is not as commonly used as supervised learning [29]. Reinforcement learning is another complex part of AI where the model is trained to make decisions sequentially [30]. The output is dependent on the state of the current input, and the following input would then be reliant on the output of the previous output [30]. In terms of this project, the most suitable type of machine learning would be supervised learning. This is because there are datasets available that contain production values of oil and gas as well as values of sensors such as temperature and pressure.

2.5 Tree-Based Algorithms

Tree-based algorithms are a well-known part of Machine Learning, more specifically predictive modelling. Predictive modelling is the process of predicting future outcomes from data gathered beforehand. Tree-based regression algorithms are commonly used for predictive analysis of numerical values [31]. This regression model works by investigating the connection between variables [31]. It will determine the value of one variable based on the other variables present [31].

A commonly used algorithm for predictive models is the random forest algorithm [31]. This is a supervised learning algorithm that is based on the ensemble learning method [31]. Ensemble learning is the process of combining the prediction results of several machine learning algorithms [31]. The goal of this is to make the prediction results more accurate. The random forest algorithm combines the predictive results of several decision trees [31]. The respective decision trees do not interfere with one another [31]. There are two steps for the random forest algorithm; the first step is building n decision tree regressors, where n is the number of decision tree regressors [31]. These trees can be modified by specified hyperparameters, such as the strategy best used to split the node into sub-nodes or the function used to measure the quality of the split [32]. The final step would be to take the average prediction values of the decision tree regressors; this average will serve as the final output of the model [31].

Another algorithm for predictive models is the gradient boosting algorithm. This algorithm is based on the concept of boosting [33]. In terms of regression, boosting is a procedure of building strong regressors by combining weak learners [33]. This algorithm has three requirements, namely loss function, weak learners, and additive model.

A loss function would measure how similar the values predicted by the algorithm are to the actual values. In terms of regression problems, the loss function used could be Mean Absolute Error (MAE), Mean Squared Error (MSE), Root Mean Square Error (RMSE), and Coefficient of Determinant (R^2). [33]. Additionally, this algorithm is based on the idea that combining multiple weak learners would result in an accurate

result. The weak learners used in gradient boosting are typically decision trees [33]. Gradient boosting is also an additive model as it adds the weak learners one by one. Every new predictor would gain new knowledge from the error of the previous predictor, and it would work to correct the error, which would result in a better model [33].

2.6 Data Analytics Pipeline

The predictive models have to be trained on a dataset so that they can learn; however, before training, it is vital to understand and clean the dataset used. The steps that can be taken to understand and clean the dataset are shown in Figure 2.3. These steps will be explained further in the upcoming sections.

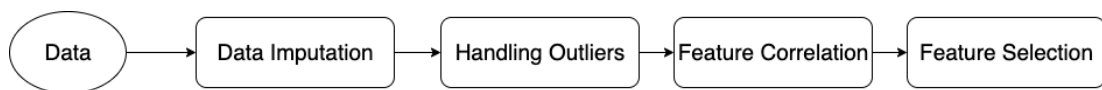


Figure 2.3 : Data Analytics Pipeline

2.7 Data Imputation

An essential part of model training is the quality of the dataset. A possible problem in a dataset is missing data. Missing data in a dataset could prove to be problematic as it could affect the model's ability to perform well.

2.7.1 Mechanism of Missingness

There are three possible mechanisms for missing data in a dataset; these mechanisms are Missing Completely at Random (MCAR), Missing at Random (MAR), Missing Not at Random (MNAR).

In the MCAR mechanism, the missing values are unrelated to the other values in the dataset, both missing and present; therefore, these missing values are random. In this situation, the missing values are considered negligible as they would not significantly impact the model performance [34].

In the MAR mechanism, the missing values are also random such as in MCAR; however, there are possibilities of the data in question being dependent on other values in the dataset. In this situation, the missing values should be considered as they could affect the model's performance. However, the effect is not extreme [34].

In the MNAR mechanism, the missing values are strongly dependent on the other values in the dataset, both missing and present. MNAR is the most serious reason mechanism for missing data as it cannot be ignored and could affect the model's performance [34]. In these cases, it is recommended to validate the data collection process [34].

In order to counteract the effects of the missing values on the dataset, data imputation methods could be used. Data imputation methods include central value imputation and forward filling.

2.7.2 Central Value Imputation

Central value imputation is the process of filling in the missing data in the dataset with their central tendencies [35]. These central tendencies could either be the mean, median, or mode. The mode is typically used to fill in the missing data for categorical

variables, whilst the mean and median are often used to fill in for numerical variables [36]. The central tendencies are deemed as reasonable estimates for filling in the missing data. However, this method would not yield ideal results if the missing data follows the MNAR mechanism, and it could also introduce bias in the dataset [37]. Additionally, filling in the missing values with the mean could reduce the variance in the data set [35].

2.7.3 Forward Filling

Forward filling is the process of filling in the missing data with the value observed before the missing value [38]. For instance, in a dataset such as Table 2.1, the forward filling method could be used to fill in the missing data. Using this method would change the dataset, as shown in Table 2.2. This method is generally used for time series datasets and is one of the easiest ways to deal with missing values. However, a disadvantage of this method is that it will not be able to fill in the missing value if there is no value prior to the missing value.

5	NaN	4
NaN	3	2
3	2	NaN

Table 2.1 : Sample Dataset

5	NaN	4
5	3	2
3	2	2

Table 2.2 : Sample Dataset after Forward Filling

2.8 Outliers

Besides missing data, another problem possible in a dataset is the presence of outliers. Outliers can be defined as *a data in a dataset that strays from the other data* [35]. It is necessary to detect these outliers as they could skew the model's training which would reduce the accuracy of the model [39]. The removal of outliers is usually one of the earliest steps in a machine learning problem [39]. There are several methods that can be utilized in order to identify these outliers. One of those methods is to use Tukey's method. The Tukey's Method is based on statistics where data is expected to follow a distribution model such as normal distribution [40]. A data is considered an outlier if it deviates from the model [40]. The Tukey's Method divides the dataset into quartiles; the quartiles commonly used are the lower quartile (Q_1), median (Q_2), and upper quartile (Q_3) [40]. The equation for a quartile is :

$$Q_r = l_1 + \frac{r\left(\frac{N}{4}\right) - c}{f}(l_2 - l_1), \quad (2.7)$$

where Q_r is the r^{th} quartile, l_1 is the lower limit, l_2 is the upper limit, f is the frequency, and c is the cumulative frequency of the class preceding the quartile class [40]. The Tukey's Method involves calculating the Interquartile Range (IQR) between the lower quartile and the upper quartile in a boxplot [40]. The equation for the IQR is

$$IQR = Q_3 - Q_1. \quad (2.8)$$

In order to accurately determine which data is an outlier, the Tukey's Method calculates the upper limit and lower limit of the data distribution. The equation for the upper limit is

$$Upper\ Limit = Q_3 + (1.5 * IQR). \quad (2.9)$$

On the other hand, the equation for the lower limit is

$$Lower\ Limit = Q_1 - (1.5 * IQR). \quad (2.10)$$

The Tukey's Method will remove any data that does not fall between the upper limit and lower limit [40].

2.9 Correlations

In order to better understand a dataset, the correlation between features in the dataset could be considered. Correlation is known as a statistical measure that describes how one feature is related to another feature [41]. It is often used during Exploratory Data Analysis (EDA) to gain a better understanding of how a feature affects other features in the dataset. There are different types of correlations, namely positive correlation, negative correlation, and no correlation [41].

A positive correlation denotes when the value of one feature increases, the value of the other feature increases as well [41]. In a graph format, a strong positive correlation would have a positive gradient, as shown in Figure 2.4.

A negative correlation indicates that when the value of one feature increases, the value of the other feature decreases [41]. A negative correlation would have a negative gradient, as shown in Figure 2.5.

No correlation indicates that the features being assessed are not related; therefore, a change in one feature would not impact the other feature [41]. In a graph format, features with no correlation would look like Figure 2.6.

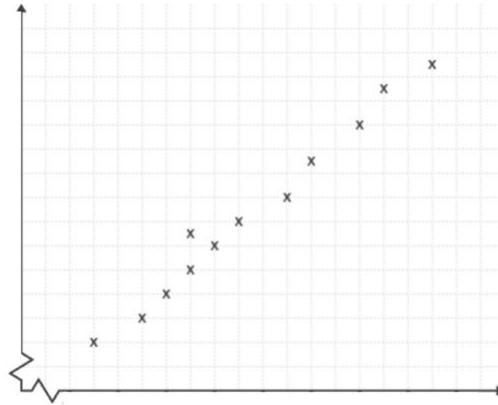


Figure 2.4 : Positive Correlation [42]

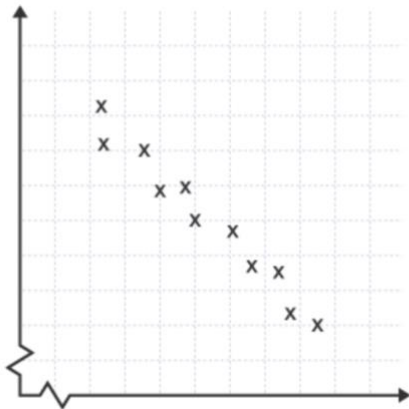


Figure 2.5 : Negative Correlation [42]

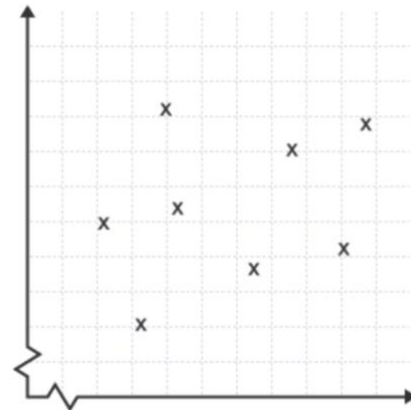


Figure 2.6 : No Correlation [42]

For numeric features, the commonly used methods for measuring the correlation between features are Pearson Correlation and Spearman Correlation.

2.9.1 Pearson Correlation

In Pearson correlation, the features being compared get assigned a value between -1 and 1 [43]. A correlation value of 1 or -1 would mean that the features being compared are strongly related to one another. A correlation value of 1 expresses that if one feature is present, then the other feature will unquestionably be present as well [43]. In addition to this, a correlation value of -1 would mean that if one feature is present, then the other feature will undeniably be absent [43]. There are also possibilities of having a correlation value of <1 or >-1. This means that the correlation is almost exactly positive or negative; however, there exists a small number of records that behave differently [43]. On the other hand, a correlation value of 0 would mean that the absence or presence of a feature is in no way related to the presence or absence of another feature [43].

The equation for Pearson correlation is :

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}, \quad (2.11)$$

where r is the Pearson correlation coefficient, x is the values in the first set of data, y is the values in the second set of data, and n is the total number of values.

2.9.2 Spearman Correlation

The Spearman correlation is a method that measures the strength and direction of the relationship between two features in a dataset [44]. Spearman correlation requires continuous data, which has a monotonic relationship. This means that when one feature increases, the other feature could either increase or decrease [44]. However, the relationship between the features does not have to be linear [44]. The correlation values

in Spearman correlation follow the same principle as those in Pearson correlation. The values range from -1 to 1 as well. If the correlation value is -1 , then as one variable increases, the other variable would decrease [44]. If the correlation value is 0 , then a change in a variable would not affect the other variable [44]. On the other hand, if the correlation value is 1 , then as one variable increases, the other variable would increase as well [44].

There are two equations that can be used to calculate Spearman's correlation. The first equation is :

$$\rho = 1 - \frac{6 \sum_i d_i^2}{n(n^2 - 1)}, \quad (2.12)$$

where ρ is the spearman correlation, d_i is the difference between the features, and n is the total number of values [45]. Equation 2.12 can only be used if there are no duplicates in the dataset. If duplicates exist in the dataset, then the second equation will be used. The second equation is :

$$\rho = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 (y_i - \bar{y})^2}}, \quad (2.13)$$

where ρ is the spearman correlation, x is the value of feature x , \bar{x} is the mean of feature x , y is the value of feature y , and \bar{y} is the mean of feature y [45].

2.10 Feature Selection

After understating the dataset, the process of feature selection could be implemented. Feature selection is the process of cutting down the input variables which will be fed into the models [46]. This is useful as it gets rid of the noise in the dataset so that the model can focus on valuable information [46]. In order to determine which features

are ideal to be used in the dataset, the Pearson correlation of the features should be taken into consideration as the values in the dataset are numerical [47]. It is ideal to add highly correlated features for the model's training. However, highly correlated parameters should not be the only features added to the model as they could reduce the model's accuracy [47]. It would lead to a lack of variation in the data or even result in data leakage, which would make the model perform unrealistically well [47].

2.11 Model Evaluation

After the dataset has been cleaned and the model has been trained, it is time to evaluate the performance of the model. Model evaluation is vital as it allows researchers to determine whether or not the model made is accurate. In order to evaluate models, researchers make use of metrics; the metrics for regression models are MAE, MSE, RMSE and, R^2 . The MAE, MSE, and RMSE metrics greatly penalize outliers as their value increases significantly in the presence of outliers [48]. For these metrics, a higher value indicates poor performance. However, RMSE is generally preferred over MAE and MSE as RMSE uses the same units as the variable in the y-axis [48]. The R^2 metric is also another ideal metric to consider as it is able to explain how well the model can predict the value compared to the original value [48].

2.11.1 Root Mean Square Error

This metric is the root squared average difference between the actual value and the predicted value [48]. RMSE is the square root of the MSE metric. The lower this value, the lower the deviations between the actual and predicted values [48].

The formula for RMSE is,

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_p - y)^2}{n}}, \quad (2.14)$$

where y_p is the predicted value, y is the actual value, and n is the number of values [48].

2.11.2 Coefficient of Determinant

This metric is the measure of how well the regression model has predicted the value based on the actual value [48]. R^2 generally ranges from 0 to 1; however, there are instances when the value could be negative [48]. A R^2 value closer to 1 would mean that the model gives an accurate prediction. The formula for R^2 is

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_p - y)^2}{\sum_{i=1}^n (\bar{y} - y)^2}, \quad (2.15)$$

where y_p is the predicted value, y is the actual value, and \bar{y} is the average of the actual values [48].

2.12 REST API

Sometimes the models developed might be used by external applications such as websites. In these cases, REST API could be used to connect the models to the external application. Representational State Transfer (REST) is a type of architectural style that specifies principles that will act as a guide for website architecture design [49]. The REST API allows users to access web services in a simple manner. Users use HTTP methods, namely GET, POST, DELETE, PUT, and PATCH, to operate the resources such as websites [49]. The GET method is mainly used to read information; this method does not allow information modification [49]. The POST method is used to

create new resources which are subordinate to another parent resource [49]. The DELETE method is used to delete an existing resource [49]. The PUT method is used to update a resource that is present; if the resource specified is not present, then a new resource could be generated [49]. The PATCH method is also used to update resources, similar to the PUT method [49]. However, the PATCH method only performs partial updates; it will not wholly change the resource [49]. Unlike the PUT method, the PATCH method is not capable of creating a new resource [49].

REST API is the ideal method to connect to an external application as, based on the adoption trend, REST API is widely accepted by many users [50]. Furthermore, REST separates the client side and server side, which is advantageous as if one component fails, it would not impact the other components [51]. In addition to this, REST is capable of adapting to any type of syntax or platform [51].

Chapter 3

SYSTEM DESIGN

This chapter will discuss the problem statement of this project as well as the proposed solution for the problem. It will also discuss existing works done in this field.

3.1 Problem Statement

As has been mentioned in Chapter 1, the main goal of this project is to create a VDR website application to help the oil and gas industry. The author's role in this project is to create a predictive model which will aim to predict oil and gas production. This model will only be available for premium users.

3.2 Related Works

In a study [52], the researcher made use of 2 models to predict the concentration of gas. The researcher used a long and short term memory (LSTM) model and a random forest model and compared the results. The models were evaluated with the R-squared score, root mean square (RMSE), and mean absolute error (MAE). The LSTM has an R-squared value of 0.31, an RMSE value of 0.45, and a MAE value of 0.56. On the other hand, the random forest model has a R-squared value of 0.95, RMSE value of 0.23, and MAE value of 0.34. The researcher concluded that the random forest model was simpler and gave better results than the LSTM model.

In an article [53], the researcher made use of Facebook's Prophet model in order to predict gas production. The dataset used by the researcher was Canadian's natural gas production; the dataset contained two columns which were the date and the volume of gas. The model was evaluated with the R-squared score and the mean absolute error (MAE). The R-squared value was 0.911, whereas the MAE score was 7782.

Another article [54] shows a researcher using Linear Regression to predict oil production. The dataset used was the Volve dataset which is located in the North Sea and was updated on a daily basis from 2005 to 2016. The linear regression model was evaluated with the R-squared score; the value was 0.55.

The same article [54] also showed the researcher using Polynomial Regression to predict oil production. The same dataset used for the linear regression model was used for this polynomial regression. The dataset features were converted into their higher orders, and the linear regression algorithm was applied to it. The model was evaluated with the R-squared score; the value was 0.95.

3.3 Proposed Solution

Figure 3.1 shows the flowchart of the proposed solution. The user will first upload the file containing the data, such as the pressure and temperature of the oil and gas wells. The model will then predict oil and gas production from this file, and the user will have the option to download the prediction results.

3.3.1 Model Selection

As discussed in Section 3.2, there are several models that have been used in the field of oil and gas production prediction. Amongst these models, the random forest algorithm was shown to achieve one of the best results. In this project, the author will use the gradient boosting algorithm as it has the capability of giving a more accurate result compared to the random forest algorithm. This is because in the gradient boosting algorithm, the trees are trained one by one; thus, the current tree is capable of correcting the error of the previous one [55].

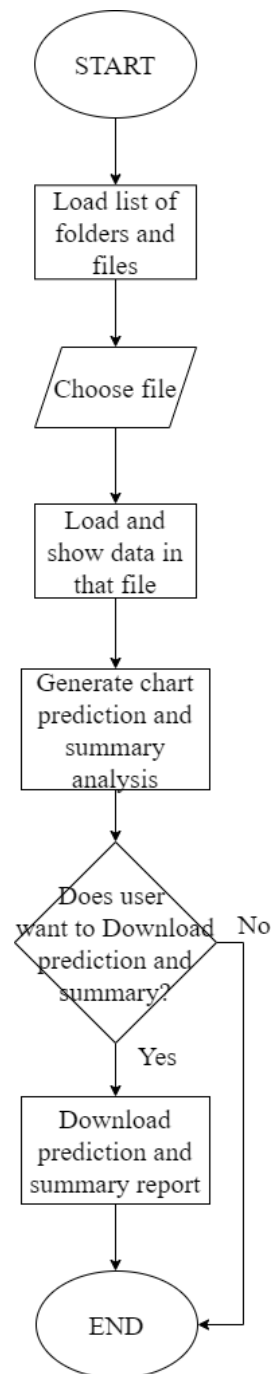


Figure 3.1 : Flowchart for Proposed Solution

Chapter 4

SOLUTION DESIGN

This chapter will explain the process of data collection, data cleaning, and pre-processing. The collected data that has been cleaned and processed will be used to train the model.

4.1 Data Collection

For the project, the author utilized two open-sourced datasets. The first dataset is entitled Volve, whilst the second dataset is entitled Kyle Master. The Volve dataset contained 15,634 rows of data and was obtained from *Kaggle*. On the other hand, the Kyle Master dataset contained 27,324 rows of data and was obtained from the online data centre of the *Oil and Gas Authority*. It is ideal to use a large dataset as it would lead to lower estimation variance, which means the model will be able to predict more accurately. Both Volve and Kyle Master datasets contain valuable information. However, in order to ensure that the data in these datasets are in better shape for a machine learning model, data cleaning and pre-processing must be done.

4.2 Data Pre-processing

This section will highlight the steps taken in order to clean and process the data for model training.

4.2.1 Empty Data Analysis

Volve and Kyle Master contained missing data; therefore, it is imperative to check the relationship between the features in the dataset. This is done so that it can be determined whether or not the presence of the missing value is correlated to other values in the dataset. In order to check this, a heatmap was used to see the correlation

values on both datasets. Table I describes the observations derived from the heatmaps. As stated in Table 1, the Volve dataset follows the MNAR mechanism, whereas the Kyle Master dataset follows the MAR mechanism. Section 2.4.1 states that these missing mechanisms imply that the missing values are dependent on one another. Thus, it should not be ignored and should either be deleted or filled in using data imputation methods.

Dataset	Volve	Kyle Master
Observation	Contains mainly “<1” and “1” feature correlation values, meaning the features are highly dependent on one another. A value of “<1” denotes that the correlation is almost exactly 1.	Feature correlation values are mostly 0.1, and some features have a correlation value of 1, meaning most of the features do not show much correlation, however, few features are highly correlated.
Missing Data Mechanism	Missing Not at Random (MNAR)	Missing at Random (MAR)

Table 4.1 : Observations for Missing Data

4.2.2 Data Imputation

As has been mentioned in section 4.2.1, both Volve and Kyle Master dataset contains missing values. Additionally, the missing data mechanisms are not MCAR as the dataset contains missing values that are dependent on one another. Therefore action should be taken to ensure the model performance will not be affected. For this project, the author will use two methods and compare the feature correlation to see which method would make the model perform better. The first method the author will use is forward filling, where the empty value is replaced by the last observed record. The second method used is central value imputation, where the author will fill in the missing values with the mean value of the feature.

4.2.3 Correlation in Dataset

In this section, this paper will explore the correlations between the features in the dataset. Table II describes the features with the highest correlation values in Volve and Kyle Master datasets, respectively, when the respective data imputation methods are

used. The feature correlation for Volve using mean imputation is not as strong as the feature correlation when forward filling is used. On the other hand, the feature correlation for Kyle Master using forward filling is similar to when mean imputation is used.

Volve				Kyle Master			
Forward Filling		Mean Imputation		Forward Filling		Mean Imputation	
Features	Correlation	Features	Correlation	Features	Correlation	Features	Correlation
BORE_OIL_VOL and BORE_GAS_VOL	0.999	BORE_OIL_VOL and BORE_GAS_VOL	0.999	Oil (m3) and Gas(m3)	0.428	Oil (m3) and Gas(m3)	0.427
AVG_DOWNHOLE_PRESSURE and AVG_DOWNHOLE_TEMPERATURE	-0.844	AVG_DOWNHOLE_PRESSURE and AVG_DP_TUBING	0.697	Av. DHT (Deg C) and Av. DHP (bar)	-0.596	Av. DHT (Deg C) and Av. DHP (bar)	-0.596
AVG_WHT_P and AVG_WHP_P	0.677	AVG_WHT_P and BORE_WAT_VAL	0.674	Av. WHT (Deg C) and Oil (m3)	0.500	Av. WHT (Deg C) and Oil (m3)	0.500

Table 4.2 : Features with High Pearson Correlation in Volve and Kyle Master

4.3 Feature Selection and Conversion

This section will explain and justify which features will be used for the model's training and discuss feature conversion. Table 4.3 shows the columns of each dataset that have the same meaning displayed side by side. For instance, *DATEPRD* in Volve is the same as *Date* in Kyle Master. As the goal is to create a model that can predict oil and gas production, it is essential to include their production values. In Volve, the first two features selected for model training are *BORE_OIL_VOL* and *BORE_GAS_VOL*. *AVG_DOWNHOLE_PRESSURE* and *AVG_DOWNHOLE_TEMPERATURE* are also included as these features have a high correlation value. Additionally, as mentioned in Section 2.1, oil and gas formation are also reliant on pressure and temperature, which makes these features ideal for the model's training. *AVG_WHP_P* is also added to the dataset as it shows a decent correlation to *BORE_OIL_VOL* and *BORE_GAS_VOL*. Oil and gas production can

also be improved by water injection [56], therefor; therefore is also added to the model's training. *ON_STREAM_HRS* will also be added as this column shows how long the machine operates. In Kyle Master, the first two features selected are *Oil (m3)* and *Gas (m3)*, as these features contain the production value of oil and gas. *Av. WHT (Deg C)* and *Av. WHP (bar)* are also included as they have a decent correlation with *Oil (m3)* and *Gas (m3)*. Furthermore, as oil and gas production is reliant on the pressure and temperature of the reservoir, thus the features *Av. DHT (Deg C)* and *Av. DHP (bar)* are added for the model's training. Lastly, *Hours Online* will also be added for training the model. These datasets both have similar columns even though the names are different. For instance, *Av. DHT (Deg C)* and *Av. DHP (bar)* in the Kyle Master dataset has the same meaning as *AVG_DOWNHOLE_PRESSURE* and *AVG_DOWNHOLE_TEMPERATURE* in the Volve dataset. Additionally, *Oil (m3)* and *Gas (m3)* in the Kyle Master dataset have the same meaning as *BORE_OIL_VOL* and *BORE_GAS_VOL*. However, the unit of measurement in each dataset is different. Therefore, it needs to be standardized so that the model will perform better. Hence, the temperatures will be standardized into °C (degrees Celsius), while the pressures will be standardized into *bar*, and the volumes will be standardized into m^3 (meter cubic).

Volve Dataset	Kyle Master Dataset	Unit of Measurement	Selected for Model Training
DATEPRD	Date	-	No
WELL_BORE_CODE	Wellbore ID	-	No
NPD_WELL_BORE_CODE		-	No
NPD_WELL_BORE_NAME		-	No
NPD_FIELD_CODE		-	No
NPD_FIELD_NAME		-	No
NPD_FACILITY_CODE		-	No
NPD_FACILITY_NAME		-	No
ON_STREAM_HRS	Hours Online	hours	Yes
AVG_DOWNHOLE_PRESSURE	Av. DHP (bar)	bar	Yes

AVG_DOWNHOLE_TEMPERATURE	Av. DHT (Deg C)	°C	Yes
AVG_DP_TUBING		-	No
AVG_ANNULUS_PRESSES		Bar	No
AVG_CHOKE_SIZE_P	Platform Choke %	%	No
AVG_CHOKE_UOM		%	No
AVG_WHP_P	Av. WHP (bar)	bar	Yes
AVG_WHT_P	Av. WHT (Deg C)	°C	Yes
DP_CHOKE_SIZE		-	No
BORE_OIL_VOL	Oil (m3)	m^3	Yes
BORE_GAS_VOL	Gas (m3)	m^3	Yes
BORE_WAT_VOL	Produced Water (m3)	m^3	No
BORE_WI_VOL		-	No
FLOW_KIND		-	No
WELL_TYPE		-	No

Table 4.3 : Column Headings in Volve and Kyle Dataset

4.4 Feature Statistics

In order to better understand the selected features in the dataset, several techniques were employed to understand how the data is distributed. Table 4.4 describes the selected features of the Volve dataset, whereas Table 4.5 describes the selected features for the Kyle Master dataset.

Table 4.4 : Feature Statistics for Volve Dataset

Feature	Range	Outlier Count	Mean	Standard Deviation	Variance
ON_STREAM_HRS	25 hours	715	23 hours	3 hours	9 hours
AVG_DOWNHOLE_PRESSURE	307 bar	144	242 bar	27 bar	729 hours
AVG_DOWNHOLE_TEMPERATURE	107.7 °C	156	104 °C	4 °C	16 °C
BORE_OIL_VOL	5900 m^3	283	1458 m^3	1463 m^3	2.140.369 m^3
BORE_GAS_VOL	86863 m^3	182	212937 m^3	207073 m^3	42.879.227.329 m^3
AVG_WHP_P	120 bar	44	48 bar	20 bar	400 bar
AVG_WHT_P	86 °C	352	73 °C	18 °C	324 °C

Table 4.5 : Feature Statistics for Kyle Master Dataset

Feature	Range	Outlier Count	Mean	Standard Deviation	Variance
Hours Online	1912 hours	1326	23 hours	27 hours	729 hours
Av. DHP (bar)	1122 bar	3	111 bar	39 bar	1521 bar
Av. DHT (Deg C)	245 °C	645	94 °C	9 °C	81 °C
Oil (m3)	3509 m^3	447	380 m^3	328 m^3	107.584 m^3
Gas (m3)	1.304.298.362.420 m^3	226	178.525.800.000 m^3	175.599.300.000 m^3	30.835.114.160.490.000.000.00 m^3
Av. WHP (bar)	325 bar	48	57 bar	35 bar	875 bar
Av. WHT (Deg C)	228 °C	597	62 °C	19 °C	361 °C

In Table 4.4 and Table 4.5, range denotes the range of the specified feature; more specifically, it is the difference between the lowest value up to the highest value of the feature. Outlier count is the number of outliers in the feature. Mean is the center point of the feature. It is the mathematical average of the feature. Standard deviation is the measure of how varied the feature is relative to the mean.

From Table 4.4 and Table 4.5, it can be seen that the range of values for all the selected features in the Kyle dataset is larger than the features in the Volve dataset. This denotes that the data in the Kyle dataset is dispersed compared to the Volve dataset. In addition to this, the standard deviation and variance of the features in the Kyle dataset are much larger than the features in the Volve dataset. This observation further supports the fact that the features in the Kyle dataset are more spread out than the features in the Volve dataset.

For the model's training, both these datasets will be combined. As shown in Table 4.3, the selected columns in the Volve dataset and the Kyle Master dataset have the same meaning. Therefore when combining the datasets, the columns in the Volve dataset were renamed to match the columns in the Kyle Master dataset. Table 4.6 describes the selected features of the combined dataset.

Table 4.6 : Feature Statistics for Volve + Kyle Master Dataset

Feature	Range	Mean	Standard Deviation	Variance
Hours Online	1912 hours	23 hours	22 hours	484 hours
Av. DHP (bar)	162 bar	102 bar	19 bar	361 bar
Av. DHT (Deg C)	344 °C	158 °C	74 °C	5476 °C
Oil (m3)	5900 m^3	801 m^3	1117 m^3	1.247.689 m^3
Gas (m3)	1.164.213,36242 m^3	176375.2 m^3	178282.6 m^3	3.178.468.546.276 m^3
Av. WHP (bar)	325 bar	49 bar	29 bar	841 bar
Av. WHT (Deg C)	228 °C	64 °C	18 °C	324 °C

REFERENCES

- [1] M. S. Vassiliou, Historical dictionary of the petroleum industry, 2018.
- [2] International Association of Oil & Gas Producers, "Oil and gas in Everyday Life," International Association of Oil & Gas Producers, [Online]. Available: <https://www.iogp.org/oil-natgas-in-everyday-life/>. .
- [3] W. P. Council, "Why are oil and gas important?," [Online]. Available: <https://www.world-petroleum.org/edu/221-why-are-oil-and-gas-important#:~:text=Oil%20is%20one%20of%20the%20most%20important%20raw,about%20two%20million%20tonnes%20of%20oil%20and%20gas.> . [Accessed March 2022].
- [4] R. Ranggasari, "Oil and gas reserves potential in eastern Indonesia reaches 9.8bn barrels," Tempo, [Online]. Available: <https://en.tempo.co/read/1536679/oil-and-gas-reserves-potential-in-eastern-indonesia-reaches-9-8bn-barrels#:~:text=Overall%2C%20the%20Energy%20Ministry%20recorded%20there%20are%2070,2.44%20billion%20barrels%20and%20gas%20of%2043.6%20TCF.> .
- [5] Indonesia Investment, "Crude Oil Indonesia," [Online]. Available: <https://www.indonesia-investments.com/business/commodities/crude-oil/item267..>
- [6] W. Kenton, "Virtual Data Room (VDR)," 23 June 2021. [Online]. Available: <https://www.investopedia.com/terms/v/virtual-data-room-vdr.asp#:~:text=Virtual%20Data%20Rooms%2C%20or%20VDRs%2C%20exist%20as%20a,joint%20venture%20that%20requires%20access%20to%20shared%20data..> [Accessed 19 03 2022].
- [7] Lynx, "License Pricing - Lynx Information Systems," Lynx Information System, [Online]. Available: <http://www.lynxinfo.co.uk/download-pricing.html>.
- [8] Intviewer, "Intviewer - Fast Geoscience Visualization, Analysis & QC,," Intviewer, 02 August 2021. [Online]. Available: <https://www.int.com/products/intviewer/#:~:text=INTViewer%20is%20a%20platform%20and%20application%20that%20allows,to%20a%20desktop%20or%20remotely%20via%20the%20cloud..>
- [9] INTViewer, "INTViewer. Geoscience Analysis and QC, Simplified,," INTViewer, [Online]. Available: <https://www.int.com/products/intviewer/>.
- [10] IBM, "Data Science," IBM, 15 May 2020. [Online]. Available: <https://www.ibm.com/cloud/learn/data-science-introduction>.
- [11] G. Gurung, R. Shah and D. P. Jaiswal, "Software development life cycle models-A comparative study," *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, pp. 33-27, 2020.

- [12] A. Mishra and D. Dubey, "A Comparative Study of Different Software Development Life Cycle Models in Different Scenarios," *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, vol. 1, no. 5, pp. 1-6, 2013.
- [13] J. Shah, "A Comparative Study of Software Development Life Cycle Models".
- [14] P. Pedamkar, "What is SDLC: Different phases and models of SDLC," *EDUCBA*, 2022.
- [15] "Principles behind the Agile Manifesto," Agile, [Online]. Available: <https://agilemanifesto.org/principles.html>.
- [16] Harrison, Bob, "The data room," 2020, pp. 21-26.
- [17] Schlumberger, "Secure, Remote Access to Field Datasets Enables Potential Investors to Complete Asset Evaluations," Schlumberger, [Online]. Available: <https://www.slb.com/resource-library/case-study/dss/delfi-virtual-data-room-generic-asia-pacific-cs>.
- [18] A. El-Banbi, A. Ahmed and E.-M. Ahmed , "Black Oils," in *PVT Property Correlations*, Elsevier, 2018, p. 147–182.
- [19] S. Mokhatab, W. A. Poe and J. Y. Mak, "Natural Gas Fundamentals," in *Handbook of Natural Gas Transmission and Processing*, Elsevier, 2019.
- [20] A. El-Banbi, A. Ahmed and E.-M. Ahmed, "Dry Gases," in *PVT Property Correlations*, Elsevier, 2018.
- [21] T. Ahmed, "Reservoir-Fluid Properties," in *Reservoir-Fluid Properties* , Elsevier, 2010, pp. 29-135.
- [22] I. Fetoui, "Hydrocarbon Phase Behavior," [Online]. Available: <https://production-technology.org/category/pvt/>.
- [23] P. Z. a. K. H. C. Janiesch, "Machine learning and deep learning," *Electronic Markets*, vol. 31, no. 3, pp. 685-695, 2021.
- [24] IBM Cloud Education, "What is machine learning," IBM, 15 July 2020. [Online]. Available: <https://www.ibm.com/cloud/learn/machine-learning>.
- [25] I. Sydorenko, "What is a dataset in Machine Learning," High quality data annotation for Machine Learning, 5 April 2021. [Online]. Available: <https://labeledyourdata.com/articles/what-is-dataset-in-machine-learning>.
- [26] W. X. P. C. M. C. a. S. D. Y. Gong, "Supervised Learning," *Machine learning techniques for multimedia*, p. 21–49., 2008.
- [27] IBM Cloud Education, "What is Supervised Learning," IBM, 2022. [Online]. Available: <https://www.ibm.com/cloud/learn/supervised-learning>.
- [28] A. M. J. A. V. M. A. A. L. a. A. A. S. A. R. van Loon, "Understanding supervised, unsupervised, and reinforcement learning," *Big Data Made Simple*, 2019.
- [29] D. N. Dimid, "Unsupervised learning algorithms cheat sheet,," 17 February 2022. [Online]. Available: <https://towardsdatascience.com/unsupervised-learning-algorithms-cheat-sheet-d391a39de44a>.
- [30] R. S. Sutton and A. G. Barto, "Reinforcement learning: An introduction," *Cambridge, MA*, 2020.
- [31] Elzain, Hussam Eldin, Chung, Sang Yong, Senapathi, Venkatramanan, Sekar, Selvam, Lee, Seung Yeop, Roy, Priyadarsi D., Hassan, Amjed and

- Sabarathinam, Chidambaram, "Comparative study of machine learning models for evaluating groundwater vulnerability to nitrate contamination," *Ecotoxicology and Environmental Safety*, vol. 229, pp. 61-113, 2022.
- [32] Scikit, "ScikitLearn," [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeRegressor.html>.
- [33] A. Kumar, "Introduction to the Gradient Boosting Algorithm," 20 May 2020. [Online]. Available: <https://medium.com/analytics-vidhya/introduction-to-the-gradient-boosting-algorithm-c25c653f826b>.
- [34] S. Rosenthal, "Data Imputation," in *The International Encyclopedia of Communication Research Methods*, Wiley, 2017, pp. 1-12.
- [35] A. Swalin, "How to handle missing data," Medium, 19 March 2018. [Online]. Available: <https://towardsdatascience.com/how-to-handle-missing-data-8646b18db0d4..>
- [36] H. Kang, "The prevention and handling of the missing data," in *Korean Journal of Anesthesiology*, vol. 64, 2013, p. 402.
- [37] W. Badr, "'6 different ways to compensate for missing data (data imputation with examples)," 12 January 2019. [Online]. Available: <https://towardsdatascience.com/6-different-ways-to-compensate-for-missing-values-data-imputation-with-examples-6022d9ca0779..>
- [38] K. N, "Part-1 : Data Preparation Made Easy with python!!," Medium, 09 May 2020. [Online]. Available: <https://medium.com/analytics-vidhya/part-1-data-preparation-made-easy-with-python-e2c024402327..>
- [39] Á. Fernández, J. R. Dorronsoro and J. Bella, "Supervised outlier detection for classification and regression," *Neurocomputing*, vol. 486, pp. 77-92, 2022.
- [40] C. M. Salgado, C. Azevedo, H. Proença and S. M. Vieira, "Noise Versus Outliers," in *Secondary Analysis of Electronic Health Records*, Cham, Springer International Publishing, 2016, pp. 163-183.
- [41] F. Malik, "Understanding value of correlations in data science projects," Medium, 10 June 2019. [Online]. Available: <https://medium.com/fintechexplained/did-you-know-the-importance-of-finding-correlations-in-data-science-1fa3943debc2#:~:text=Correlation%20is%20a%20statistical%20measure.%20Correlation%20explains%20how,%28variables%29%20can%20be%20positively%20correlated%20>
- [42] BBC Bitesize, "Types of correlation - scattergraphs - national 4 application of Maths Revision," BBC News, [Online]. Available: <https://www.bbc.co.uk/bitesize/guides/zmt9q6f/revision/2..>
- [43] Nettleton, David, "Selection of Variables and Factor Derivation," in *Commercial Data Mining*, Elsevier, 2014, pp. 79-104.
- [44] "Spearman correlation coefficient: Definition, formula and calculation with example," QuestionPro, 15 January 2020. [Online]. Available: <https://www.questionpro.com/blog/spearmans-rank-coefficient-of-correlation/>.
- [45] Swapnilbobe, "Spearman's correlation," 13 April 2021. [Online]. Available: <https://medium.com/analytics-vidhya/spearmans-correlation-f34c094d99d8#:~:text=Here%2C%20we%20are%20calculating%20spearma>

n%E2%80%99s%20correlation%20using%20the,of%20relationship%20between%20ranks%20of%20two%20individual%20features. .

- [46] V. Kumar and S. Minz, "Feature Selection: A literature Review," *Smart Computing Review*, vol. 4, pp. 1-19, 2014.
- [47] K. Menon, "Feature selection in machine learning," Simplilearn, 16 September 2021. [Online]. Available: https://www.simplilearn.com/tutorials/machine-learning-tutorial/feature-selection-in-machine-learning#what_is_feature_selection..
- [48] A. Chugh, "MAE, MSE, RMSE, Coefficient of Determination, Adjusted R Squared — Which Metric is Better?," 8 December 2020. [Online]. Available: <https://medium.com/analytics-vidhya/mae-mse-rmse-coefficient-of-determination-adjusted-r-squared-which-metric-is-better-cd0326a5697e>.
- [49] Zhao, Yunwei and Wan, Xin, "The Design of Embedded Web System based on REST Architecture," in *IEEE*, 2019.
- [50] D. Bryant, "GraphQL-ultimate-guide," [Online]. Available: <https://www.infoq.com/articles/GraphQL-ultimate-guide/>.
- [51] BBVA API Market, "REST API: What is it, and what are its advantages in project development?," 2016. [Online]. Available: <https://www.bbvaapimarket.com/en/api-world/rest-api-what-it-and-what-are-its-advantages-project-development/>.
- [52] C. Xie, L. Chao, Y. Qin, J. Cao and Y. Li, "Using a stochastic forest prediction model to predict the hazardous gas concentration in a one-way roadway," *AIP Advances*, vol. 10, no. 11, 2020.
- [53] R. Sharma, "Using Facebook Prophet for Forecasting Natural Gas Production," Medium, 13 March 2021. [Online]. Available: <https://medium.com/mllearning-ai/forecast-using-prophet-canadian-natural-gas-production-dataset-b1f9c57548d8>.
- [54] J. Chahar, "Prediction of Oil Production by applying Machine Learning on Volve Field Production Data.," 17 December 2020. [Online]. Available: <https://www.linkedin.com/pulse/prediction-oil-production-applying-machine-learning-volve-chahar/>.
- [55] EDUCBA, "Difference Between Random forest vs Gradient boosting," [Online]. Available: <https://www.educba.com/random-forest-vs-gradient-boosting/>.
- [56] D. P. Nolan, "Overview of Oil, Gas, and Petrochemical Facilities," in *Handbook of Fire and Explosion Protection Engineering Principles for Oil, Gas, Chemical, and Related Facilities*, Elsevier, 2019, pp. 33-50.

APPENDICES

Appendix A

I. Feature correlation in Volve and Kyle Master datasets when forward filling is used

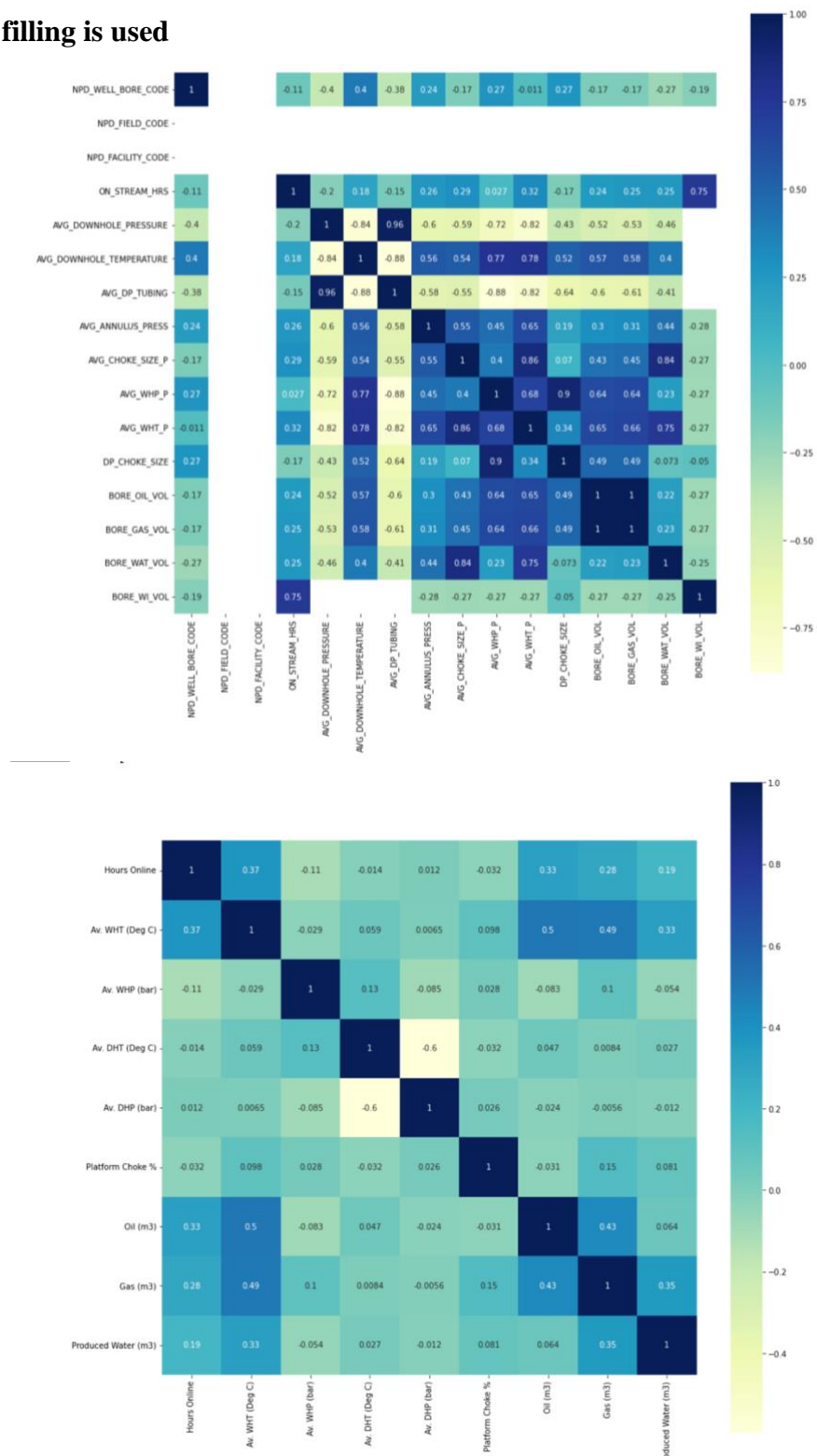


Figure A.1 : Feature correlation for Volve and Kyle Master datasets with forward filling

II. Feature correlation in Volve and Kyle Master datasets when mean imputation is used

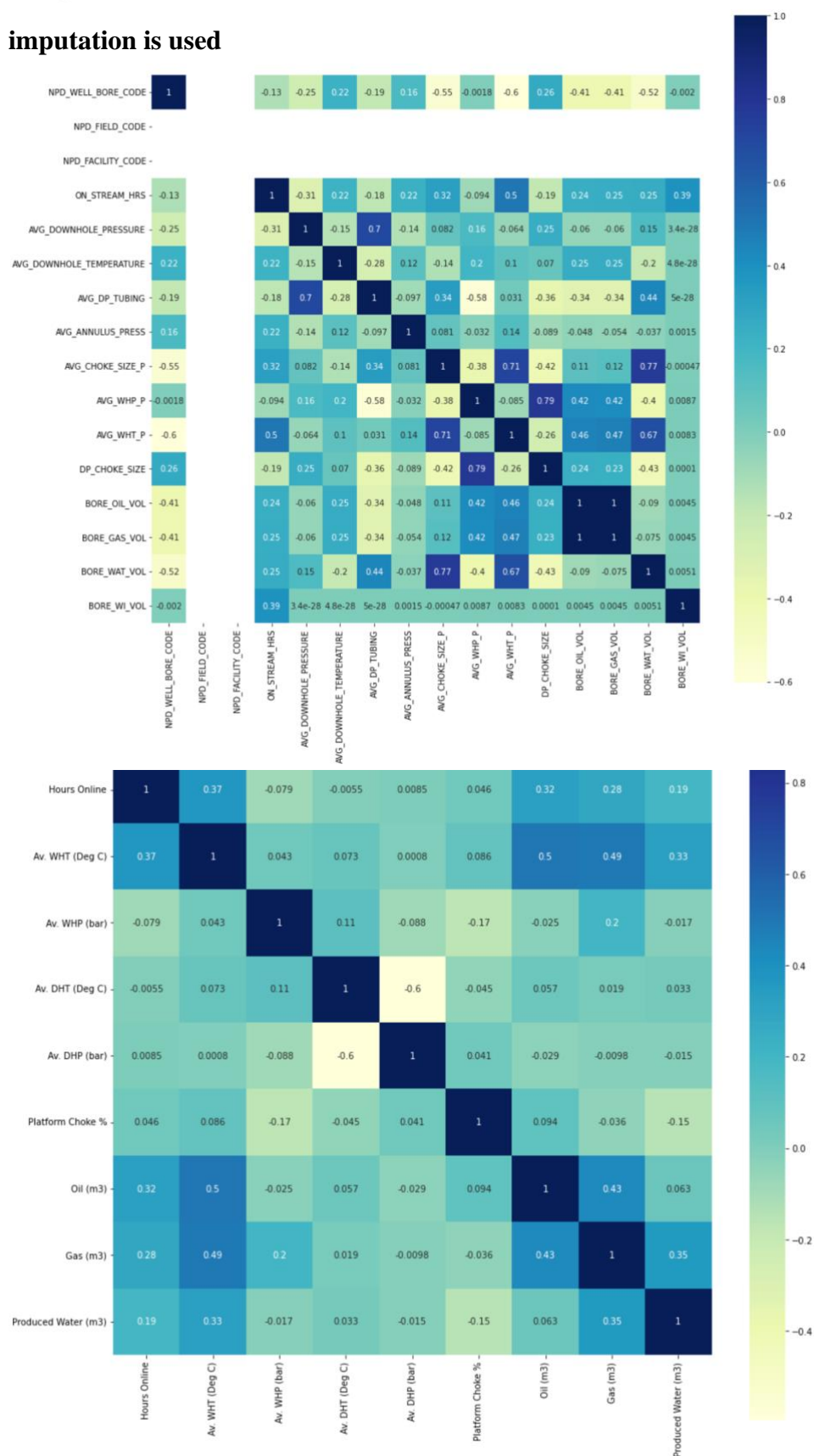


Figure A.2 : Feature correlation for Volve and Kyle Master datasets with mean imputation

III. Feature statistics in Volve dataset

ON_STREAM_HRS

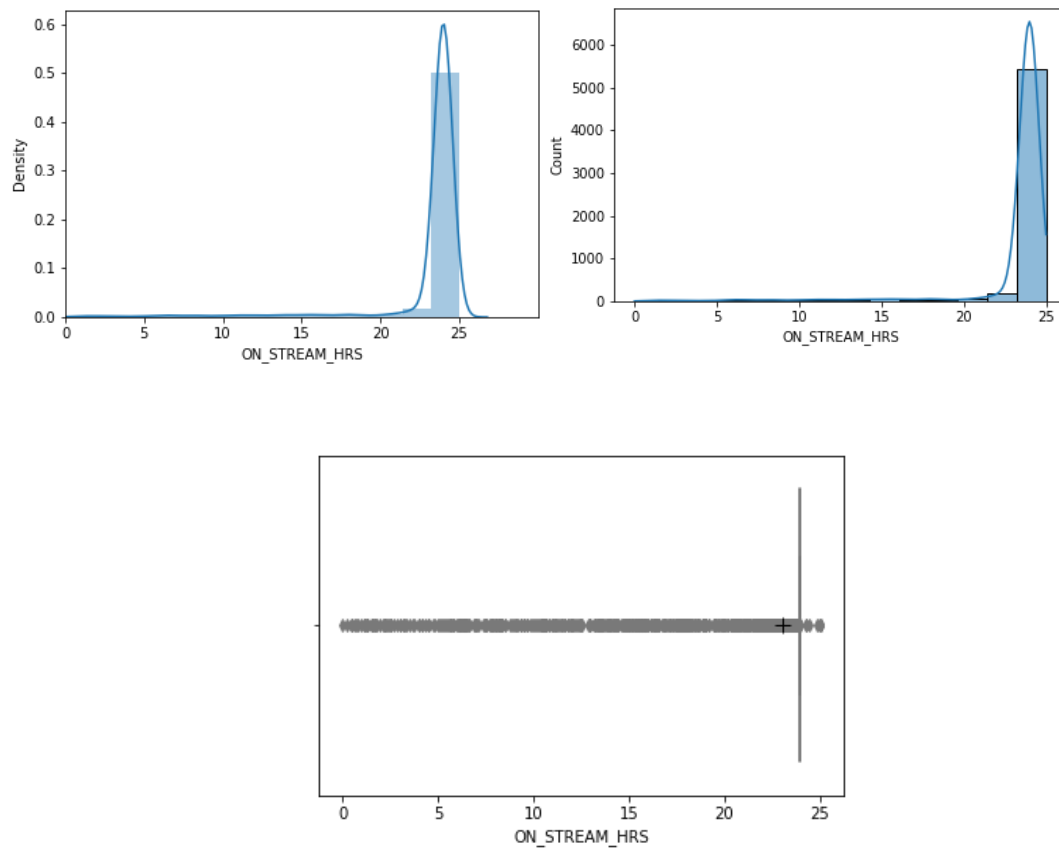
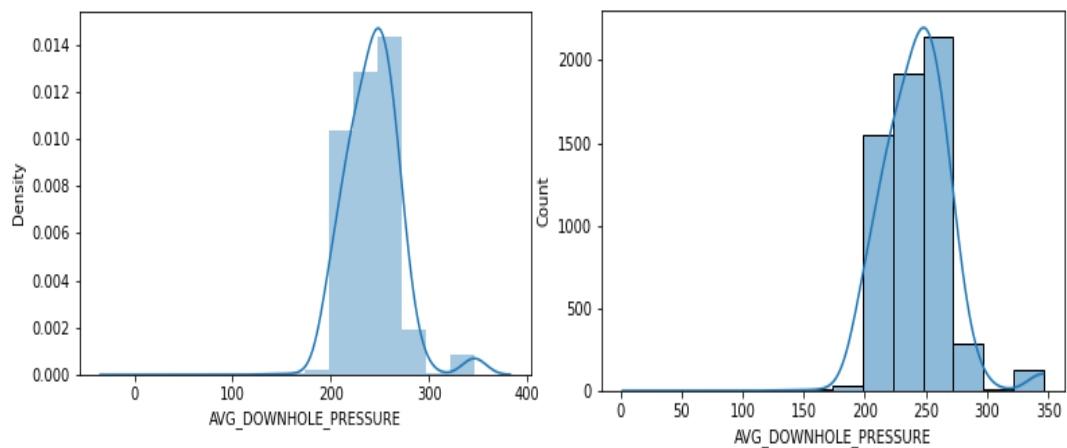


Figure A.3: Kernel Density Estimation plot, histogram, and boxplot for `ON_STREAM_HRS`

AVG_DOWNHOLE_PRESSURE



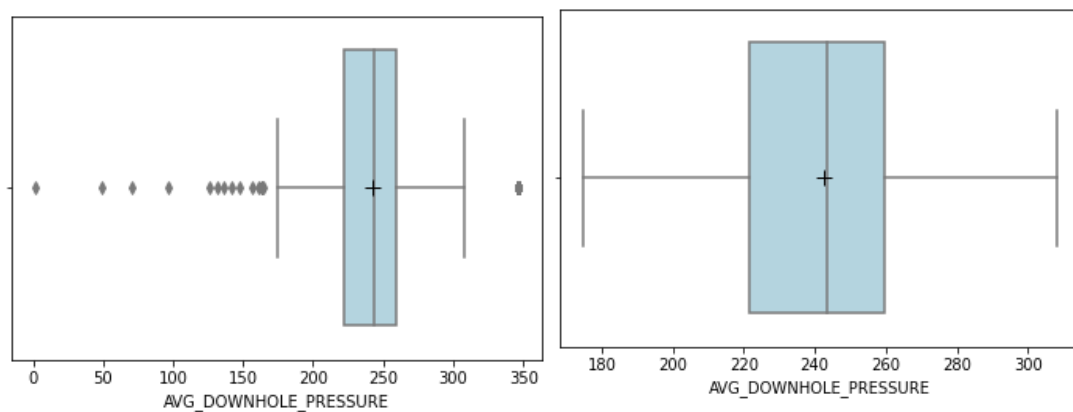
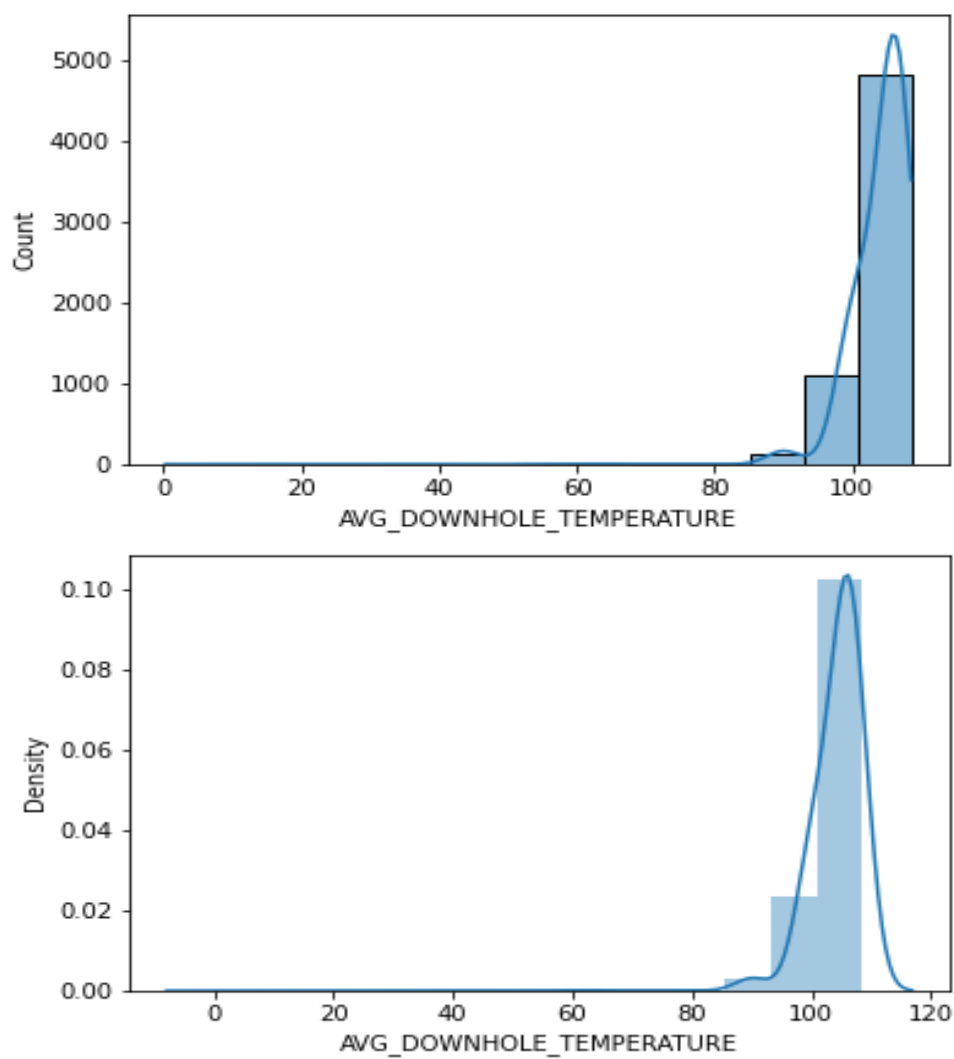


Figure A.4 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for `AVG_DOWNHOLE_PRESSURE`

`AVG_DOWNHOLE_TEMPERATURE`



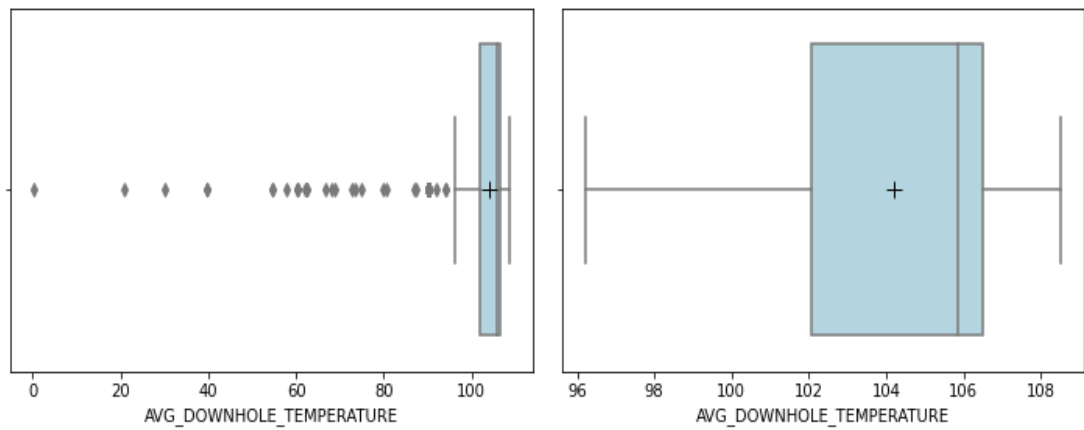
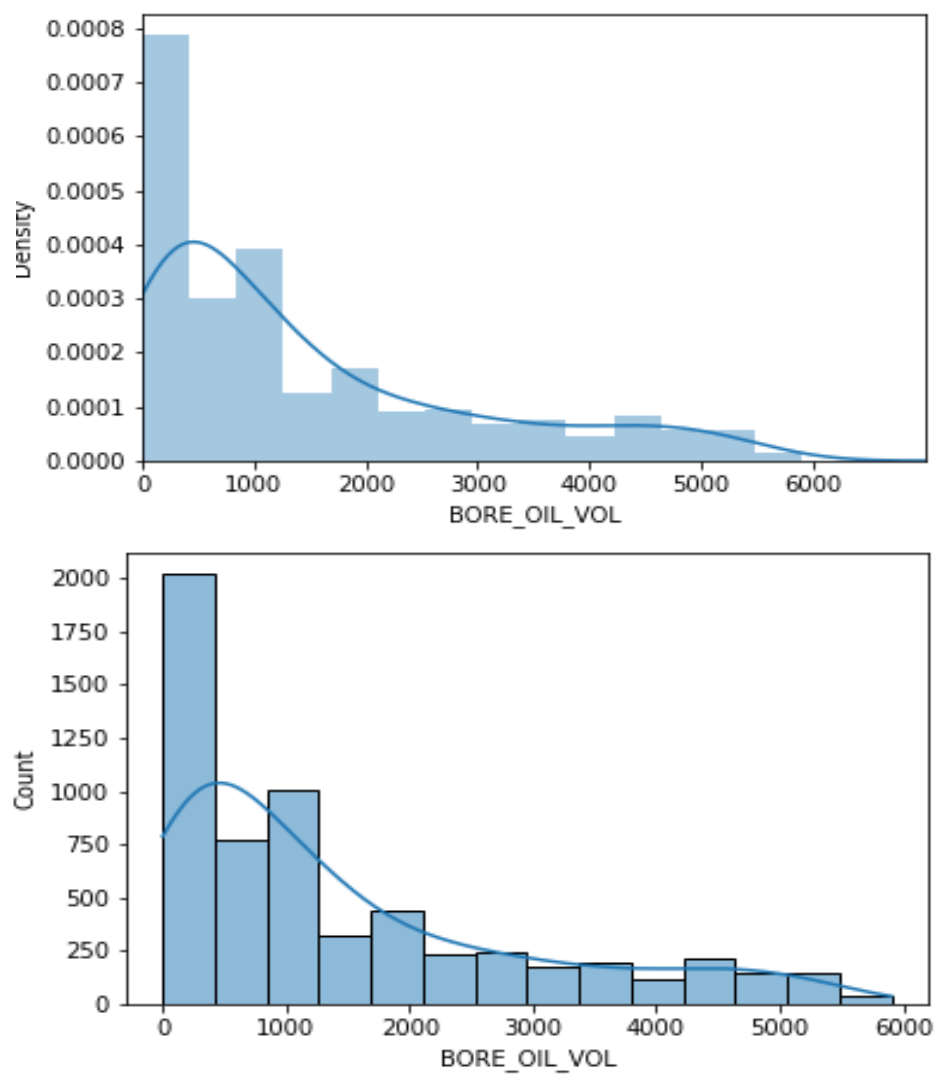


Figure A.5 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for `AVG_DOWNHOLE_TEMPERATURE`

BORE_OIL_VOL



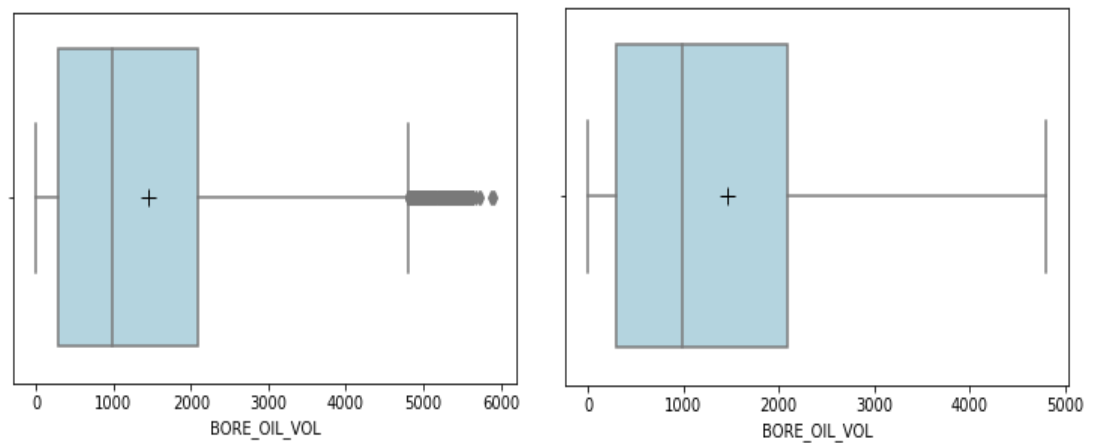
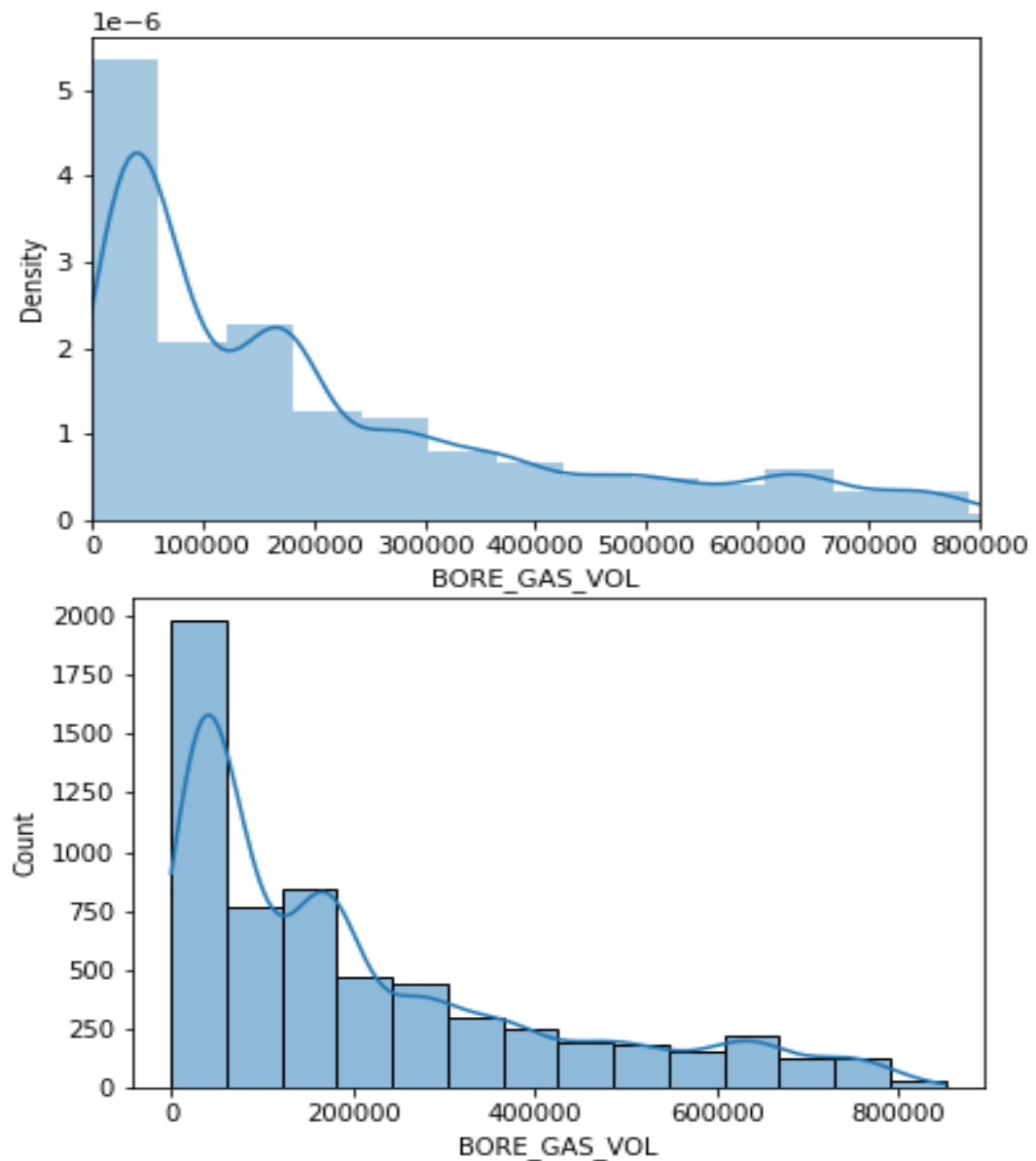


Figure A.6 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for `BORE_OIL_VOL`

BORE_GAS_VOL



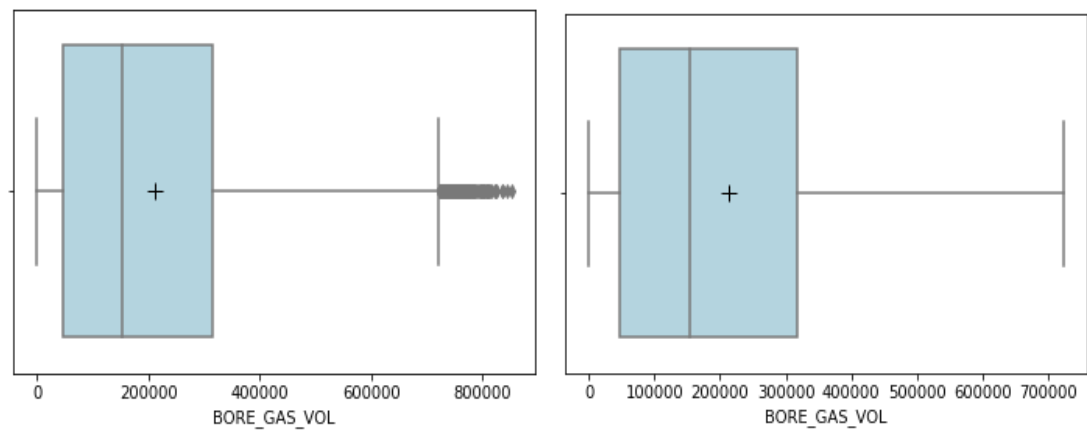
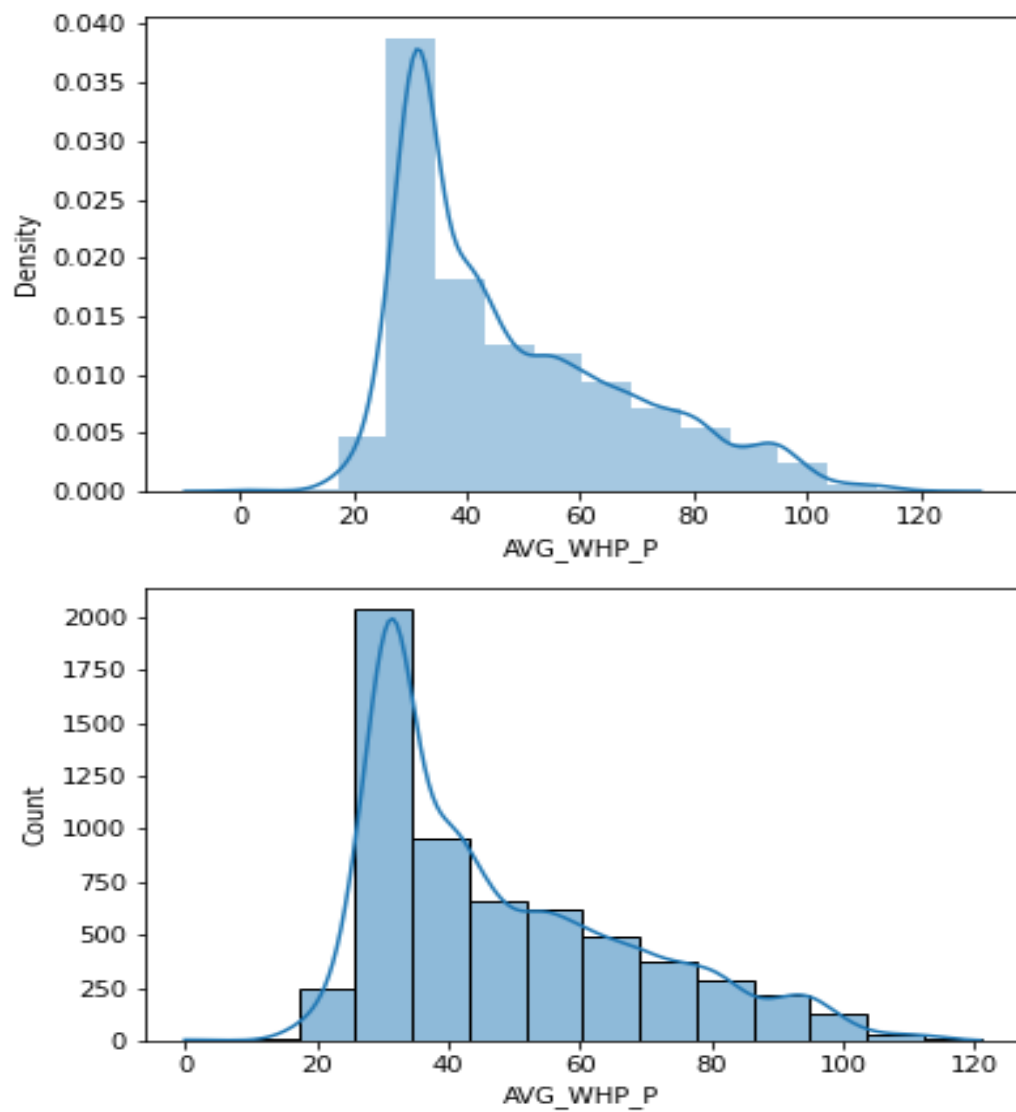


Figure A.7 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for `BORE_GAS_VOL`

AVG_WHP_P



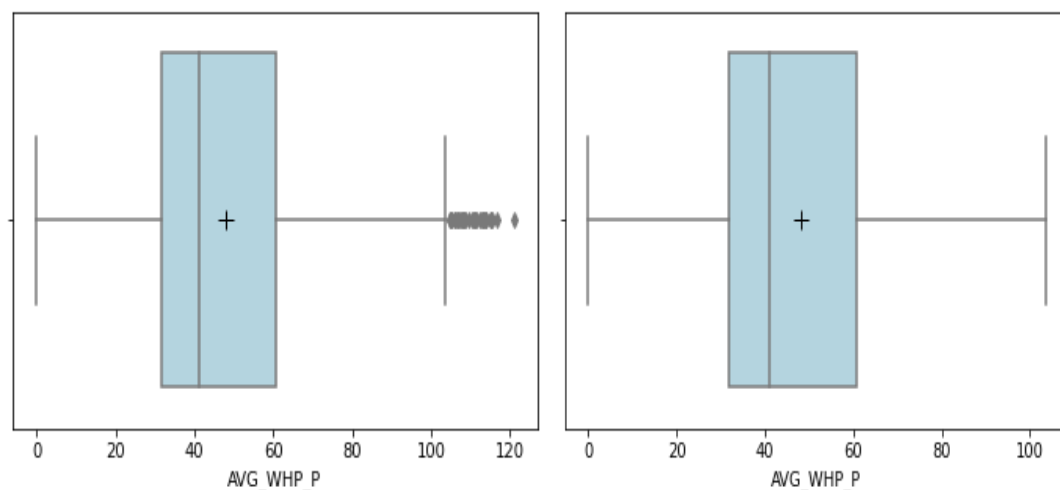
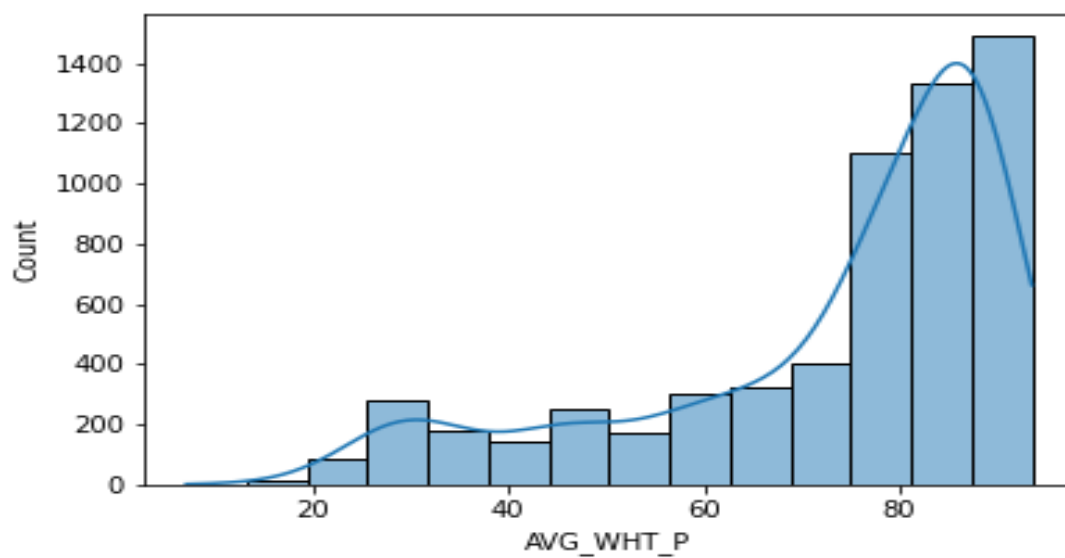
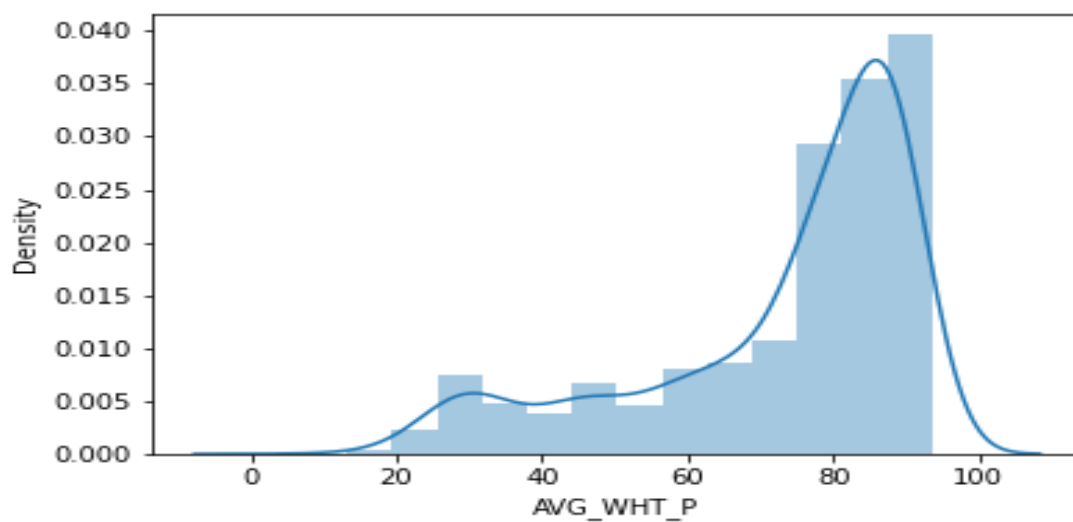


Figure A.8 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for AVG_WHP_P

AVG_WHT_P



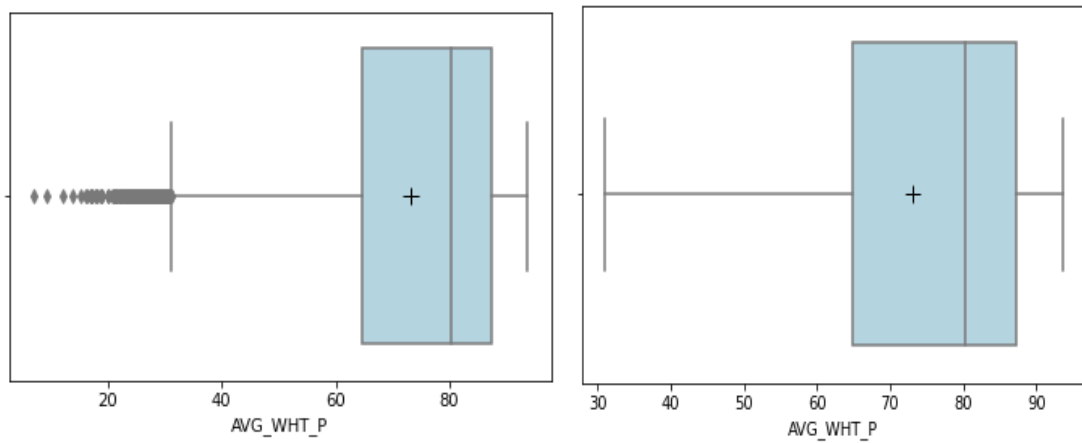
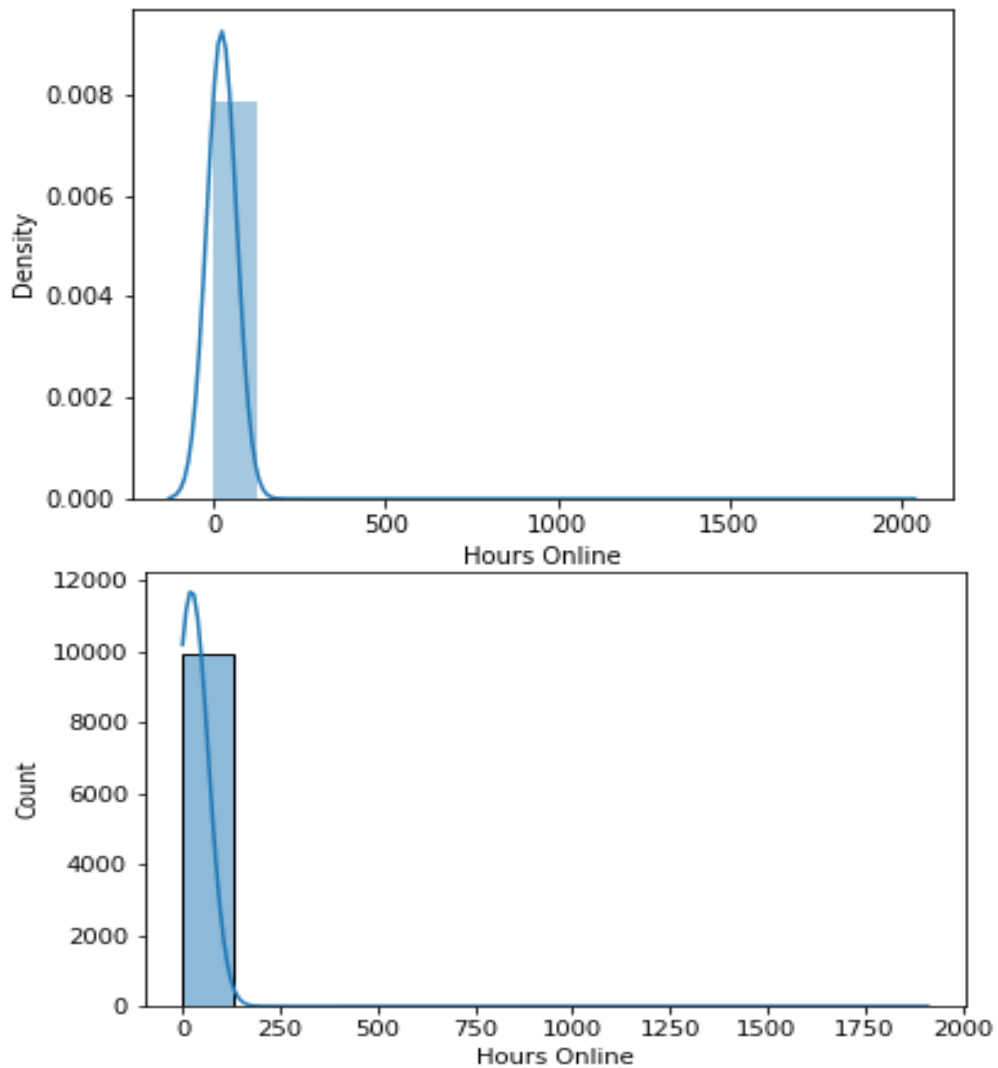


Figure A.9 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for AVG_WHT_P

IV. Feature statistics in Kyle Master dataset

Hours Online



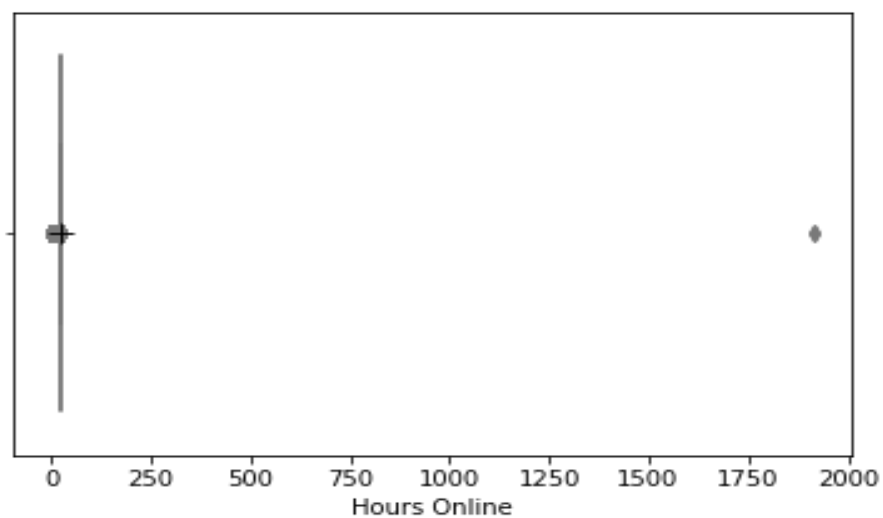
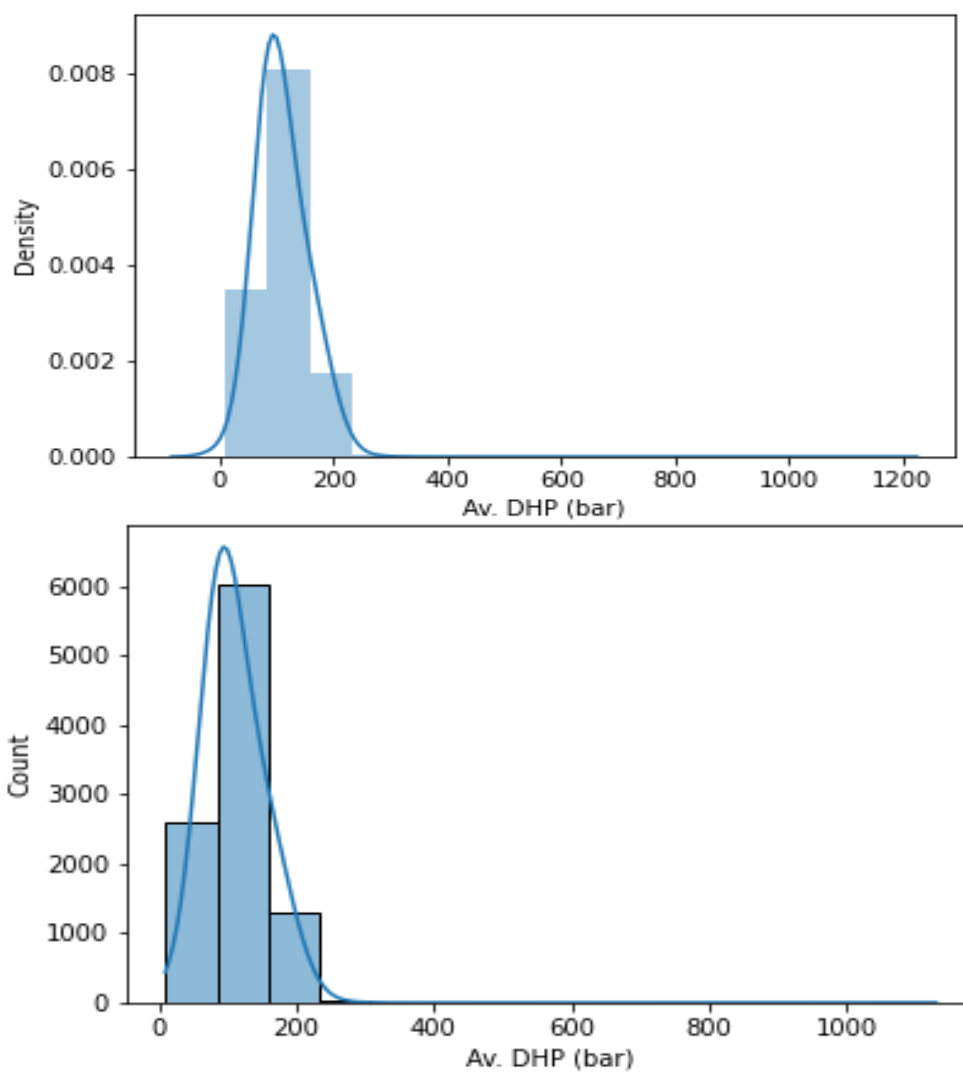


Figure A.10 : Kernel Density Estimation plot, histogram, and boxplot for Hours Online

Av. DHP (bar)



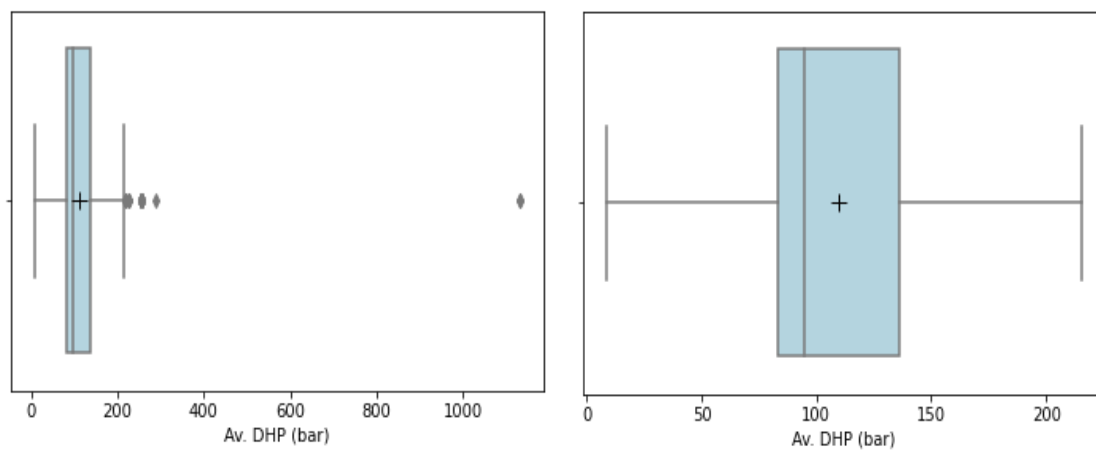
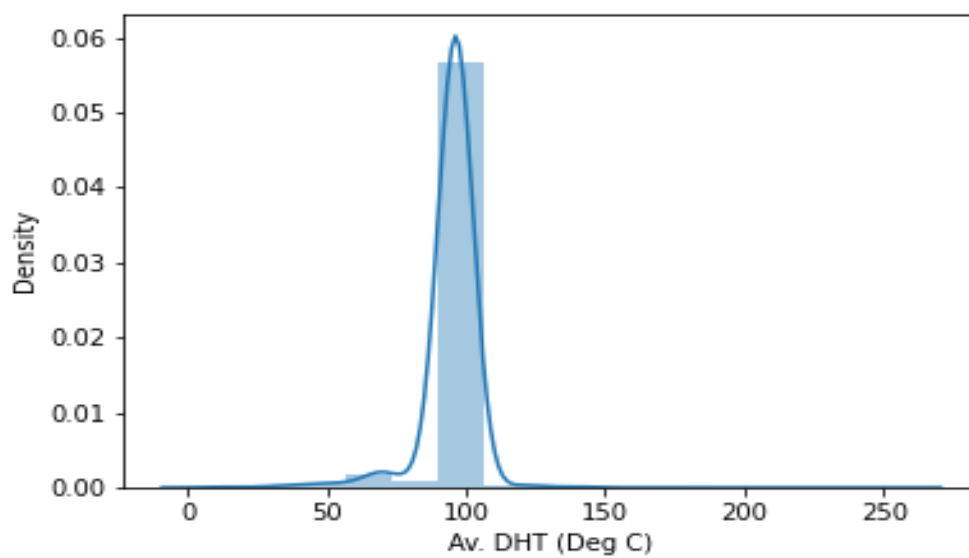
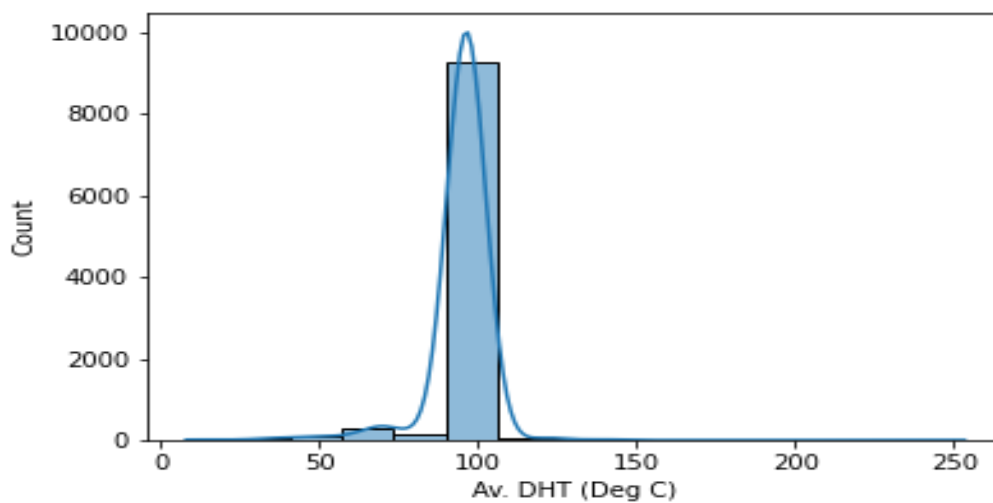


Figure A.11 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for Av. DHP (bar)

Av. DHT (Deg C)



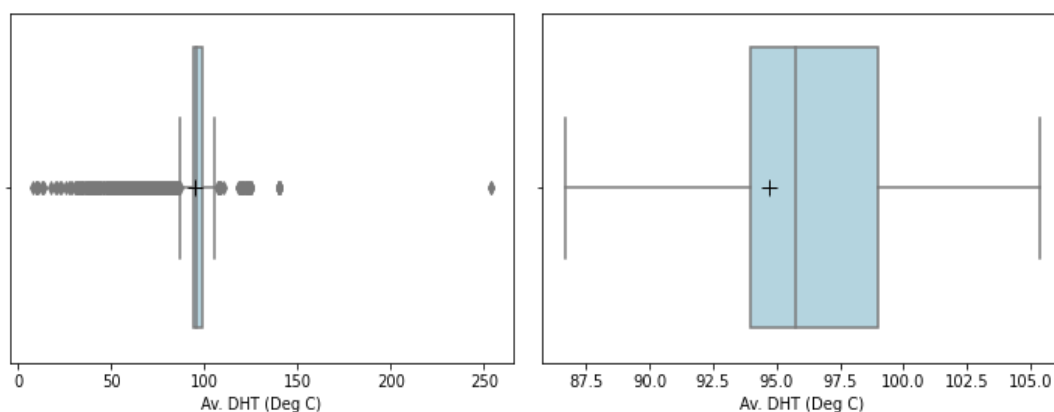
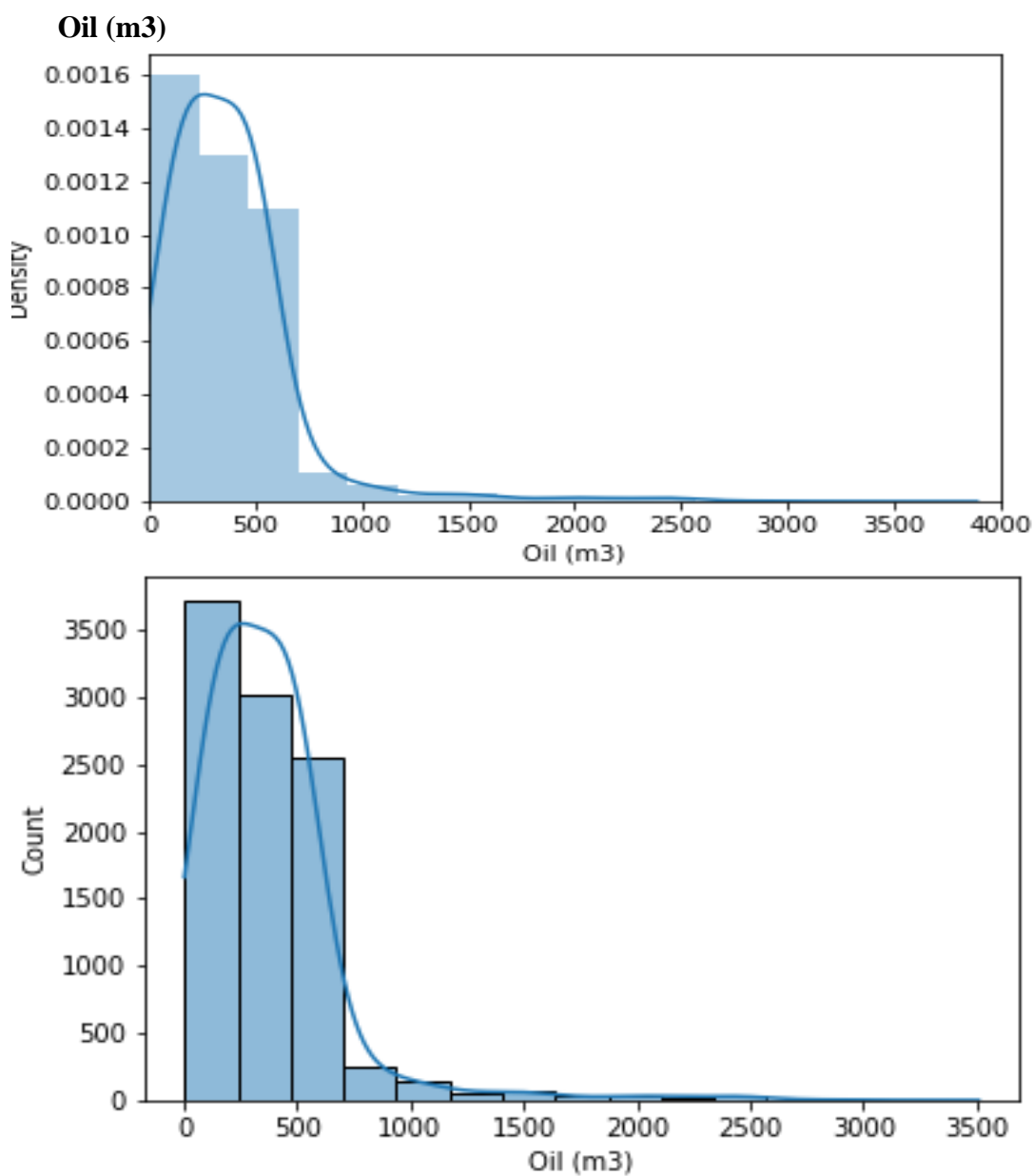


Figure A.12 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for Av. DHT (Deg C)



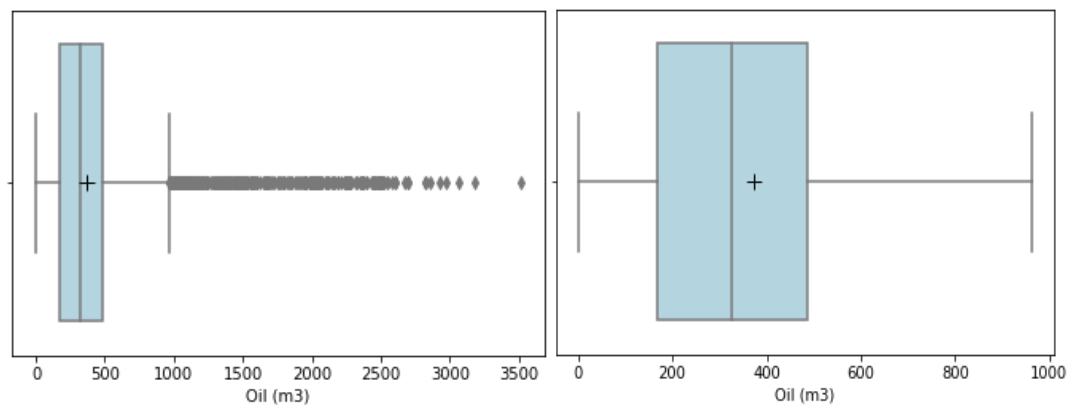
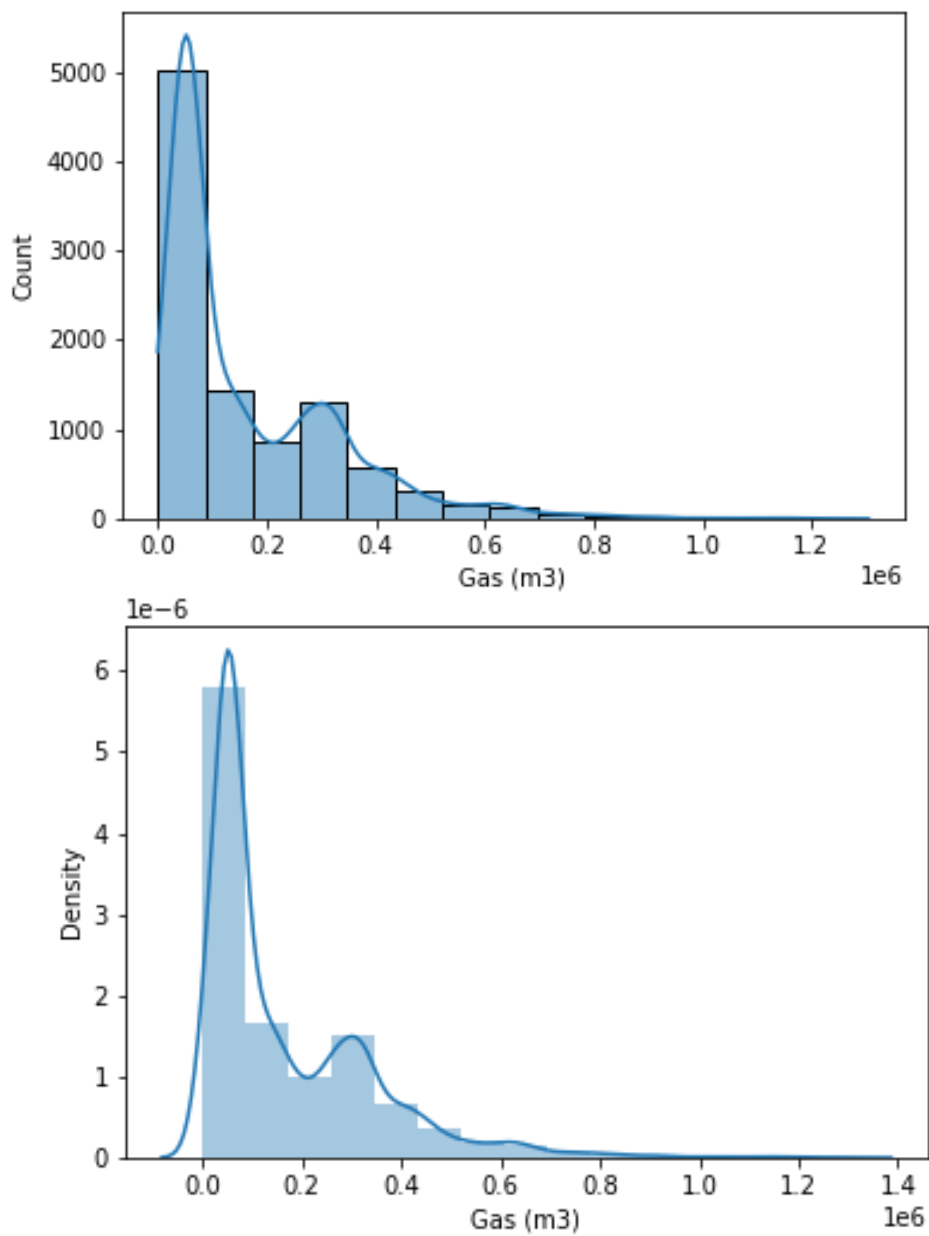


Figure A.13 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for Oil (m3)

Gas (m3)



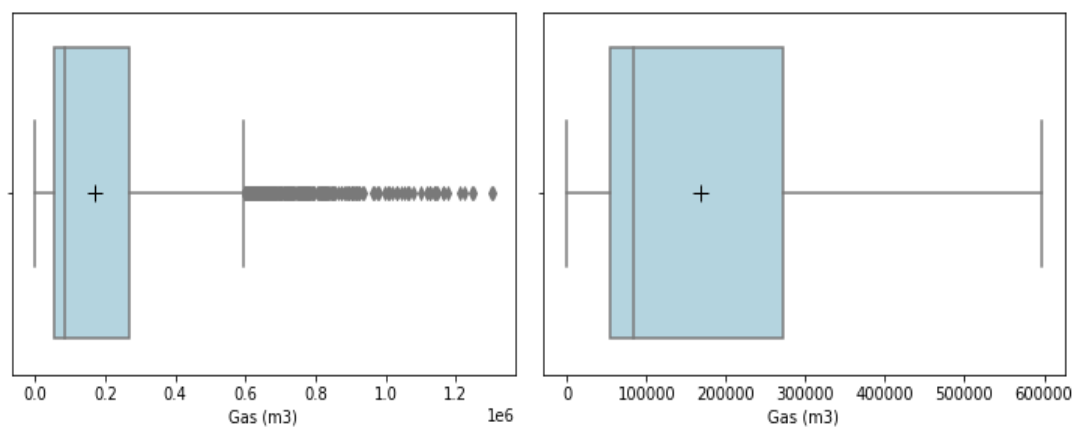
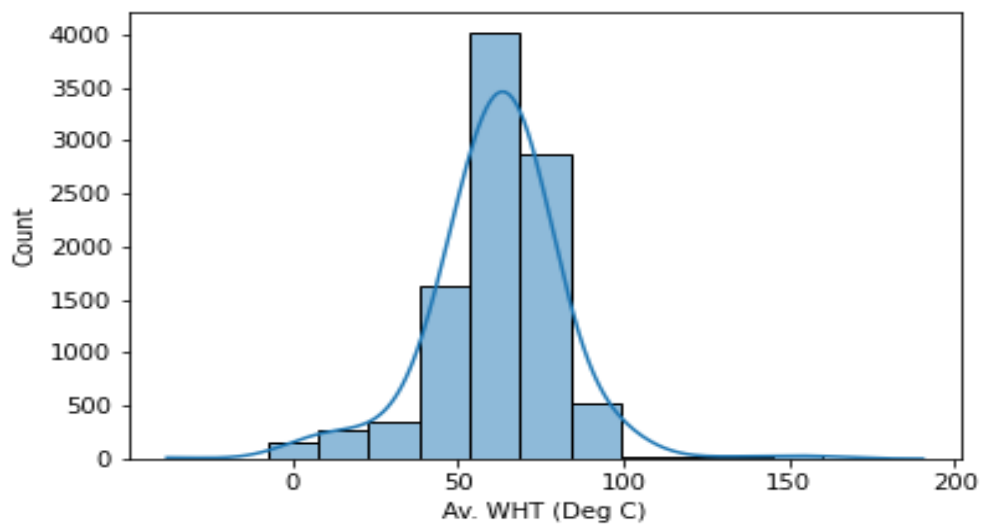
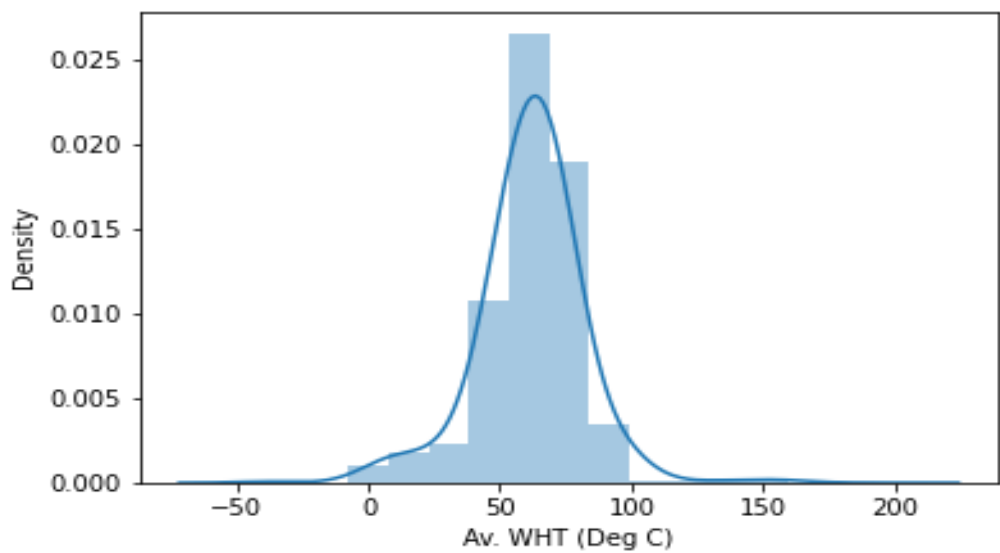


Figure A.14 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for Gas (m3)

Av. WHT (Deg C)



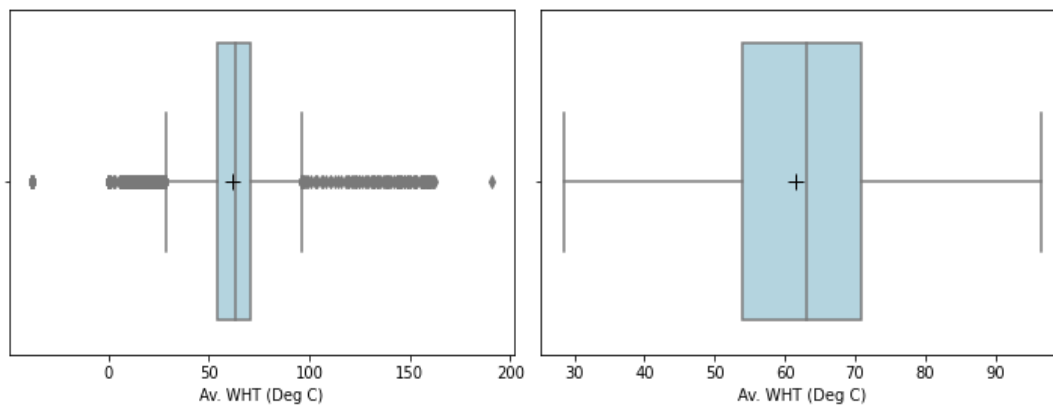
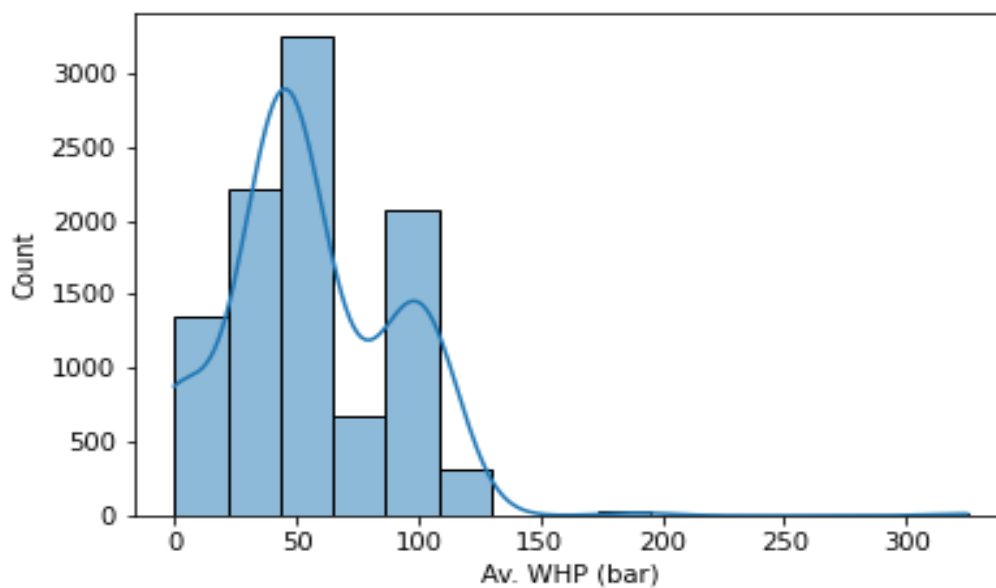
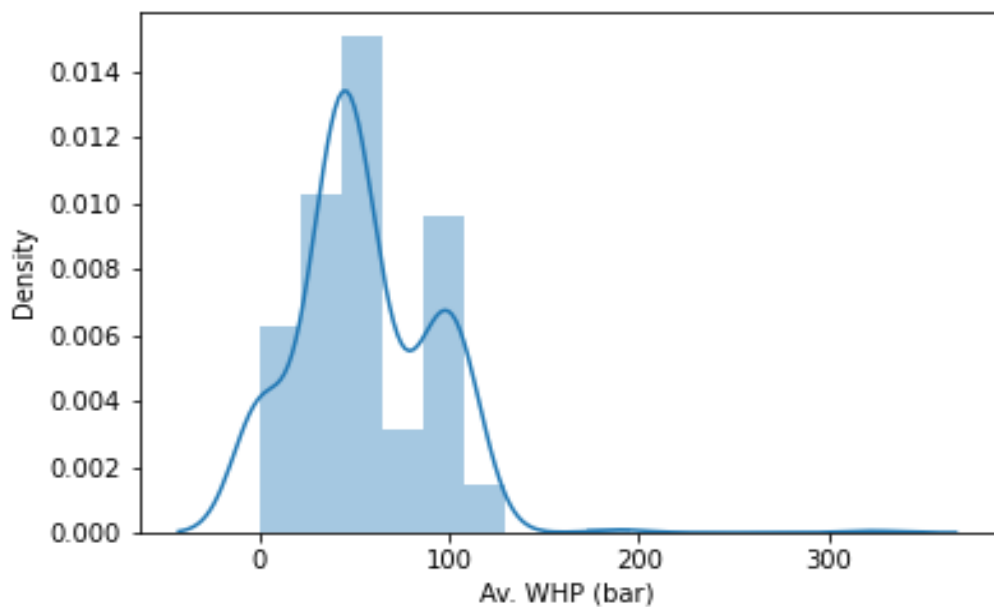


Figure A.15 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for Av. WHT (Deg C)

Av. WHP (bar)



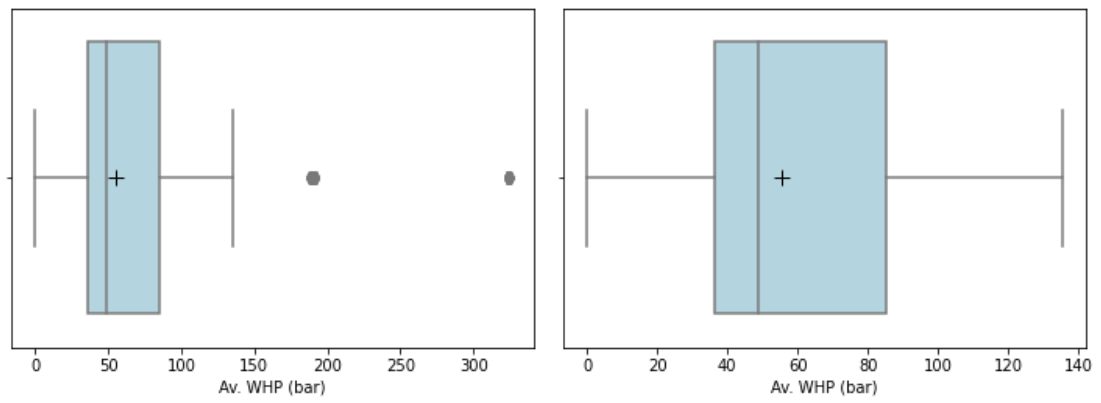


Figure A.16 : Kernel Density Estimation plot, histogram, and boxplot with and without outliers for Av. WHP (bar)

V.